

Linguistically-Informed Multilingual Instruction Tuning: Is There an Optimal Set of Languages to Tune?

Gürkan Soykan*
Information Technology Group,
Wageningen University and Research,
Wageningen, Netherlands

Gözde Gül Şahin**
Department of Computer Engineering,
Koç University, İstanbul, Türkiye
KUIS AI, Koç University, İstanbul,
Türkiye

Multilingual language models often perform unevenly across different languages due to limited generalization capabilities for some languages. This issue is significant because of the growing interest in making universal language models that work well for all languages. Instruction tuning with multilingual instruction-response pairs has been used to improve model performance across various languages. However, this approach is challenged by high computational costs, a lack of quality tuning data for all languages, and the “curse of multilinguality”—the performance drop per language after adding many languages. Recent studies have found that working with datasets with few languages and a smaller number of instances can be beneficial. Yet, there exists no systematic investigation into how choosing different languages affects multilingual instruction tuning. Our study proposes a method to select languages for instruction tuning in a linguistically informed way, aiming to boost model performance across languages and tasks. We use a simple algorithm to choose diverse languages and test their effectiveness on various benchmarks and open-ended questions. Our results show that this careful selection generally leads to better outcomes than choosing languages at random. We suggest a new and simple way of enhancing multilingual models by selecting diverse languages based on linguistic features that could help develop better multilingual systems and guide dataset creation efforts. All resources, including the code for language selection and multilingual instruction tuning, are made available in our official repository at <https://github.com/GGLAB-KU/ling-informed-mit>, enabling reproducibility and further research in this area.

1. Introduction

Having high-performance multilingual language models is essential for an inclusive and diverse NLP community. However, multilingual language models exhibit varying proficiency across languages (Joshi et al. 2020), with some languages benefiting less from these models due to disparities in model generalization (Xue et al. 2021; Shli-azhko et al. 2022). To address this, instruction tuning using multilingual instruction-response pairs has been proposed to enhance model performance across diverse languages (Muennighoff et al. 2022). However, this approach is challenged by high computational costs, the scarcity of high-quality tuning data for all languages, and the “curse

* E-mail: gurkan.soykan@wur.nl, Work is done during PhD studies at Koç University.

** Website: <https://gglab-ku.github.io/>

of multilinguality”—a phenomenon where adding more languages to the training mix diminishes per-language performance (Conneau et al. 2019).

Lifting the curse of multilinguality is a challenging task that has been approached through various methods. One line of approach is the use of architectural design variants, such as training language or language-family specific adapters (Chronopoulou, Stojanovski, and Fraser 2023; Pfeiffer et al. 2022). For instance, X-Mod (Pfeiffer et al. 2022) trains language-specific adapters for masked language models (e.g., mBERT (Devlin et al. 2019), XLM-R (Conneau et al. 2019)), and Hyper-X (Üstün et al. 2022) proposes a hypernetwork that uses task-language combinations to facilitate multi-task multilingual transfer. Despite the encouraging results, their application to multilingual instruction tuning (MIT) scenarios that mostly involve larger language models might not be straightforward; hence yet unknown.

From a data-centric perspective, recent works focus on instruction tuning methods to overcome both the scarcity of high-quality datasets and the expense of MIT. For instance, Shaham et al. (2024) perform instruction tuning using language groups of four with few examples, and Chen et al. (2023) evaluate the impact of language subsets within a constrained computational budget across various evaluation benchmarks within the scope of MIT. Both of these works provide evidence that multilingual instruction tuning, even with downsampled data, matches or exceeds the performance of monolingual tuning for each language.

Another research line investigates the incorporation of derived language features into a variety of multilingual NLP tasks. For instance, Üstün et al. (2020) perform adapter tuning using additional linguistic typology features for universal dependency parsing. Chen and Ritter (2020) use learned language vectors that are trained with typology prediction auxiliary task for model selection of better cross-lingual transfer. Finally, LangRank (Lin et al. 2019) employ multiple linguistic features such as geographic, genetic, syntactic etc. for better cross-lingual learning to select the best languages for downstream tasks e.g., NMT, POS Tagging, and entity linking. All aforementioned work shows improved performance on this wide range of multilingual tasks, hinting at the potential benefit of using language features for MIT.

Language features have been tested for various multilingual downstream tasks and have achieved competitive performance. However, the majority of the studies either do not use a systematic approach for integration, or require additional model training (Lin et al. 2019). Most importantly, there exists no prior work that investigates the impact of language features within the scope of MIT. The closest studies to us that perform multilingual instruction tuning (Shaham et al. 2024; Chen et al. 2023; Kew, Schottmann, and Sennrich 2023), perform language selection based on random subsets or consider only one-dimensional (categorical) aspects such as language family or script.

In this work, we hypothesize that optimum language subset selection for multilingual instruction tuning can be achieved in a linguistically informed manner. To test this hypothesis, we investigate various language selection techniques and how they compare to baselines such as randomly selected subsets, subsets curated by picking a single language from each language family, and all languages available in the dataset. Our linguistically informed selection techniques range from common methods, such as typological feature vectors (Littell et al. 2017), to less conventional but promising approaches, including language embeddings derived from semantic typology features (Chen, Biswas, and Bjerva 2023). We then apply a simple k-means clustering algorithm to select a fixed number of languages based on each technique. Finally, we instruction-tune models from three different model families using LoRA (Hu et al. 2021), focusing on specific language subsets, and measure their cross-lingual and cross-

task performance on natural language understanding and commonsense reasoning tasks under constrained computational resources. The experiments are conducted on the following model families: mGPT (Shliazhko et al. 2022), mT5 (Xue et al. 2021), and BLOOM (Scao et al. 2022), encompassing both relatively small and large decoder-only models, as well as an encoder-decoder model. We evaluate them on five common multilingual benchmarks. We systematically analyze the impact of each language subset and provide further insight on performance for unseen languages (§ 7.1), the effect of model size and varying the number of languages (§ 7.2, § 7.3), and compare the monolingual versus multilingual instruction tuning in a case study for Vietnamese (§ 7.4). Our experimental results indicate that:

- Linguistically informed language selection for multilingual instruction tuning generally outperforms the random baseline in terms of average performance across different model families and sizes, with statistically significant improvements observed in specific tasks and models.
- The top-performing language subset varies by specific tasks and model families, suggesting that the effectiveness of language selection strategies is not uniform and must be evaluated for the task and model combination.
- Cross-lingual generalization capabilities vary significantly across different language subsets, indicating that linguistically informed selections can aid in better generalization to unseen languages.
- Varying the number of languages in multilingual instruction tuning leads to task and model dependent performance effects. Increasing the number of languages beyond a certain threshold results in performance degradation, resembling the curse of multilinguality under a fixed computational budget.

Given our findings, feature-based language selection for multilingual instruction tuning could lead to improvements in cross-lingual performance. Additionally, it may inform researchers curating datasets on how to allocate their efforts effectively to encompass various aspects of linguistic diversity.

2. Related Work

Language Selection. Lin et al. (2019) formulate language selection as a ranking problem to facilitate optimal cross-lingual transfer by making use of language features such as geographic, syntactic, and phonological features for various downstream tasks. Glavaš and Vulić (2021) perform hierarchical clustering of languages based on syntactic similarity from URIEL (Littell et al. 2017) to select training treebanks for knowledge transfer on dependency parsing. Finally, several works (Yong et al. 2022; FitzGerald et al. 2023) among many others, consider various factors such as typological diversity (e.g., by leveraging the WALS (Dryer and Haspelmath 2013)), script diversity (Yong et al. 2022), and availability and internet influence FitzGerald et al. (2023) while curating datasets or modeling languages.

Language Clustering. Another line of research leverages the features of multilingual pretrained language models to form language representations. For instance, mPLM-Sim (Lin et al. 2023) extract embeddings using multilingual pretrained language models

on parallel corpora and calculate language similarity by comparing embeddings from different languages. This method allows for informed source language selections that may enhance cross-lingual transfer. It has been observed that mPLM-Sim correlates with certain linguistic measures such as lexical and syntactic similarity. Similarly, [Fan et al. \(2021\)](#) attempt to obtain language representations using multilingual PLMs by calculating the centroid of CLS embeddings for each language in a multilingual dataset and clustering based on these representations. Each cluster, referred to as a *Sprachbund* or *language federation*, shows that languages within the same cluster, when trained together, exhibit improvements in cross-lingual performance by mitigating cross-lingual contradictions. Although these embedding-based approaches have shown improvements in cross-lingual performance by mitigating cross-lingual contradictions, their effectiveness in cross-lingual generalization remains to be thoroughly evaluated. Their direct application to our setting is limited due to their reliance on high-quality embeddings extracted from the models. Furthermore, obtaining optimal embeddings from different model architectures, such as decoder-only models, is not straightforward and remains an active area of research. Furthermore, compared to our method, these embedding-based approaches can be computationally intensive and might depend heavily on the quality of the parallel corpora used. Apart from embedding-based methods, [Wang et al. \(2023\)](#) introduce a method for optimizing multilingual model training by grouping languages based on gradient similarities. The approach involves three steps: measuring gradient similarities across languages while training a multilingual model, using these similarities to form optimal language groups, and then training individual models for each group. A drawback of this approach is the need to train a separate model for each group.

On the other hand, The ACQDIV project ([Jancso, Moran, and Stoll 2020](#)) focuses on understanding the universal cognitive processes involved in child language acquisition. This project compiles a database of 15 child language acquisition corpora from 14 typologically diverse languages, applying a fuzzy clustering algorithm with typological feature values to ensure linguistic diversity. The methodology of the ACQDIV project aligns closely with our approach in creating maximally diverse language sets. However, while ACQDIV focuses on curating a database and developing a corpus aggregation pipeline, we use these diverse language sets specifically to optimize cross-lingual model performance. Our approach extends beyond data curation to practical application in multilingual language models, as our primary aim is to identify the optimal set of languages using linguistic features represented by universal language feature vectors. Clustering these features enables a model-agnostic language subset selection strategy.

Multilingual Instruction Tuning (MIT). MIT has emerged as a powerful technique to enhance the cross-lingual and cross-task capabilities of LLMs ([Muennighoff et al. 2022](#)). Recent research explores various strategies for leveraging language diversity in MIT, with a focus on efficiency and effectiveness within constrained computational budgets. [Chen et al. \(2023\)](#) investigate the use of the Alpaca dataset ([Taori et al. 2023](#))¹ and find that MIT with LoRA ([Hu et al. 2021](#)) outperforms monolingual tuning in all eight languages, even with downsampled data. [Shaham et al. \(2024\)](#) further demonstrate that including even a small amount of multilingual data (40 examples) during instruction tuning improves the LLM’s ability to follow instructions in those languages. [Kew, Schottmann, and Sennrich \(2023\)](#) explore the minimal amount of multilingual data needed for effective MIT for English-centric LLMs. They find that fine-tuning with just

three languages can enhance cross-lingual transfer for generative tasks, implying that including all potential target languages might not be necessary.

Previous works provide encouraging results on strategic language selection for MIT within a budget. However, our work differs by using multilingual language models as the base and we focus not on the amount of the multilingual data but the composition of it. Hence, our work builds upon these existing studies by systematically analyzing the impact of language selection on MIT performance within a constrained computational budget.

3. Methodology

We first survey the available language subset selection techniques that are *distinct* and *commonly used* in the literature for variety of multilingual downstream tasks (e.g., dependency parsing) and pre-training of multilingual language models (Glavaš and Vulić 2021; Chang et al. 2023). We find that researchers either use predefined categories; or richer, high dimensional feature vectors based on linguistic databases such as URIEL (Littell et al. 2017). Predefined categories are mostly chosen as the language family (Kew, Schottmann, and Sennrich 2023; Ógúnremí, Jurafsky, and Manning 2023), script (Fujinuma, Boyd-Graber, and Kann 2022; Yong et al. 2022) and availability of resources (i.e., high-resource, low-resource) (FitzGerald et al. 2023); whereas feature vectors mostly contain a richer set of linguistic properties on various dimensions such as syntactic, phonetic and geographic.

Typological Feature Vector (TYPO): Similar to previous works using typological feature vectors for universal dependency parsing (Üstün et al. 2020), we employ 289-dimensional binary features containing syntactic, phonological, and phonetic inventory features extracted from the URIEL database (Littell et al. 2017).²

Learned Language Vector (LEARN): We use the learned language vectors by Malaviya, Neubig, and Littell (2017), which utilize parallel text for 1017 languages to train a many-to-one neural machine translation system with a typology prediction auxiliary task. These vectors, termed *MTBoth*, comprise language embeddings combined with the mean of hidden cell states of the encoder LSTM. They are commonly utilized for various tasks including model selection to enhance cross-lingual transfer capabilities (Chen and Ritter 2020).

Geographical Feature Vectors (GEO): These vectors articulate the orthodromic distance of a language to specific points on the globe, scaled as a fraction of the Earth’s antipodal distance. Data points for languages are derived from Glottolog (Hammarström et al. 2023), WALS (Dryer and Haspelmath 2013), and SSWL’s (Collins and Kayne 2011) language location information. Researchers (Chang et al. 2023; Lin et al. 2019) have extensively used them for

1 A collection of English instruction-response pairs, translated into eight languages (Bulgarian, Czech, Chinese, German, Finnish, French, Russian, Spanish) using machine translation systems

2 These binary features incorporate data from databases like WALS (Dryer and Haspelmath 2013), PHOIBLE (Moran and McCloy 2019), among others. Missing annotations in syntax, phonology, and inventory features are predicted via k-nearest neighbors on genetic, geographical, and feature distance averages by previous work (Littell et al. 2017).

tasks such as language modeling, neural machine translation (NMT), part-of-speech (POS) tagging, and entity linking.

Semantic Typology (SEM): More recently, [Chen, Biswas, and Bjerva \(2023\)](#) introduce the idea of utilizing colexification—words with multiple meanings—across languages to learn language representations. To do so, they use a pipeline where they first build a synset graph based on WordNet synsets, and train synset embeddings using several node embeddings methods. Then, these node embeddings are combined to first form colexification embeddings and then language embeddings by various techniques (e.g., max-pooling, summing etc. . .)³. Despite its recency, SEM is already used for tasks like offensive language detection ([Zhou et al. 2023](#)).

Baselines: As a categorical baseline, we use the language family (FAM). Additionally, we select random subsets of languages (RND) and use all available languages (ALL) for comparison purposes.

We used the `lang2vec` library⁴ to acquire the vectors for majority of the aforementioned methods (TYPO, LEARN, GEO). For SEM embeddings, we used the provided repo⁵ by the [Chen, Biswas, and Bjerva \(2023\)](#).

3.1 Language Selection

Algorithm 1 Language Selection Algorithm

Input:

features
seed

▷ Language vectors
▷ Seed value

Output: *selected_langs*

▷ Selected languages

```
1: selected_langs ← an empty list
2: SETSEED(seed)
3: kmeans ← APPLYKMEANS(features)
4: centroids ← kmeans.cluster_centers
5: for each centroid in centroids do
6:   Find the closest feature vector to centroid
7:   Add the matching lang. to selected_langs
8: end for
9: return selected_langs
```

To select languages using language features, we use an embarrassingly simple k-means clustering algorithm given in Alg. 1. It operates as follows: specific features for each language are first extracted for each *subset*, and then subjected to clustered with k-means, where $k = 14$. The choice of 14 is deliberate—aligning with the number of language families and also avoiding the curse of multilinguality. This number ensures a balance between diversity and manageability in our analysis. Following the clustering,

³ We use the CLICS prone embeddings with the concat+max setting due to its generally superior performance on realistic scenarios, such as complex sentences.

⁴ <https://github.com/antonisa/lang2vec>

⁵ <https://github.com/siebeniris/Colex2Lang/>

the centroid of each cluster is identified. Next, we determine the closest feature to each centroid. Finally, the language associated with this nearest feature is added to the list of selected languages. To ensure robust results, we run this process over different random seeds. This method allows us to systematically and objectively select languages that represent the central tendencies of each cluster, ensuring a broad and representative sample.

3.1.1 Clustering Methodology. We employed the KMeans algorithm from the scikit-learn library (1.4.1)⁶ to categorize language feature vectors into 14 distinct groups. The KMeans algorithm was utilized with its default configuration settings as defined in the documentation, except for the number of clusters. The initialization method ‘k-means++’ was chosen. We used the ‘auto’ setting for the number of times the algorithm is run with different centroid seeds. The seed values are 66, 42, and 10. The algorithm was allowed to run up to 300 iterations for a single initialization. We adhered to the tolerance level of 0.0001. The ‘Lloyd’ algorithm was selected for our experiments. When KMeans was not applicable (e.g., for RANDOM and FAMILY subsets), different seeds were set before sampling to ensure the robustness of our method.

4. Experimental Setup

We use Bactrian-X (Li et al. 2023), since it is notably the most comprehensive multilingual instruction dataset to date, which is further enriched by translated English instructions and responses generated by ChatGPT⁷. Following the budget-controlled instruction tuning concept and downsampled multilingual setting (Chen et al. 2023), we merge and shuffle data from all languages. This standardizes the computation budget across experiments and facilitates fair comparisons by ensuring equal representation in the training set. With a typical epoch encompassing 67k instances, we use 4786 instances per language⁸. In contrast, for our monolingual case study, we train exclusively in Vietnamese, utilizing all 67k instances derived from this single language (see § 7.4).

Selected Languages. We restrict our language selection to the 52 available languages in Bactrian-X. In our main experiments, we run the language selection algorithm with 3 different random seeds specifically for the KMeans clustering, ensuring that the selection of languages is not biased by the initial seed and thus enhancing the generalizability and robustness of the results. For non-feature-based subsets like FAM and RND, we also repeat the selection process 3 times with different random seeds. RND is selected through random sampling from all available languages, and for the FAM subset, we choose one language per language family, using random sampling within each family. Each subset selected with these different seeds is then used individually for fine-tuning the models. It is important to note that these seeds are separate from those used during the fine-tuning process, which ensures that the variation in language selection does not interfere with the model training itself. Table 1 shows the language subsets selected by

⁶ <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>

⁷ It is a parallel collection with 3.4 million instruction–response pairs across 52 languages. In our experiments, Telugu was excluded because it was not selected by our algorithm with any of the different random seeds. As a result, we utilized the remaining 51 languages for our experiments, except in the ALL setup

⁸ For each 14 language, each is equally proportionately represented. Furthermore, we maintain data parallelism by training on the same data pairs across languages

Table 1

Composition of language subsets for instruction tuning main results with respect to different seeds outlining the specific languages included in each subset, categorized based on selection criteria: language family (FAM), typological features (TYPO), learned language vectors (LEARN), geographical feature vector (GEO), Language Embeddings from Semantic Typology (SEM), and a randomly selected set (RND).

Subset	Languages
FAM	az,en,fi,he,ja,ka,vi,ko,ml,mn,my,th,tl,xh ar,en,fi,id,ja,ka,ko,mn,sw,ta,th,tr,vi,zh af,az,et,he,ja,ka,ko,ml,mn,my,th,tl,vi,xh
GEO	af,bn,et,fr,gu,he,hr,id,ja,kk,mn,sw,ta,vi af,bn,et,fr,he,hr,id,ja,ka,kk,mn,ta,ur,vi af,bn,et,fr,he,hr,id,ja,kk,mn,sw,ta,ur,vi
LEARN	cs,hi,id,km,ko,lt,lv,my,nl,pt,sl,ta,tl,vi cs,fi,hr,km,ko,lt,lv,my,nl,pt,sl,ta,uk,vi cs,hi,id,ja,lt,lv,nl,pt,sl,sv,ta,tl,tr,vi
RND	cs,gu,hi,id,ko,lv,mk,ml,ps,pt,si,ta,vi,zh cs,en,et,fr,hi,hr,it,ja,ml,mn,mr,ne,pl,pt ar,az,de,en,es,et,he,hi,id,ml,ps,ru,sv,uk
SEM	gl,gu,ka,kk,ko,ml,pl,ps,si,sl,sv,tl,uk,vi en,fr,gl,gu,ka,kk,ko,pl,ps,sl,sv,tl,uk,vi en,gl,gu,ka,kk,ko,ml,pl,ps,sl,sv,tl,uk,vi
TYPO	az,bn,de,et,fa,hi,it,ja,mk,sw,ta,tl,vi,zh ar,az,bn,de,et,fa,it,ja,mk,sw,ta,th,ur,zh ar,az,bn,en,fi,he,hi,it,ja,mk,nl,sw,ta,th

different seeds using our algorithm for the main results (see Table 2 for characteristics of all used languages).

Models. We experiment with various model architectures, each with unique pretraining backgrounds and capabilities. The mGPT 1.3B model (Shliazhko et al. 2022), despite its relatively small size, stands out as a highly performant, decoder-only architecture. It is a multilingual adaptation of GPT-3, pre-trained on an extensive selection of 61 languages from 25 diverse language families using sources like Wikipedia and the C4 Corpus without introducing architectural changes. For encoder-decoder architectures, we use the mT5 3.7B XL variant (Xue et al. 2021). mT5 is a multilingual extension of the T5 model, pre-trained on a Common Crawl-based dataset spanning 101 languages. We also use various configurations of Bloom (Scao et al. 2022), specifically the 1.7B, 3B, and 7B1 models, to examine scaling laws and the behavior of a 7B level decoder-only model. BLOOM is trained on the ROOTS corpus with content in 59 languages—including both natural and programming languages—demonstrates robust performance across benchmarks, with enhanced results after multitask prompted finetuning.

Table 2
Characteristics of all languages included in our subsets for instruction tuning.

Language	ISO 639-1	ISO 639-3	Family	SVO Order	Script
Afrikaans	af	afr	Indo-European	SVO/SOV	Latin
Arabic	ar	ara	Afro-Asiatic	SVO/VSO	Arabic
Azerbaijani	az	aze	Turkic	SOV	Latin
Bengali	bn	ben	Indo-European	SOV	Bengali
Czech	cs	ces	Indo-European	SVO	Latin
German	de	deu	Indo-European	SOV	Latin
English	en	eng	Indo-European	SOV	Latin
Estonian	et	est	Uralic	SVO	Latin
Persian	fa	fas	Indo-European	SVO	Persian (Arabic)
Finnish	fi	fin	Uralic	SVO	Latin
French	fr	fra	Indo-European	SVO/VSO/VOS/OVS	Latin
Galician	gl	glg	Indo-European	SOV	Latin
Gujarati	gu	guj	Indo-European	SVO	Gujarati
Hebrew	he	heb	Afro-Asiatic	SVO	Hebrew
Hindi	hi	hin	Indo-European	SVO	Devanagari
Croatian	hr	hrv	Indo-European	SVO	Latin
Indonesian	id	ind	Austronesian	SOV	Latin
Italian	it	ita	Indo-European	SVO/SOV	Latin
Japanese	ja	jpn	Japonic	SOV	Kanji,Hiragana,Katakana
Georgian	ka	kat	Kartvelian	SVO	Georgian
Kazakh	kk	kaz	Turkic	SVO	Cyrillic
Khmer	km	khm	Austroasiatic	SOV	Khmer
Korean	ko	kor	Koreanic	SVO	Hangul
Lithuanian	lt	lit	Indo-European	SVO	Latin
Latvian	lv	lav	Indo-European	SVO	Latin
Macedonian	mk	mkd	Indo-European	SOV	Cyrillic
Malayalam	ml	mal	Dravidian	SOV	Malayalam
Marathi	mr	mar	Indo-European	SOV	Devanagari
Mongolian	mn	mon	Mongolic	SOV	Cyrillic
Nepali	ne	nep	Indo-European	SOV	Devanagari
Burmese	my	mya	Sino-Tibetan	SVO	Burmese
Dutch	nl	nld	Indo-European	SVO	Latin
Polish	pl	pol	Indo-European	SOV	Latin
Pashto	ps	pus	Indo-European	SVO/SOV	Arabic (Naskh)
Portuguese	pt	por	Indo-European	SVO	Latin
Romanian	ro	ron	Indo-European	SVO	Latin
Russian	ru	rus	Indo-European	SVO	Cyrillic
Sinhala	si	sin	Indo-European	SOV	Sinhala
Slovenian	sl	slv	Indo-European	SVO	Latin
Spanish	es	spa	Indo-European	SVO	Latin
Swedish	sv	swe	Indo-European	SOV	Latin
Swahili	sw	swa	Atlantic-Congo	SVO	Latin
Tamil	ta	tam	Dravidian	SVO	Tamil
Thai	th	tha	Tai-Kadai	SOV	Thai
Turkish	tr	tur	Turkic	SOV	Latin
Tagalog	tl	tgl	Austronesian	SOV	Latin
Ukrainian	uk	ukr	Indo-European	SOV	Cyrillic
Urdu	ur	urd	Indo-European	SVO	Arabic (Nastaliq)
Vietnamese	vi	vie	Austroasiatic	SVO	Latin (Quoc ngu)
Xhosa	xh	xho	Atlantic-Congo	SVO	Latin
Chinese	zh	zho	Sino-Tibetan	SVO/SOV	Simplified Chinese

4.1 Instruction Tuning Settings

During fine-tuning, only the responses were subject to loss computation, with instructions being masked. We fine-tuned with hyperparameters set to 7748 steps over 4 epochs, a learning rate of 3×10^{-4} , and a maximum sequence length of 768. For the mT5 encoder-decoder model, we set the maximum source length to 256 and maintained an effective batch size of 32.

We adopted LoRA (Hu et al. 2021) as Parameter-Efficient Fine-Tuning (PEFT) approach instead of full fine-tuning. LoRA retains the pre-trained model weights intact and introduces trainable rank decomposition matrices into the Transformer (Vaswani et al. 2017) layers and it is formalized by the following equation:

$$\mathbf{W}' = \mathbf{W} + \mathbf{B}\mathbf{A} \quad (1)$$

where \mathbf{W} represents the original weights, \mathbf{W}' the updated weights, and \mathbf{A} and \mathbf{B} are the low-rank matrices introduced by LoRA. In our implementation, we set the rank $r = 64$ and the scaling factor $\alpha = 16$, with an additional dropout rate of 0.05.

All experiments were conducted on a single GPU, with the NVIDIA RTX A6000 48 GB being the baseline. Training durations ranged from approximately 9.5 hours for the mGPT 1.3B model to around 39 hours for the largest BLOOM 7B model.

5. Evaluation

To rigorously assess our models and the effects of language subset selection on instruction tuning, we conduct a comprehensive set of experiments across a variety of multilingual NLP tasks, focusing on common natural language understanding and commonsense reasoning: XNLI, PAWS-X, XCOPA, XStoryCloze, and XWinograd explained below:

XCOPA (Ponti et al. 2020): Evaluates causal commonsense reasoning in 11 languages by predicting the correct next sentence from two options.

XStoryCloze (Lin et al. 2021): Assesses the ability to choose a plausible ending to a four-sentence story from two alternatives, across 11 languages.

XWinograd (Tikhonov and Ryabinin 2021; Muennighoff et al. 2022): A multilingual benchmark comprising Winograd Schema Challenge problems in six languages, testing commonsense reasoning.

XNLI (Conneau et al. 2018): Extends the MultiNLI corpus with translations into 14 languages, aiming to classify the entailment relation between sentence pairs.

PAWS-X (Yang et al. 2019): A multilingual paraphrase detection dataset based on English PAWS, requiring classification of sentence pairs as paraphrases or not, available in seven languages.

For consistency and reproducibility, we employ the lm-evaluation harness framework (Gao et al. 2023)⁹, primarily due to its widespread use and comprehensive coverage. However, we acknowledge its limitations in aligning with real-world language

⁹ <https://github.com/EleutherAI/lm-evaluation-harness>

model usage, as noted in recent critiques (Lyu, Wu, and Aji 2024; Biderman et al. 2024). We use the framework with default prompt templates in a zero-shot setting. We ensure that all selected datasets represent a diverse range of linguistic phenomena and reasoning tasks, providing a well-rounded evaluation of the models’ performance as much as possible. We report the performance of our models along with the 95% confidence intervals, which are given in Table 3. We calculated the margins of error, denoted by the ‘±’ symbol in the tables, based on three runs for each model and language subset combination. Additionally, we determined the p-values (highlighted) through paired t-tests. For the exact formulation of the confidence interval refer to Appendix B.

6. Results

Table 3

Our main results with BLOOM 7B1, mGPT, and mT5-xl models across different language selection strategies on multilingual benchmarks. The results with the highest mean for each model and task are shown in **bold**. We also compare language selection criteria to the random baseline using a paired t-test. We use **dark green** for $p \leq 0.05$ and **green** for $0.05 < p \leq 0.1$, and **light green** for $0.1 < p \leq 0.15$ to highlight the cases where the improvements over the baselines are statistically significant or worse. Similarly, we use **dark red**, **red**, **light red** to denote significantly lower scores. All experiments were conducted three times with different seeds. The “-” symbol denotes the base model, which has not been instruction-tuned.

Model	Subset	XNLI	XCOPA	XStoryCloze	XWinograd	PAWS-X	Avg.
Bloom	ALL	40.88 ± 0.25	58.18 ± 0.54	62.03 ± 0.13	73.71 ± 0.23	47.46 ± 0.82	56.45 ± 0.26
	RND	41.63 ± 0.25	57.72 ± 0.86	61.73 ± 0.30	73.91 ± 0.31	47.36 ± 2.85	56.47 ± 0.69
	FAM	41.18 ± 0.67	57.95 ± 0.19	61.49 ± 0.72	74.00 ± 0.81	47.88 ± 2.29	56.50 ± 0.56
	GEO	41.24 ± 0.66	57.97 ± 0.39	62.15 ± 0.31	73.96 ± 1.69	47.99 ± 0.85	56.66 ± 0.35
	LEARN	41.26 ± 0.66	58.08 ± 0.22	61.96 ± 0.06	74.17 ± 0.29	46.88 ± 1.85	56.47 ± 0.21
	SEM	41.19 ± 0.54	58.03 ± 0.66	61.44 ± 0.69	74.16 ± 0.18	47.35 ± 0.47	56.43 ± 0.31
	TYPO	41.63 ± 0.80	58.09 ± 0.54	61.94 ± 0.28	74.17 ± 1.08	47.75 ± 0.31	56.71 ± 0.24
	-	41.12	56.87	59.30	73.97	49.37	56.13
mGPT	ALL	41.23 ± 0.12	55.82 ± 0.14	55.98 ± 0.07	60.44 ± 0.30	50.14 ± 0.43	52.72 ± 0.07
	RND	40.71 ± 0.23	55.91 ± 0.62	55.97 ± 0.36	60.78 ± 0.20	48.71 ± 0.99	52.42 ± 0.36
	FAM	40.27 ± 0.08	55.82 ± 0.92	55.81 ± 0.85	60.88 ± 1.04	49.69 ± 1.72	52.49 ± 0.30
	GEO	41.01 ± 0.29	55.73 ± 0.56	56.35 ± 0.13	61.05 ± 0.44	49.97 ± 2.64	52.82 ± 0.55
	LEARN	40.64 ± 0.42	55.95 ± 0.47	56.15 ± 0.27	60.66 ± 1.03	48.92 ± 0.84	52.46 ± 0.38
	SEM	40.77 ± 0.21	55.61 ± 0.42	56.14 ± 0.34	60.57 ± 0.47	49.17 ± 2.17	52.45 ± 0.44
	TYPO	40.84 ± 0.31	55.93 ± 0.57	56.24 ± 0.44	60.39 ± 0.83	49.17 ± 2.06	52.51 ± 0.64
	-	40.90	55.04	54.43	60.55	50.30	52.24
mT5-xl	ALL	36.49 ± 0.69	53.64 ± 1.01	52.31 ± 0.21	50.39 ± 0.81	50.76 ± 1.15	48.72 ± 0.31
	RND	36.86 ± 0.29	53.10 ± 0.48	52.46 ± 0.52	49.88 ± 0.75	51.74 ± 1.74	48.81 ± 0.55
	FAM	36.12 ± 0.44	53.72 ± 0.34	52.59 ± 0.34	51.01 ± 1.69	52.28 ± 0.76	49.14 ± 0.19
	GEO	36.91 ± 0.36	53.23 ± 0.23	52.61 ± 1.03	50.20 ± 0.76	52.11 ± 0.91	49.01 ± 0.50
	LEARN	36.69 ± 0.55	53.48 ± 0.53	52.62 ± 0.21	49.65 ± 1.49	51.69 ± 1.05	48.82 ± 0.13
	SEM	36.41 ± 0.57	53.35 ± 0.29	52.38 ± 0.33	50.62 ± 1.89	51.63 ± 0.66	48.88 ± 0.39
	TYPO	36.62 ± 0.37	53.22 ± 0.54	52.58 ± 0.34	50.09 ± 0.86	51.40 ± 0.35	48.78 ± 0.11
	-	33.32	52.35	51.36	50.57	50.14	47.55

Our results, summarized in Table 3, demonstrate the impact of language subset selection on model performance across a range of NLP tasks. We observe the highest average scores with geographical feature vectors (GEO) and typological feature vectors

Table 4

Language Intersections for XNLI Task: This table shows the intersections of languages across models, subsets, and tasks for XNLI. The columns labeled M-S-T, S-T, M-T, and M-S represent the intersections between the model (M), subset (S), and task (T), with the denominator reflecting the total number of possible intersections. The languages that intersect across all three dimensions (M-S-T) are listed for reference. The table is sorted by columns M-S-T, S-T, M-T, and M-S, in descending order of intersections.

Task	Model	Subset	M-S-T	S-T	M-T	M-S	M-S-T Languages
XNLI	mT5-xl	ALL	13/15	13/15	15/15	50/52	ar,en,es,hi,ur,sw,tr,de,fr,ru,vi,zh,th
	mGPT	ALL	12/15	13/15	14/15	39/52	ar,en,es,hi,ur,sw,tr,de,fr,ru,vi,th
	mT5-xl	RND	9/15	9/15	15/15	31/32	ar,en,es,hi,fr,de,ru,vi,zh
	mT5-xl	TYPO	9/15	9/15	15/15	20/21	ar,en,hi,ur,sw,de,vi,zh,th
	Bloom	ALL	8/15	13/15	8/15	18/46	ar,en,es,hi,ur,sw,fr,vi
	mGPT	RND	8/15	9/15	14/15	24/32	ar,en,es,hi,fr,de,ru,vi
	mGPT	TYPO	8/15	9/15	14/15	18/21	ar,en,hi,ur,sw,de,vi,th
	mT5-xl	FAM	7/15	7/15	15/15	21/22	ar,en,sw,tr,vi,zh,th
	Bloom	RND	6/15	9/15	8/15	13/32	ar,en,es,hi,fr,vi
	Bloom	TYPO	6/15	9/15	8/15	8/21	ar,en,hi,ur,sw,vi
	mGPT	FAM	6/15	7/15	14/15	20/22	ar,en,sw,tr,vi,th
	Bloom	FAM	4/15	7/15	8/15	8/22	sw,ar,vi,en
	mT5-xl	GEO	4/15	4/15	15/15	15/16	sw,fr,vi,ur
	mGPT	GEO	4/15	4/15	14/15	14/16	sw,fr,vi,ur
	Bloom	GEO	4/15	4/15	8/15	8/16	sw,fr,vi,ur
	mT5-xl	LEARN	3/15	3/15	15/15	18/20	tr,hi,vi
	mT5-xl	SEM	3/15	3/15	15/15	15/16	fr,vi,en
	mGPT	LEARN	3/15	3/15	14/15	15/20	tr,hi,vi
	mGPT	SEM	3/15	3/15	14/15	11/16	fr,vi,en
	Bloom	SEM	3/15	3/15	8/15	5/16	fr,vi,en
Bloom	LEARN	2/15	3/15	8/15	5/20	hi,vi	

(TYPO) in the mGPT and BLOOM models respectively, highlighting the importance of these criteria in enhancing model performance. In contrast, the highest average performance for the mT5-xl model is found in the subsets created by the language family (FAM). Although the best performing selection slightly differs, we observe that linguistically-informed selection generally outperforms the random baseline on average. Additionally, there exists a linguistically-informed subset that is better than using all of the languages. From a task-specific perspective, several key patterns emerge in our current results.

For the **XNLI task**, we observe distinct patterns across the models. In BLOOM and mT5-xl, most language subsets yield significantly lower results compared to the random baseline, as indicated by the extensive red coloring in Table 3. Considering BLOOM, if we examine the intersections between the subset languages and the model’s pretraining languages in Table 4 (For the language intersections of XCOPA and XStoryCloze, refer to Table 5, and for XWinograd and PAWS-X, see Table 6.), we observe that BLOOM has the least overlap with its subsets and the highest number of languages not present in its pretraining corpus. Therefore, the results might be as seen. However, specific to XNLI, when we compare the pretraining corpora with XNLI language intersections, BLOOM has the least overlap with only 8 languages from XNLI’s total of 15 languages. In contrast, mT5-xl and mGPT cover almost all XNLI languages with 15 and 14 overlaps, respectively. Conversely, in mGPT, the trend is reversed; subsets such as GEO, SEM, and TYPO show significantly higher results, while only the FAM subset statistically underperforms. We believe this is because mGPT is exposed to 14 out of 15 languages

Table 5
Language Intersections for XCOPA and XStoryCloze Tasks

Task	Model	Subset	M-S-T	S-T	M-T	M-S	M-S-T Languages
XCOPA	mT5-xl	ALL	9/11	9/11	10/11	50/52	ta,et,it,tr,sw,id,vi,zh,th
	mGPT	ALL	8/11	9/11	8/11	39/52	ta,et,it,tr,sw,id,vi,th
	mT5-xl	FAM	8/11	8/11	10/11	21/22	ta,et,tr,sw,id,vi,zh,th
	mGPT	FAM	7/11	8/11	8/11	20/22	ta,et,tr,sw,id,vi,th
	mT5-xl	TYPO	7/11	7/11	10/11	20/21	ta,et,it,sw,vi,zh,th
	mGPT	TYPO	6/11	7/11	8/11	18/21	ta,et,it,sw,vi,th
	mT5-xl	RND	6/11	6/11	10/11	31/32	ta,et,it,id,vi,zh
	mGPT	RND	5/11	6/11	8/11	24/32	ta,et,it,id,vi
	mT5-xl	GEO	5/11	5/11	10/11	15/16	ta,et,sw,id,vi
	mGPT	GEO	5/11	5/11	8/11	14/16	ta,et,sw,id,vi
	Bloom	ALL	4/11	9/11	4/11	18/46	sw,ta,id,vi
	Bloom	FAM	4/11	8/11	4/11	8/22	sw,ta,id,vi
	Bloom	GEO	4/11	5/11	4/11	8/16	sw,ta,id,vi
	mT5-xl	LEARN	4/11	4/11	10/11	18/20	tr,ta,id,vi
	mGPT	LEARN	4/11	4/11	8/11	15/20	tr,ta,id,vi
	Bloom	TYPO	3/11	7/11	4/11	8/21	sw,ta,vi
	Bloom	RND	3/11	6/11	4/11	13/32	ta,id,vi
	Bloom	LEARN	3/11	4/11	4/11	5/20	ta,id,vi
	mT5-xl	SEM	1/11	1/11	10/11	15/16	vi
	mGPT	SEM	1/11	1/11	8/11	11/16	vi
Bloom	SEM	1/11	1/11	4/11	5/16	vi	
XStoryCloze	mT5-xl	ALL	10/11	10/11	11/11	50/52	ar,en,es,te,my,hi,sw,id,ru,zh
	mGPT	ALL	9/11	10/11	10/11	39/52	ar,en,es,te,my,hi,sw,id,ru
	Bloom	ALL	7/11	10/11	8/11	18/46	ar,en,es,te,hi,sw,id
	mT5-xl	RND	7/11	7/11	11/11	31/32	ar,en,es,hi,id,ru,zh
	mGPT	RND	6/11	7/11	10/11	24/32	ar,en,es,hi,id,ru
	mT5-xl	FAM	6/11	6/11	11/11	21/22	ar,en,my,sw,id,zh
	Bloom	RND	5/11	7/11	8/11	13/32	ar,en,es,hi,id
	mGPT	FAM	5/11	6/11	10/11	20/22	ar,en,my,sw,id
	mT5-xl	TYPO	5/11	5/11	11/11	20/21	ar,en,hi,sw,zh
	Bloom	FAM	4/11	6/11	8/11	8/22	sw,id,ar,en
	Bloom	TYPO	4/11	5/11	8/11	8/21	sw,ar,hi,en
	mGPT	TYPO	4/11	5/11	10/11	18/21	sw,ar,hi,en
	mT5-xl	LEARN	3/11	3/11	11/11	18/20	id,my,hi
	mGPT	LEARN	3/11	3/11	10/11	15/20	id,my,hi
	Bloom	LEARN	2/11	3/11	8/11	5/20	id,hi
	mT5-xl	GEO	2/11	2/11	11/11	15/16	sw,id
	mGPT	GEO	2/11	2/11	10/11	14/16	sw,id
	Bloom	GEO	2/11	2/11	8/11	8/16	sw,id
	mT5-xl	SEM	1/11	1/11	11/11	15/16	en
	mGPT	SEM	1/11	1/11	10/11	15/16	en
Bloom	SEM	1/11	1/11	8/11	5/16	en	

of XNLI, hence it is able to benefit more from the language subsets. However, following this logic, mT5-xl should also perform on par than mGPT, but this is not the case. This might be due to mGPT being pretrained on fewer total languages, with mGPT covering 61 languages while mT5-xl covers 101. It should be noted that all three models are pretrained on different datasets (Scao et al. 2022; Xue et al. 2021; Shliazhko et al. 2022). We observe that these differences are further reflected and emphasized by the ALL subset, which exhibits a divergent pattern: it results in significantly lower scores in BLOOM due to its low intersection with XNLI languages, significantly higher scores in mGPT due to its high intersection rate, and average performance in mT5-xl.

For the **XCOPA task**, mGPT shows only a slight drop with the SEM subset with p-value of 0.14 while no subset stands out in BLOOM, indicating on par performance

Table 6
Language Intersections for XWinograd and PAWS-X Tasks

Task	Model	Subset	M-S-T	S-T	M-T	M-S	M-S-T Languages
XWinograd	mT5-xl	ALL	6/6	6/6	6/6	50/52	en,pt,ja,fr,ru,zh
	mT5-xl	RND	6/6	6/6	6/6	31/32	en,pt,ja,fr,ru,zh
	mGPT	ALL	5/6	6/6	5/6	39/52	en,pt,ja,fr,ru
	mGPT	RND	5/6	6/6	5/6	24/32	en,pt,ja,fr,ru
	Bloom	ALL	3/6	6/6	3/6	18/46	fr,en,pt
	Bloom	RND	3/6	6/6	3/6	13/32	fr,en,pt
	mT5-xl	FAM	3/6	3/6	6/6	21/22	ja,zh,en
	mT5-xl	TYPO	3/6	3/6	6/6	20/21	ja,zh,en
	mGPT	FAM	2/6	3/6	5/6	20/22	ja,en
	mGPT	TYPO	2/6	3/6	5/6	18/21	ja,en
	mT5-xl	LEARN	2/6	2/6	6/6	18/20	ja,pt
	mT5-xl	GEO	2/6	2/6	6/6	15/16	fr,ja
	mT5-xl	SEM	2/6	2/6	6/6	15/16	fr,en
	mGPT	LEARN	2/6	2/6	5/6	15/20	ja,pt
	mGPT	GEO	2/6	2/6	5/6	14/16	fr,ja
	mGPT	SEM	2/6	2/6	5/6	11/16	fr,en
	Bloom	SEM	2/6	2/6	3/6	5/16	fr,en
	Bloom	FAM	1/6	3/6	3/6	8/22	en
	Bloom	TYPO	1/6	3/6	3/6	8/21	en
	Bloom	GEO	1/6	2/6	3/6	8/16	fr
Bloom	LEARN	1/6	2/6	3/6	5/20	pt	
PAWS-X	mT5-xl	ALL	7/7	7/7	7/7	50/52	en,es,ja,fr,de,ko,zh
	mT5-xl	RND	7/7	7/7	7/7	31/32	en,es,ja,fr,de,ko,zh
	mGPT	ALL	6/7	7/7	6/7	39/52	en,es,ja,fr,de,ko
	mGPT	RND	6/7	7/7	6/7	24/32	en,es,ja,fr,de,ko
	mT5-xl	FAM	4/7	4/7	7/7	21/22	ja,ko,zh,en
	mT5-xl	TYPO	4/7	4/7	7/7	20/21	de,ja,zh,en
	Bloom	ALL	3/7	7/7	3/7	18/46	es,fr,en
	Bloom	RND	3/7	7/7	3/7	13/32	es,fr,en
	mGPT	FAM	3/7	4/7	6/7	20/22	ja,ko,en
	mGPT	TYPO	3/7	4/7	6/7	18/21	de,ja,en
	mT5-xl	SEM	3/7	3/7	7/7	15/16	fr,ko,en
	mGPT	SEM	3/7	3/7	6/7	11/16	fr,ko,en
	Bloom	SEM	2/7	3/7	3/7	5/16	fr,en
	mT5-xl	LEARN	2/7	2/7	7/7	18/20	ja,ko
	mT5-xl	GEO	2/7	2/7	7/7	15/16	fr,ja
	mGPT	LEARN	2/7	2/7	6/7	15/20	ja,ko
	mGPT	GEO	2/7	2/7	6/7	14/16	fr,ja
	Bloom	FAM	1/7	4/7	3/7	8/22	en
	Bloom	TYPO	1/7	4/7	3/7	8/21	en
	Bloom	GEO	1/7	2/7	3/7	8/16	fr
Bloom	LEARN	0/7	2/7	3/7	5/20	-	

with random baseline. In contrast, mT5-xl performs statistically better compared to the random baseline with the FAM, LEARN, and ALL subsets, suggesting that it benefits more from diverse language subsets, possibly due to its alignment with the languages present in XCOPA. This can be attributed to the fact that mT5-xl covers 10 out of the 11 languages in XCOPA, which is not the case for mGPT and BLOOM. Therefore, we can infer that mT5-xl demonstrates cross-lingual generalization, even when the languages at the intersection of pretraining corpora and the task are not part of the instruction-tuning set. For instance, although the LEARN subset shares fewer languages (4 out of 11) compared to FAM and ALL, it still shows a performance improvement. This supports our hypothesis that multilingual instruction tuning can be optimized by linguistically

informed language subset selection, as LEARN outperforms random subset despite having fewer intersecting languages.

In the **XStoryCloze task**, the GEO and LEARN subsets consistently outperform others, particularly in BLOOM and mGPT, where they show significantly better results, with GEO achieving the highest average performance in both models. In mT5-xl, while GEO and LEARN subsets do not show statistically significant improvements, they still lead in average scores. This suggests that these subsets may provide robust cross-lingual transfer, possibly due to their broader coverage of linguistic diversity, which aligns well with the typological diversity claimed in XStoryCloze dataset (Lin et al. 2021). This is particularly noteworthy because the GEO and LEARN subsets have very few language intersections with the XStoryCloze task (only 2/11 and 3/11, respectively). For mGPT, the languages in the GEO subset are highly represented in the pretraining corpus (14/16), but the same cannot be said for BLOOM, where more than half of the GEO and LEARN languages are ones the model has not encountered during pretraining. Despite this, both subsets show performance improvements compared to the random baseline. This suggests that there may have been a positive interaction between the languages in XStoryCloze and those in the GEO and LEARN subsets, leading to cross-lingual generalization.

For the **XWinograd task**, we observe distinct patterns across the models. The SEM subset performs significantly better in BLOOM with p value of 0.14, but interestingly, it shows significantly worse performance in mGPT, along with the ALL subset. This disparity may be attributed not only to language intersections but also to the underlying datasets used during pretraining. BLOOM, pre-trained on the *ROOTS* dataset (Laurençon et al. 2022), benefits from a more diverse and filtered data source, likely enabling better cross-lingual generalization. Despite fewer task-model language intersections (3/6), BLOOM shows stronger performance, particularly in languages like *French* (fr), see Table 10, where high-quality pretraining data and having programming languages in the pretraining dataset could compensate for fewer overlapping languages for reasoning tasks such as XWinograd. On the other hand, the GEO subset yields significantly better results, achieving p-values of 0.13 and 0.005, in both mGPT and mT5-xl, indicating its potential for enhancing cross-lingual reasoning in this task. The difference between the SEM and GEO subsets in terms of language intersections could be attributed to the higher overlap between the pretraining corpus and the languages in the GEO subset. Notably, the FAM subset achieves the best performance in mT5-xl. This is particularly interesting because the random subset covers all six languages in XStoryCloze (6/6), whereas the FAM subset covers only three (3/6). However, mT5-xl still demonstrates strong cross-lingual generalization with the FAM subset. In fact, we observe this even in the model’s average cross-task performance.

For the **PAWS-X task**, we observe that the results are generally consistent across different subsets, with no significant differences. Interestingly, even the base, non-fine-tuned models achieve better average scores compared to the instruction-tuned models. Given that PAWS-X is a binary decision task, the results hover around 50%, with mGPT and BLOOM performing worse than expected. One possible explanation could be that the multilingual instruction tuning might negatively impact specific tasks like paraphrase identification.

Overall, our findings suggest that strategic selection of languages based on typology, geographical features, and language family features might provide benefits for crosslingual transfer, depending on the model and task, with statistically significant improvements observed in the mGPT-GEO and mT5-xl-FAM model-subset combinations for the average results. This aligns with our hypothesis that efficient language selection

can enhance model training *under resource constraints*. However, we also find that there is no universally “best selection” technique that consistently achieves the highest scores across tasks and models. Instead, subset selection strategies seem to be influenced by a range of factors, including the specific model architecture, the task at hand, the model’s pretraining objectives, the diversity and quality of the pretraining datasets, and even the inclusion of non-natural languages such as programming languages. These findings emphasize the necessity of tailored approaches for multilingual instruction tuning, where the choice of language subsets is adapted to the model, task, and pretraining conditions (see Tables 7, 8, 9, 10, and 11 for the complete breakdown of results per task and language). While these strategies show potential, it’s also crucial to acknowledge the limitations of the evaluation framework used. The lm-evaluation harness (Gao et al. 2023), despite its widespread use, has been critiqued for not fully aligning with real-world language model applications. This misalignment may partly explain the variability in our results, as it may not accurately capture the models’ capabilities across all linguistic phenomena and reasoning tasks (Lyu, Wu, and Aji 2024; Biderman et al. 2024).¹⁰

Table 7
Zero-shot accuracy performance across XNLI for all models.

Model	Subset	ar	hi	es	ru	sw	tr	bg	de	en	fr	th	ur	vi	zh	Avg		
BLOOM-1.7B	Fam	33.36	33.01	47.96	38.62	33.57	34.31	35.41	38.27	33.90	51.33	47.44	33.46	40.43	44.68	38.50	39.66	
	Geo	33.37	43.59	48.19	39.01	34.88	34.02	36.64	38.63	34.03	53.40	47.07	33.53	39.32	44.10	40.98	40.05	
	Learn	33.25	44.02	47.55	40.00	33.59	34.22	36.37	37.94	33.53	52.88	47.15	33.67	39.93	44.81	41.48	40.02	
	Rnd	33.32	42.77	47.56	39.49	34.70	34.30	36.52	38.62	33.32	53.20	46.79	33.63	39.68	45.41	40.32	39.97	
	Sem	33.45	43.98	47.38	39.08	33.82	34.21	36.35	37.54	33.60	52.60	46.48	33.63	41.37	45.45	38.89	39.85	
	Typo	33.36	43.45	47.87	38.69	34.95	34.54	36.95	38.69	33.45	53.01	47.15	33.90	39.45	45.27	40.39	40.10	
	All	33.41	43.29	47.58	40.16	34.30	33.76	36.71	37.78	33.84	53.82	47.15	35.63	39.99	45.59	39.68	40.16	
	-	33.49	44.02	48.19	39.48	35.18	35.22	35.78	39.44	33.82	51.49	47.11	33.49	38.67	46.18	35.98	39.83	
	-	33.62	44.70	49.93	40.60	36.17	33.83	36.56	39.56	34.61	54.29	50.60	35.34	39.04	43.76	41.74	40.96	
	-	33.62	45.26	49.22	40.99	36.39	33.17	37.71	40.33	33.45	54.63	50.00	33.44	39.14	43.47	42.69	40.86	
BLOOM-3B	Fam	33.57	45.25	49.79	41.42	34.67	33.81	38.10	41.28	34.07	54.50	49.85	33.72	39.01	44.35	40.71	40.94	
	Geo	33.67	48.03	49.08	41.62	35.22	33.93	37.70	40.83	33.55	54.42	49.61	33.45	39.16	44.67	41.45	40.88	
	Learn	33.38	45.61	48.88	42.48	35.03	34.38	37.70	40.09	33.71	54.33	49.27	33.78	39.45	43.64	39.41	40.74	
	Rnd	33.57	45.75	49.56	42.07	37.16	34.10	36.75	38.72	33.72	54.38	50.12	35.10	38.61	44.86	41.93	41.19	
	Sem	33.45	44.02	49.20	42.17	34.63	33.36	36.25	39.40	33.92	55.34	50.99	36.50	39.37	43.16	40.94	41.16	
	Typo	33.33	44.26	47.95	41.33	36.10	34.66	37.15	40.96	33.57	53.25	48.19	33.61	39.48	45.18	37.63	40.44	
	All	33.85	46.00	47.84	42.49	36.60	34.31	38.33	41.15	34.74	53.72	50.41	35.64	41.66	43.37	37.60	41.18	
	-	33.34	46.09	47.32	41.87	37.56	34.60	38.03	40.37	34.77	53.98	49.56	35.96	42.14	43.60	39.41	41.26	
	-	33.55	46.53	47.35	41.74	37.10	33.13	38.46	40.94	34.43	53.98	50.37	35.97	42.49	43.56	39.69	41.26	
	-	33.82	46.88	47.58	43.55	36.07	34.82	38.57	41.57	34.69	54.61	50.29	36.15	42.18	44.71	39.01	41.63	
BLOOM-7B	Fam	33.56	46.64	47.72	42.63	36.16	33.28	37.39	42.37	34.51	54.23	49.65	35.78	44.73	43.98	37.18	41.19	
	Geo	33.61	46.12	47.59	41.41	33.33	37.93	34.93	38.71	33.81	54.07	50.61	37.48	42.09	43.21	38.14	41.63	
	Learn	33.59	45.30	46.06	42.92	36.82	32.61	38.58	40.39	34.24	54.34	48.29	37.42	42.24	43.67	37.68	40.88	
	Rnd	34.06	46.71	48.47	42.49	38.27	35.22	38.67	39.24	36.99	53.53	48.80	33.65	40.84	44.62	35.22	41.11	
	Sem	33.67	45.38	47.81	44.41	36.48	42.57	40.45	41.26	35.57	54.14	49.14	44.43	35.79	43.92	34.61	40.27	
	Typo	33.82	46.00	45.34	44.47	41.00	41.31	40.83	40.10	46.10	54.34	50.37	45.03	38.06	43.52	33.16	41.01	
	All	33.68	46.09	44.30	44.65	40.47	41.42	41.94	41.15	40.35	53.37	49.74	44.87	36.55	44.31	38.37	40.64	
	-	33.80	46.52	44.50	44.77	41.40	40.34	42.62	40.48	44.69	53.65	49.20	44.85	36.80	43.75	33.66	40.71	
	-	34.19	46.68	45.11	44.83	41.61	42.45	42.01	40.48	45.02	55.70	48.81	44.26	36.98	44.21	33.50	40.77	
	-	33.99	46.05	40.70	44.56	41.94	40.60	42.82	39.22	41.41	54.82	49.41	44.82	36.83	44.45	33.45	40.84	
mGPT-1.3B	Fam	33.93	40.80	45.94	45.02	41.64	42.50	41.33	41.33	45.87	35.85	49.71	44.79	37.95	35.29	44.22	33.58	41.23
	Geo	33.53	41.00	44.62	43.37	41.89	41.81	41.73	45.38	35.74	49.52	44.94	36.83	34.74	44.66	33.78	41.02	
	Learn	33.51	36.72	36.15	38.39	35.11	34.06	39.17	35.88	33.37	42.01	36.56	35.16	34.82	33.57	33.38	36.12	
	Rnd	33.48	37.51	35.41	38.37	35.23	37.03	33.95	39.60	40.79	34.03	43.26	39.69	38.30	34.62	33.47	36.91	
	Sem	33.34	36.93	36.28	39.06	36.19	36.05	33.92	39.53	39.34	33.12	45.83	38.39	37.82	34.97	33.68	36.69	
	Typo	33.78	36.43	36.39	39.73	35.98	34.18	39.88	39.88	33.55	43.17	39.93	37.32	35.01	33.95	33.79	36.86	
	All	33.62	36.53	34.73	39.79	35.71	34.31	38.74	39.44	33.44	42.82	38.89	35.96	35.13	34.48	33.41	36.41	
	-	33.49	35.98	36.29	39.09	37.15	33.96	38.22	38.62	33.65	42.42	38.63	39.09	34.97	34.09	33.67	36.62	
	-	33.46	36.75	35.65	39.51	36.10	33.97	38.73	37.56	34.59	43.69	39.23	36.19	34.77	33.56	33.57	36.49	
	-	33.29	33.25	33.33	33.37	33.17	33.29	33.73	33.33	32.13	33.29	33.05	34.54	33.33	33.33	33.37	33.32	

10 While we wished to compare our findings with the reported results of the BLOOM, BXBLOOM (model trained with all the languages), and BMBLOOM (monolingual model) in the Bactrian-X study (Li et al. 2023), possible differences in the evaluation methodologies make that inappropriate. Even our baseline results for the BLOOM model showed significant differences. To ensure replicability and comparability, we opted for the lm-evaluation harness (Gao et al. 2023) instead of crafting custom prompts for the assessments.

Table 8
Zero-shot accuracy performance over XCOPA for all models.

Model	Subset	tr	sw	et	id	it	ta	th	vi	ht	qu	zh	Avg
BLOOM-1.7B	Fam	54.20 ± 1.09	52.67 ± 1.20	47.93 ± 2.35	65.73 ± 2.07	52.00 ± 1.13	55.67 ± 1.68	53.13 ± 1.74	68.27 ± 1.82	52.47 ± 4.66	<u>51.73</u> ± 1.03	64.33 ± 1.34	56.19 ± 0.81
	Geo	53.60 ± 1.24	53.53 ± 1.21	49.07 ± 0.76	68.00 ± 0.99	52.00 ± 1.31	53.93 ± 1.25	52.47 ± 1.88	67.40 ± 1.09	50.40 ± 1.31	51.27 ± 1.52	65.07 ± 0.76	56.07 ± 0.58
	Learn	54.00 ± 1.25	<u>53.87</u> ± 1.15	48.00 ± 1.31	66.13 ± 0.57	52.33 ± 0.57	54.80 ± 0.86	53.20 ± 2.17	67.33 ± 1.03	51.40 ± 1.49	50.73 ± 0.76	<u>65.13</u> ± 1.21	56.09 ± 0.56
	Rnd	54.67 ± 1.25	52.80 ± 0.76	<u>48.67</u> ± 0.29	65.87 ± 0.23	52.47 ± 2.04	55.60 ± 2.28	52.73 ± 1.15	<u>67.47</u> ± 1.52	50.13 ± 1.03	50.20 ± 3.10	65.07 ± 4.23	55.85 ± 0.32
	Sem	<u>54.60</u> ± 1.20	52.87 ± 1.15	47.33 ± 1.15	65.00 ± 0.00	52.93 ± 0.29	<u>56.33</u> ± 0.76	52.87 ± 1.52	67.13 ± 0.57	<u>52.33</u> ± 1.88	51.33 ± 1.88	64.13 ± 0.29	55.92 ± 0.30
	Typo	53.07 ± 1.19	53.87 ± 1.14	48.13 ± 1.52	65.93 ± 1.60	52.87 ± 1.43	54.87 ± 0.29	52.67 ± 2.91	67.33 ± 0.76	50.73 ± 0.76	50.27 ± 3.52	64.67 ± 1.60	55.85 ± 0.17
	All	53.93 ± 0.76	54.20 ± 2.28	47.33 ± 1.74	<u>67.07</u> ± 0.76	<u>52.93</u> ± 1.74	55.13 ± 0.76	52.73 ± 1.03	67.00 ± 1.31	51.13 ± 0.76	52.40 ± 0.50	65.27 ± 1.74	56.29 ± 0.28
-	52.80	51.80	47.40	63.20	52.60	56.60	53.20	65.80	50.40	50.60	61.40	55.07	
BLOOM-3B	Fam	53.60 ± 0.99	51.60 ± 3.88	49.13 ± 2.50	68.93 ± 1.03	52.87 ± 2.74	56.00 ± 1.31	<u>54.47</u> ± 1.52	71.33 ± 0.76	52.00 ± 3.30	51.87 ± 2.81	65.60 ± 0.86	57.04 ± 0.62
	Geo	55.67 ± 1.88	52.33 ± 1.74	<u>49.53</u> ± 2.07	70.93 ± 1.60	<u>53.73</u> ± 1.25	55.53 ± 1.25	54.00 ± 2.77	70.20 ± 1.49	52.87 ± 1.82	51.40 ± 0.50	<u>65.73</u> ± 0.29	<u>57.45</u> ± 0.88
	Learn	54.80 ± 1.79	<u>52.67</u> ± 1.03	48.80 ± 1.57	70.67 ± 1.74	53.20 ± 1.49	54.40 ± 1.79	52.47 ± 1.74	<u>71.00</u> ± 0.86	51.40 ± 1.72	51.20 ± 0.50	65.13 ± 0.76	56.88 ± 0.93
	Rnd	53.93 ± 2.83	52.53 ± 2.24	49.47 ± 0.76	69.33 ± 1.25	53.73 ± 3.31	55.73 ± 3.31	53.27 ± 2.35	70.33 ± 1.15	52.00 ± 2.28	<u>51.73</u> ± 1.52	65.13 ± 1.52	57.02 ± 0.48
	Sem	53.20 ± 1.49	52.80 ± 1.49	48.33 ± 1.74	69.60 ± 1.72	51.87 ± 0.76	<u>56.67</u> ± 2.51	53.47 ± 2.45	70.27 ± 0.29	51.93 ± 1.25	51.67 ± 1.03	64.87 ± 2.94	56.79 ± 0.31
	Typo	54.00 ± 2.63	52.40 ± 1.79	48.60 ± 2.28	70.00 ± 1.49	53.47 ± 1.10	55.67 ± 2.01	54.13 ± 1.25	70.60 ± 0.56	<u>52.33</u> ± 3.32	51.47 ± 2.91	65.67 ± 1.43	57.12 ± 0.59
	All	53.20 ± 1.49	52.20 ± 0.99	51.53 ± 1.88	<u>70.67</u> ± 2.76	53.80 ± 0.99	55.80 ± 0.99	54.47 ± 1.88	71.00 ± 0.56	51.87 ± 0.59	51.53 ± 1.43	66.07 ± 1.52	57.47 ± 0.39
-	53.40	51.40	49.20	69.20	51.60	58.20	52.60	68.80	50.20	50.60	62.00	56.11	
BLOOM-7B	Fam	51.67 ± 3.23	54.07 ± 0.29	<u>48.93</u> ± 1.03	<u>72.13</u> ± 1.88	53.73 ± 1.03	57.73 ± 2.91	55.80 ± 1.31	73.13 ± 1.74	52.40 ± 1.49	50.60 ± 1.49	67.27 ± 2.29	57.95 ± 0.19
	Geo	<u>53.33</u> ± 1.25	<u>54.40</u> ± 1.72	48.47 ± 0.76	72.20 ± 0.88	54.67 ± 1.74	57.13 ± 1.03	54.27 ± 2.07	73.20 ± 2.28	52.13 ± 0.29	50.47 ± 0.57	67.33 ± 1.25	57.97 ± 0.39
	Learn	52.47 ± 2.83	53.93 ± 1.52	48.00 ± 3.48	72.00 ± 4.24	54.93 ± 0.76	57.47 ± 1.03	<u>55.47</u> ± 1.88	72.80 ± 2.63	52.67 ± 1.03	51.00 ± 0.99	68.13 ± 1.82	58.08 ± 0.22
	Rnd	52.47 ± 2.74	53.40 ± 1.31	<u>49.20</u> ± 1.79	71.67 ± 3.49	54.13 ± 0.76	56.40 ± 2.98	54.27 ± 3.86	73.47 ± 1.88	52.40 ± 0.99	50.40 ± 1.49	67.20 ± 2.17	57.72 ± 0.86
	Sem	53.73 ± 2.87	54.67 ± 1.25	48.27 ± 1.25	72.00 ± 0.86	53.40 ± 2.48	<u>58.00</u> ± 1.49	53.80 ± 2.28	72.67 ± 1.03	52.93 ± 1.15	50.60 ± 1.79	67.47 ± 1.52	58.03 ± 0.66
	Typo	53.20 ± 4.74	52.87 ± 1.15	48.73 ± 2.45	72.07 ± 1.03	<u>54.80</u> ± 1.99	57.00 ± 1.49	54.80 ± 1.31	<u>73.33</u> ± 0.76	52.47 ± 1.25	51.67 ± 1.82	<u>68.00</u> ± 1.31	<u>58.09</u> ± 0.54
	All	53.40 ± 1.49	53.13 ± 0.57	51.07 ± 0.76	71.67 ± 1.60	54.00 ± 0.86	56.40 ± 1.31	55.07 ± 1.03	73.20 ± 0.50	<u>52.93</u> ± 1.60	<u>51.33</u> ± 1.60	67.80 ± 0.86	58.18 ± 0.84
-	51.20	51.40	48.20	69.80	52.80	59.20	55.40	70.80	50.80	50.80	65.20	56.87	
mGPT-1.3B	Fam	56.07 ± 1.25	56.27 ± 4.62	<u>52.40</u> ± 2.48	60.27 ± 1.15	<u>60.00</u> ± 1.31	54.20 ± 0.86	56.13 ± 0.57	61.73 ± 1.25	49.80 ± 3.48	51.93 ± 3.38	55.20 ± 2.28	55.82 ± 0.47
	Geo	56.07 ± 1.25	<u>58.07</u> ± 1.03	50.73 ± 1.52	60.87 ± 4.87	59.40 ± 0.99	54.27 ± 1.25	56.07 ± 2.50	62.47 ± 1.80	50.53 ± 1.03	50.00 ± 2.28	54.87 ± 1.43	55.73 ± 0.56
	Learn	55.87 ± 1.25	57.40 ± 2.28	<u>52.07</u> ± 1.15	59.87 ± 2.91	60.40 ± 0.86	54.00 ± 0.50	55.60 ± 0.50	62.27 ± 1.88	50.53 ± 0.76	51.60 ± 2.17	55.87 ± 2.27	55.95 ± 0.47
	Rnd	56.40 ± 0.99	57.00 ± 1.99	51.13 ± 1.60	<u>60.60</u> ± 0.86	59.27 ± 2.01	<u>54.73</u> ± 1.43	56.67 ± 1.03	61.33 ± 1.74	<u>51.00</u> ± 0.86	50.53 ± 1.25	56.33 ± 3.86	55.91 ± 0.62
	Sem	55.33 ± 1.60	57.00 ± 1.49	52.00 ± 1.31	59.93 ± 1.88	59.67 ± 2.74	54.20 ± 1.31	55.53 ± 0.29	61.20 ± 1.31	51.33 ± 1.60	50.80 ± 0.86	54.73 ± 1.52	55.61 ± 0.42
	Typo	55.93 ± 2.07	58.80 ± 0.86	51.40 ± 2.98	59.93 ± 1.60	59.13 ± 1.60	53.60 ± 0.86	57.33 ± 2.87	62.00 ± 0.86	51.00 ± 0.86	50.13 ± 1.25	<u>56.00</u> ± 2.28	<u>55.93</u> ± 0.57
	All	<u>56.13</u> ± 1.15	57.53 ± 1.25	51.80 ± 0.50	59.07 ± 0.76	58.47 ± 1.03	55.00 ± 0.86	<u>57.13</u> ± 1.15	<u>62.27</u> ± 1.60	50.07 ± 0.57	<u>51.80</u> ± 0.50	54.73 ± 1.60	55.82 ± 0.14
-	56.00	56.40	53.00	58.80	58.20	53.20	55.20	60.20	49.80	50.60	54.00	55.04	
mT5-xl	Fam	55.40 ± 2.28	54.40 ± 0.86	50.47 ± 3.66	<u>52.47</u> ± 0.76	54.40 ± 3.58	<u>55.07</u> ± 3.86	57.40 ± 2.63	54.20 ± 0.86	53.13 ± 1.43	51.00 ± 1.72	53.00 ± 2.77	<u>53.72</u> ± 0.34
	Geo	55.13 ± 1.60	55.07 ± 1.52	50.73 ± 0.76	51.60 ± 0.50	53.73 ± 1.74	54.67 ± 2.07	56.13 ± 0.29	51.67 ± 1.60	52.53 ± 5.74	51.07 ± 1.52	53.20 ± 1.79	53.23 ± 0.53
	Learn	55.13 ± 0.29	<u>54.93</u> ± 2.01	50.80 ± 1.79	50.93 ± 3.31	54.93 ± 1.60	55.20 ± 0.86	57.53 ± 1.25	51.93 ± 2.01	52.47 ± 2.35	51.27 ± 1.15	53.13 ± 1.03	53.48 ± 0.53
	Rnd	57.73 ± 1.43	53.73 ± 2.07	50.07 ± 2.83	51.93 ± 3.52	53.33 ± 0.57	54.73 ± 2.55	<u>57.73</u> ± 1.43	52.00 ± 2.28	52.53 ± 4.45	50.67 ± 2.01	52.60 ± 2.17	53.10 ± 0.48
	Sem	55.27 ± 3.16	53.33 ± 1.52	49.87 ± 2.74	51.87 ± 1.03	53.93 ± 1.60	54.07 ± 2.83	58.13 ± 0.29	<u>53.00</u> ± 2.58	<u>53.07</u> ± 1.52	<u>51.60</u> ± 0.50	52.73 ± 2.35	53.35 ± 0.29
	Typo	<u>55.93</u> ± 2.07	56.07 ± 1.03	52.60 ± 1.31	49.53 ± 2.45	53.07 ± 1.25	54.13 ± 2.24	56.67 ± 1.25	51.93 ± 0.57	52.40 ± 0.50	50.67 ± 4.70	<u>53.33</u> ± 1.03	53.22 ± 0.54
	All	55.60 ± 3.48	54.20 ± 1.79	<u>51.33</u> ± 4.62	53.40 ± 1.44	<u>54.33</u> ± 1.12	54.53 ± 1.55	57.07 ± 2.24	52.87 ± 0.76	52.60 ± 0.50	50.27 ± 1.74	53.87 ± 1.60	<u>53.64</u> ± 1.04
-	53.60	50.40	50.20	49.80	53.00	54.60	55.20	52.60	53.00	52.00	51.40	52.35	

7. Analysis and Discussion

Here, we aim to gain insights on the cross-lingual transfer capabilities of models by examining their generalization to unseen languages, akin to (Li et al. 2023). We also assess how model sizes influence the performance across language subsets, reflecting scaling law effects similar to (Yong et al. 2022). Following that, we inspected the effect of varying the number of languages using the same KMeans-based algorithm on a specific subset (Geo). Lastly, we provide a case study on monolingual versus multilingual instruction tuning using Vietnamese, inspired by (Chen et al. 2023).

To realize these analyses (except for the effect of varying the number of languages), specifically the analysis on unseen languages and the case study for Vietnamese, we kept the language subset fixed and repeated the experiments three times. This ensured that the languages remained constant for each subset, avoiding variations. For the FAM subset, which does not include language feature vectors, we selected one language per language family, excluding those languages present in our evaluation benchmarks (except for English) to maximize the observation of cross-lingual generalization. If a language family had two languages and one was part of the evaluation benchmarks, the other was selected. Consequently, this approach yielded 14 languages across the available language families. The final list of selected languages for each language subset in the analysis is given in Table 12.

Table 9
Zero-shot accuracy performance over XStoryCloze for all models.

Model	Subset	ar	es	hi	ru	sw	en	eu	id	my	th	zh	Avg
BLOOM-1.7B	Fam	55.68 ± 0.02	51.14 ± 0.86	57.25 ± 1.82	62.96 ± 2.67	55.68 ± 0.05	56.65 ± 1.02	58.15 ± 0.84	48.58 ± 2.90	67.81 ± 0.90	50.23 ± 1.62	61.84 ± 0.78	57.70 ± 0.21
	Geo	55.90 ± 0.19	52.08 ± 1.61	57.40 ± 1.38	64.64 ± 0.59	55.90 ± 0.19	56.89 ± 0.57	58.57 ± 0.29	47.41 ± 0.25	69.12 ± 0.62	50.23 ± 0.86	62.94 ± 0.60	58.24 ± 0.02
	Learn	54.62 ± 0.41	51.29 ± 1.19	57.27 ± 0.39	63.64 ± 0.05	54.62 ± 0.41	56.36 ± 0.62	59.03 ± 0.76	48.47 ± 3.92	68.65 ± 0.52	50.45 ± 0.84	62.48 ± 0.29	57.98 ± 0.44
	Rnd	55.86 ± 0.76	51.71 ± 0.25	57.38 ± 1.19	63.38 ± 0.47	55.86 ± 0.76	56.67 ± 0.41	58.35 ± 1.42	46.79 ± 0.75	68.06 ± 1.27	50.67 ± 1.33	62.30 ± 0.81	57.79 ± 0.31
	Sem	55.06 ± 0.91	51.69 ± 0.92	56.48 ± 1.01	62.87 ± 0.98	55.06 ± 0.91	57.11 ± 0.85	58.17 ± 0.50	46.66 ± 0.59	67.68 ± 1.98	50.78 ± 0.22	61.95 ± 1.18	57.57 ± 0.52
	Typo	55.68 ± 1.00	52.95 ± 0.34	56.92 ± 2.39	62.85 ± 1.24	55.68 ± 1.00	56.54 ± 0.77	58.90 ± 1.15	47.08 ± 0.84	68.58 ± 1.22	50.56 ± 0.82	62.43 ± 1.82	57.95 ± 0.71
	All	55.13 ± 0.17	52.66 ± 0.20	57.27 ± 0.35	63.38 ± 0.25	55.13 ± 0.17	56.91 ± 0.29	59.01 ± 0.62	48.51 ± 0.60	68.83 ± 0.72	50.26 ± 0.50	62.52 ± 0.39	58.25 ± 0.13
-	55.00	60.75	56.78	56.36	52.28	64.46	54.93	59.76	47.19	56.52	58.24	56.02	
BLOOM-3B	Fam	59.52 ± 1.25	70.02 ± 2.77	67.04 ± 0.72	56.12 ± 1.91	59.21 ± 0.84	64.26 ± 1.64	48.07 ± 3.32	50.98 ± 0.52	53.14 ± 2.97	57.53 ± 0.10	64.46 ± 1.68	59.13 ± 0.68
	Geo	60.25 ± 0.41	71.79 ± 0.73	68.48 ± 0.91	56.72 ± 0.59	60.18 ± 0.52	65.72 ± 0.29	46.30 ± 1.28	51.22 ± 0.87	54.56 ± 0.53	57.82 ± 0.25	65.08 ± 0.69	59.83 ± 0.22
	Learn	59.56 ± 0.57	71.72 ± 0.20	68.43 ± 0.32	56.61 ± 0.37	60.27 ± 2.41	64.97 ± 0.77	48.27 ± 3.12	52.13 ± 0.75	53.41 ± 1.14	57.80 ± 0.49	65.19 ± 1.19	59.85 ± 0.71
	Rnd	59.96 ± 0.92	70.57 ± 2.70	67.84 ± 0.83	56.96 ± 0.82	60.40 ± 0.47	65.12 ± 1.40	46.35 ± 0.91	52.04 ± 1.40	53.32 ± 0.52	58.08 ± 0.53	64.59 ± 2.49	59.57 ± 0.44
	Sem	59.61 ± 0.54	70.37 ± 2.24	67.70 ± 1.89	56.25 ± 0.16	59.94 ± 0.20	64.81 ± 0.81	47.21 ± 0.53	51.82 ± 0.72	52.93 ± 0.91	58.35 ± 1.35	64.13 ± 1.90	59.38 ± 0.73
	Typo	59.48 ± 1.55	71.12 ± 2.30	67.97 ± 0.72	56.61 ± 0.49	60.34 ± 1.22	64.81 ± 0.95	46.84 ± 2.14	51.14 ± 0.34	54.75 ± 0.10	57.38 ± 0.50	64.59 ± 1.25	59.55 ± 0.44
	All	59.54 ± 0.35	70.26 ± 0.41	67.83 ± 1.03	56.47 ± 0.73	60.27 ± 0.33	64.68 ± 0.41	48.27 ± 0.85	52.15 ± 0.39	54.60 ± 0.00	57.38 ± 0.12	64.90 ± 0.75	59.67 ± 0.10
-	56.65	64.00	57.51	50.69	53.01	66.84	55.66	61.02	46.66	58.17	60.95	57.38	
BLOOM-7B	Fam	62.19 ± 0.67	74.10 ± 2.96	69.60 ± 0.81	58.68 ± 1.00	62.30 ± 0.35	68.23 ± 2.71	49.17 ± 3.28	52.64 ± 1.96	54.87 ± 2.26	58.72 ± 1.43	65.94 ± 1.39	61.49 ± 0.72
	Geo	63.11 ± 0.19	76.07 ± 0.80	70.73 ± 0.68	59.10 ± 0.99	62.23 ± 0.78	69.01 ± 0.51	49.17 ± 0.16	53.23 ± 0.34	56.34 ± 2.34	58.31 ± 0.59	66.36 ± 0.83	62.15 ± 0.31
	Learn	62.10 ± 0.41	75.51 ± 0.45	71.15 ± 0.43	59.14 ± 0.67	62.34 ± 0.16	68.32 ± 0.34	49.77 ± 2.21	53.56 ± 0.62	55.28 ± 1.75	58.53 ± 1.01	65.87 ± 0.34	61.96 ± 0.06
	Rnd	62.56 ± 0.47	74.56 ± 1.76	70.48 ± 1.02	59.30 ± 1.78	62.48 ± 0.29	68.45 ± 1.36	48.89 ± 2.24	52.95 ± 2.32	54.75 ± 0.73	58.33 ± 0.69	66.29 ± 1.68	61.73 ± 0.30
	Sem	62.28 ± 1.00	74.12 ± 3.17	69.82 ± 2.11	58.44 ± 1.57	61.99 ± 0.77	67.70 ± 2.39	48.85 ± 1.99	53.70 ± 0.78	54.75 ± 1.12	58.84 ± 0.72	65.34 ± 0.69	61.44 ± 0.69
	Typo	62.63 ± 1.33	75.23 ± 2.73	70.35 ± 1.71	59.17 ± 0.66	62.32 ± 0.09	68.46 ± 1.33	48.95 ± 0.35	53.16 ± 0.10	56.04 ± 1.60	58.77 ± 0.50	66.25 ± 0.88	61.94 ± 0.28
	All	62.41 ± 0.88	75.25 ± 0.66	71.52 ± 0.82	58.99 ± 0.80	61.94 ± 0.76	68.92 ± 0.62	49.52 ± 0.41	53.43 ± 0.41	55.86 ± 0.34	57.99 ± 0.10	66.53 ± 0.10	62.03 ± 0.13
-	58.44	66.12	60.56	52.81	53.94	70.62	57.18	64.53	48.97	57.31	61.88	59.31	
mCPT-1.3B	Fam	51.05 ± 0.52	62.32 ± 0.09	57.51 ± 1.00	55.13 ± 1.28	53.72 ± 0.39	56.30 ± 0.34	51.23 ± 3.71	59.81 ± 0.94	54.29 ± 3.81	57.20 ± 0.51	55.33 ± 1.02	55.81 ± 0.85
	Geo	51.16 ± 0.44	62.61 ± 0.50	57.98 ± 0.59	55.37 ± 0.62	54.36 ± 0.69	57.20 ± 0.25	52.52 ± 0.53	59.87 ± 0.25	56.06 ± 1.18	57.36 ± 0.41	55.35 ± 0.86	56.35 ± 0.13
	Learn	50.63 ± 0.67	62.90 ± 1.16	57.71 ± 0.59	55.66 ± 0.72	54.14 ± 0.39	56.87 ± 0.41	51.03 ± 2.85	60.00 ± 0.32	56.15 ± 0.44	57.27 ± 0.39	55.26 ± 0.57	56.15 ± 0.27
	Rnd	50.83 ± 0.16	62.10 ± 2.24	57.44 ± 0.43	55.50 ± 0.35	53.89 ± 0.59	56.38 ± 1.08	52.26 ± 0.83	59.70 ± 0.34	55.99 ± 1.46	56.85 ± 0.56	54.73 ± 0.75	55.97 ± 0.36
	Sem	50.65 ± 0.33	62.72 ± 1.32	58.19 ± 0.96	55.35 ± 0.66	53.81 ± 0.16	56.30 ± 1.32	52.62 ± 0.57	59.19 ± 0.35	55.99 ± 0.17	57.66 ± 0.80	55.04 ± 0.35	56.14 ± 0.34
	Typo	50.91 ± 1.05	62.81 ± 1.11	58.13 ± 0.80	55.55 ± 1.39	53.96 ± 0.66	56.47 ± 0.59	53.03 ± 0.41	59.74 ± 0.09	56.30 ± 0.10	56.96 ± 0.90	54.84 ± 0.62	56.24 ± 0.41
	All	50.94 ± 0.34	62.06 ± 0.34	58.02 ± 0.10	55.53 ± 0.75	54.03 ± 0.19	56.87 ± 0.52	50.41 ± 0.25	60.53 ± 0.47	55.48 ± 0.52	56.63 ± 0.41	55.28 ± 0.10	55.98 ± 0.07
-	49.31	55.46	52.75	56.65	55.00	59.96	54.60	53.14	51.22	57.25	53.41	54.43	
mT5-xl	Fam	50.87 ± 1.10	54.45 ± 1.25	52.63 ± 1.12	52.33 ± 1.17	52.19 ± 0.90	53.12 ± 1.89	50.98 ± 1.66	53.1 ± 0.50	53.54 ± 0.17	54.77 ± 0.37	50.47 ± 1.64	52.59 ± 0.34
	Geo	50.52 ± 2.60	54.56 ± 2.47	53.47 ± 0.52	52.06 ± 0.35	52.35 ± 0.60	52.72 ± 2.17	51.44 ± 1.23	52.57 ± 0.82	53.61 ± 1.29	54.98 ± 0.10	50.39 ± 1.07	52.61 ± 1.60
	Learn	50.25 ± 0.91	54.95 ± 0.81	53.43 ± 2.01	52.41 ± 0.89	52.77 ± 0.81	53.08 ± 0.59	50.63 ± 2.01	53.34 ± 1.92	53.80 ± 0.76	54.67 ± 1.47	49.53 ± 2.42	52.62 ± 0.21
	Rnd	50.74 ± 0.41	54.09 ± 0.62	53.05 ± 0.66	51.84 ± 1.96	52.28 ± 1.56	52.20 ± 2.46	51.53 ± 0.47	53.50 ± 1.39	53.19 ± 0.20	54.91 ± 1.29	49.77 ± 0.16	52.46 ± 0.52
	Sem	50.43 ± 0.59	54.62 ± 0.52	53.01 ± 0.50	51.14 ± 1.23	52.44 ± 1.19	52.22 ± 0.98	51.45 ± 1.37	53.01 ± 1.15	53.21 ± 0.32	54.62 ± 0.41	50.06 ± 1.99	52.38 ± 0.53
	Typo	50.21 ± 0.82	54.40 ± 2.23	52.95 ± 0.16	52.37 ± 1.62	52.57 ± 1.00	52.33 ± 0.76	51.27 ± 1.33	52.81 ± 0.34	53.59 ± 0.81	55.09 ± 0.82	50.80 ± 0.78	52.58 ± 0.34
	All	50.1 ± 0.99	54.45 ± 1.40	52.99 ± 0.47	52.06 ± 0.49	52.02 ± 0.59	52.22 ± 0.50	51.05 ± 1.25	52.08 ± 0.43	53.83 ± 0.84	54.75 ± 0.34	49.86 ± 0.80	52.31 ± 0.21
-	49.50	52.08	50.96	49.77	52.22	52.68	51.16	51.36	50.83	53.67	50.76	51.36	

7.1 Performance on Unseen Languages

We examine the average performance of all models on languages *which were not included* in any of the language subsets, as shown in Figure 1. We identify the set of unseen languages for each task in our 44-language subset as follows: Arabic, Bulgarian, Greek, Spanish, Russian, and Turkish for XNLI; Haitian Creole, Quechua, and Turkish for XCOQA; Arabic, Spanish, Basque, Russian, and Telugu for XStoryCloze; Russian for XWinograd; and Spanish for PAWS-X.

Our results indicate the cross-lingual transfer abilities of the models. Our results indicate that our models are capable of cross-lingual transfer. There are two main comparisons to consider: base model (i.e., not tuned with instructions) and random subset. For the mCPT, mT5, and BLOOM-3B models, language subsets generally improve zero-shot performance compared to the base models. However, for the BLOOM family models, the impact varies by size. The relative improvement or decline in performance differs for different model sizes. Notably, for the BLOOM-7B, only two of the language family subsets show improvement over the base model. In comparison to the random subset, at least one linguistically-informed subset outperforms random selection in all models, with language family, Geographical Feature Vector and Semantic Typology subsets standing out in performance.

Regarding the results, we can speculate that the variations in the pretraining corpora of different models, as well as the amount and quality of training data available

Table 10
Zero-shot accuracy performance over XWinograd for all models.

Model	Subset	ja	ru	en	fr	pt	zh	Avg
BLOOM-1.7B	Fam	56.83 ± 2.99	<u>52.38</u> ± 1.36	74.11 ± 1.71	65.06 ± 2.99	67.81 ± 0.55	70.11 ± 1.03	67.85 ± 0.98
	Geo	55.41 ± 0.39	51.21 ± 1.82	74.88 ± 1.39	66.27 ± 8.98	<u>68.06</u> ± 4.33	68.19 ± 1.24	<u>67.69</u> ± 1.29
	Learn	55.61 ± 2.53	51.53 ± 1.20	74.14 ± 0.44	68.67 ± 7.92	69.07 ± 4.75	68.38 ± 2.48	<u>67.49</u> ± 1.15
	Rnd	54.81 ± 1.73	52.28 ± 3.88	<u>74.46</u> ± 0.33	65.06 ± 2.99	66.79 ± 4.84	68.59 ± 2.05	<u>67.36</u> ± 0.31
	Sem	55.27 ± 0.26	53.02 ± 5.17	74.18 ± 1.08	66.67 ± 1.72	67.17 ± 1.44	69.18 ± 1.24	<u>67.49</u> ± 0.99
	Typo	55.19 ± 2.65	51.54 ± 5.07	<u>74.84</u> ± 1.66	65.06 ± 5.99	66.54 ± 1.89	68.59 ± 1.24	<u>67.57</u> ± 1.16
	All	<u>56.06</u> ± 0.84	52.06 ± 3.94	74.05 ± 0.13	65.06 ± 0.00	66.92 ± 2.83	68.12 ± 1.50	<u>67.35</u> ± 0.47
	-	54.12	52.70	74.67	<u>68.67</u>	63.88	<u>69.44</u>	62.81
BLOOM-3B	Fam	<u>56.55</u> ± 6.12	52.59 ± 3.56	<u>79.44</u> ± 0.00	75.50 ± 1.72	<u>68.82</u> ± 2.50	72.55 ± 0.29	<u>71.12</u> ± 0.96
	Geo	56.20 ± 1.19	53.55 ± 1.64	<u>79.44</u> ± 0.88	<u>77.11</u> ± 10.79	67.81 ± 3.82	<u>73.02</u> ± 1.97	71.14 ± 0.79
	Learn	57.42 ± 2.77	52.27 ± 0.46	79.16 ± 0.80	75.50 ± 3.46	67.05 ± 3.57	72.29 ± 0.76	71.01 ± 0.99
	Rnd	55.37 ± 1.95	51.74 ± 6.26	<u>79.44</u> ± 1.48	74.70 ± 5.99	68.57 ± 1.44	71.63 ± 1.30	<u>70.68</u> ± 0.82
	Sem	56.31 ± 1.13	52.80 ± 1.20	79.10 ± 0.70	73.89 ± 7.54	68.44 ± 2.50	72.49 ± 1.25	70.85 ± 0.53
	Typo	55.16 ± 1.38	52.17 ± 2.53	79.91 ± 2.40	77.11 ± 2.99	67.43 ± 0.55	72.42 ± 1.96	70.98 ± 1.01
	All	54.71 ± 1.64	55.24 ± 1.58	78.99 ± 0.67	73.09 ± 1.72	66.79 ± 2.73	72.55 ± 1.25	<u>70.52</u> ± 0.28
	-	56.31	<u>53.65</u>	79.01	71.08	70.34	73.21	66.08
BLOOM-7B	Fam	<u>61.14</u> ± 0.53	55.03 ± 4.48	81.88 ± 0.43	71.08 ± 5.99	75.79 ± 2.90	73.54 ± 3.01	74.00 ± 0.81
	Geo	60.10 ± 1.07	53.54 ± 3.56	82.09 ± 1.50	74.30 ± 1.74	76.94 ± 3.57	74.01 ± 3.23	73.96 ± 1.69
	Learn	61.94 ± 1.04	55.24 ± 0.79	81.65 ± 0.78	<u>73.89</u> ± 4.57	<u>77.19</u> ± 1.89	73.28 ± 0.75	74.17 ± 0.29
	Rnd	60.06 ± 1.55	54.82 ± 2.54	82.34 ± 0.33	70.68 ± 3.46	75.03 ± 4.86	73.28 ± 0.29	73.91 ± 0.31
	Sem	60.55 ± 1.73	<u>56.93</u> ± 3.96	81.84 ± 0.33	72.29 ± 0.00	76.56 ± 4.66	74.47 ± 2.33	74.16 ± 0.18
	Typo	60.20 ± 1.17	56.72 ± 5.14	82.05 ± 1.21	71.08 ± 0.00	77.70 ± 1.97	73.94 ± 3.01	<u>74.17</u> ± 1.08
	All	59.86 ± 1.37	55.35 ± 3.19	81.86 ± 1.06	71.88 ± 3.46	76.43 ± 1.64	72.82 ± 2.26	73.71 ± 0.23
	-	58.92	57.14	<u>82.19</u>	71.08	76.43	<u>74.40</u>	69.15
mGPT 1.3B	Fam	55.65 ± 1.22	57.67 ± 1.64	62.78 ± 0.89	55.82 ± 4.57	59.70 ± 0.94	65.54 ± 2.49	<u>60.88</u> ± 1.04
	Geo	<u>55.51</u> ± 0.80	57.78 ± 0.79	<u>63.18</u> ± 0.75	57.03 ± 1.72	<u>59.07</u> ± 1.44	65.48 ± 1.30	61.05 ± 0.44
	Learn	54.82 ± 1.19	58.31 ± 2.54	63.08 ± 0.75	57.03 ± 6.23	56.40 ± 2.88	64.88 ± 0.50	60.66 ± 1.03
	Rnd	55.06 ± 0.44	56.40 ± 1.98	63.30 ± 0.22	57.43 ± 4.58	58.55 ± 2.51	64.48 ± 2.75	60.78 ± 0.20
	Sem	54.92 ± 1.76	57.25 ± 2.41	63.00 ± 0.51	57.03 ± 4.58	57.66 ± 0.55	64.28 ± 1.30	60.57 ± 0.47
	Typo	54.92 ± 1.42	56.93 ± 1.20	62.32 ± 0.84	<u>59.04</u> ± 5.19	57.67 ± 3.04	<u>65.74</u> ± 2.88	60.39 ± 0.83
	All	54.88 ± 0.84	57.04 ± 0.90	62.55 ± 0.44	58.23 ± 1.74	57.28 ± 1.44	65.41 ± 1.24	60.44 ± 0.30
	-	53.49	<u>58.10</u>	62.67	60.24	57.41	67.46	58.38
mT5-xl	Fam	50.37 ± 5.17	<u>52.80</u> ± 7.07	50.47 ± 0.72	53.41 ± 4.57	52.47 ± 3.27	52.45 ± 4.00	51.01 ± 1.69
	Geo	49.50 ± 0.60	50.27 ± 9.86	49.96 ± 0.83	51.40 ± 3.46	<u>51.71</u> ± 2.83	51.59 ± 4.69	50.20 ± 0.76
	Learn	49.50 ± 2.41	48.15 ± 7.16	50.01 ± 0.54	50.60 ± 2.99	49.05 ± 2.50	49.40 ± 4.30	49.65 ± 1.49
	Rnd	49.57 ± 2.32	50.48 ± 4.38	50.21 ± 1.09	46.58 ± 6.23	50.19 ± 2.83	49.01 ± 2.99	49.88 ± 0.75
	Sem	49.60 ± 3.30	51.96 ± 5.55	50.71 ± 2.73	50.20 ± 1.72	50.06 ± 4.66	<u>51.65</u> ± 1.13	50.62 ± 1.89
	Typo	49.84 ± 1.45	50.06 ± 1.82	50.26 ± 0.69	47.79 ± 12.46	49.68 ± 1.09	50.40 ± 4.70	50.09 ± 0.86
	All	49.84 ± 0.68	52.59 ± 4.34	50.37 ± 1.76	<u>52.61</u> ± 10.52	50.06 ± 2.88	49.93 ± 1.15	<u>50.39</u> ± 0.81
	-	<u>50.36</u>	53.02	<u>50.58</u>	48.19	49.43	50.40	50.32

for each language, could significantly affect model performance. However, this does not seem to apply uniformly across BLOOM models. Therefore, conducting similar experiments with models that have varying sizes, like the BLOOM family, might lead to different insights and interpretations.

Table 11
Zero-shot accuracy performance over PAWS-X for all models.

Model	Subset	es	ja	ko	de	en	fr	zh	Avg
BLOOM-1.7B	Fam	48.37 ± 3.65	55.42 ± 0.68	53.7 ± 0.43	49.45 ± 1.86	49.67 ± 5.02	54.63 ± 0.40	52.72 ± 1.33	51.99 ± 1.02
	Geo	46.52 ± 1.94	54.98 ± 2.34	54.15 ± 1.06	48.47 ± 1.79	46.97 ± 1.93	54.07 ± 0.56	50.92 ± 0.50	50.87 ± 0.29
	Learn	46.98 ± 1.15	55.78 ± 0.94	54.78 ± 0.07	50.52 ± 0.56	47.6 ± 2.04	54.32 ± 0.19	52.22 ± 1.99	51.74 ± 0.68
	Rnd	45.82 ± 2.84	54.88 ± 5.27	53.97 ± 4.16	49.25 ± 3.34	47.73 ± 4.49	54.27 ± 1.27	52.15 ± 1.30	51.15 ± 1.94
	Sem	49.07 ± 2.22	55.93 ± 0.07	54.85 ± 0.65	52.03 ± 0.75	49.23 ± 2.26	54.43 ± 0.50	52.88 ± 1.22	52.63 ± 0.66
	Typo	47.05 ± 1.10	54.33 ± 2.05	53.03 ± 2.04	50.27 ± 2.21	48.6 ± 3.81	54.48 ± 0.83	51.23 ± 1.06	51.29 ± 0.96
	All	46.2 ± 1.31	55.92 ± 0.19	54.45 ± 0.76	51.5 ± 0.75	47.57 ± 0.59	54.02 ± 0.40	52.83 ± 0.40	51.79 ± 0.43
	-	48.65	56.05	54.25	50.85	51.1	54.8	54.65	52.91
BLOOM-3B	Fam	44.78 ± 3.59	46.65 ± 1.22	52.73 ± 3.80	45.75 ± 2.90	43.72 ± 5.63	54 ± 1.18	51.35 ± 2.72	48.43 ± 2.49
	Geo	42.38 ± 2.31	47.37 ± 5.26	52.42 ± 4.23	45.62 ± 1.98	41.97 ± 1.69	52.7 ± 0.65	51.32 ± 1.13	47.68 ± 1.39
	Learn	43.17 ± 1.41	55.22 ± 1.56	52.35 ± 2.98	46.43 ± 2.17	41.13 ± 1.79	53.67 ± 1.17	50.48 ± 0.07	48.92 ± 1.15
	Rnd	41.5 ± 0.54	54.33 ± 4.09	52.22 ± 6.15	45.02 ± 1.15	40.07 ± 4.95	52.58 ± 1.19	49.42 ± 4.13	47.88 ± 0.23
	Sem	42.92 ± 3.24	54.83 ± 2.44	51.2 ± 3.73	45.72 ± 2.73	41.98 ± 0.75	53.73 ± 0.68	51.2 ± 1.55	48.8 ± 1.25
	Typo	42.12 ± 1.41	48.95 ± 3.67	52.27 ± 3.16	44.88 ± 1.00	39.57 ± 1.37	53.33 ± 1.62	49.53 ± 3.99	47.24 ± 0.42
	All	42.35 ± 1.20	54.62 ± 1.27	52.67 ± 0.29	46.23 ± 1.58	43.5 ± 0.81	52.9 ± 0.12	49.4 ± 1.22	48.81 ± 0.55
	-	42.4	55.45	53.95	45.1	43.8	54.25	52.95	49.7
BLOOM-7B	Fam	41.32 ± 3.23	53.73 ± 4.07	52.77 ± 2.52	47.13 ± 3.48	38.53 ± 4.29	52.02 ± 1.48	49.63 ± 3.80	47.88 ± 2.29
	Geo	40.4 ± 1.40	53.5 ± 3.68	54.03 ± 0.61	49.87 ± 5.39	37.38 ± 0.31	52.03 ± 2.35	48.73 ± 1.13	47.99 ± 0.85
	Learn	39.43 ± 0.19	50.3 ± 7.55	53.33 ± 2.42	47.93 ± 3.96	37.17 ± 0.50	51.35 ± 0.81	48.67 ± 0.59	46.88 ± 1.85
	Rnd	39.63 ± 1.99	53.07 ± 4.45	53.27 ± 0.96	46.5 ± 8.02	37.55 ± 5.12	52.7 ± 1.46	48.77 ± 0.63	47.36 ± 2.85
	Sem	41 ± 1.83	49.98 ± 2.51	53.38 ± 0.31	47.35 ± 1.88	37.28 ± 0.68	52.7 ± 1.79	49.72 ± 2.13	47.35 ± 0.47
	Typo	39.68 ± 0.83	54.55 ± 0.57	53.68 ± 1.21	49.18 ± 4.26	36.78 ± 1.12	52.22 ± 2.45	48.17 ± 0.80	47.75 ± 0.31
	All	39.65 ± 2.05	52.12 ± 1.07	53.7 ± 0.76	51.57 ± 1.69	35.85 ± 0.78	51.45 ± 0.54	47.9 ± 1.51	47.46 ± 0.82
	-	41.75	55.25	55.05	47.55	39.9	53.9	52.2	49.37
mGPT 1.3B	Fam	49.82 ± 3.74	54.87 ± 1.65	49.87 ± 2.81	46.77 ± 2.21	42.73 ± 4.77	54.88 ± 0.07	48.87 ± 2.67	49.69 ± 1.72
	Geo	50.10 ± 3.08	55.27 ± 1.04	49.90 ± 4.91	47.52 ± 4.29	44.30 ± 2.04	54.87 ± 0.14	47.80 ± 4.09	49.97 ± 2.64
	Learn	47.02 ± 1.67	52.48 ± 6.93	51.15 ± 4.95	47.05 ± 1.80	42.30 ± 1.24	54.87 ± 0.07	47.55 ± 3.36	48.92 ± 0.84
	Rnd	47.72 ± 1.86	54.23 ± 2.42	47.85 ± 5.70	45.57 ± 4.72	42.35 ± 0.66	54.83 ± 0.14	48.42 ± 1.55	48.71 ± 0.99
	Sem	49.62 ± 3.80	53.45 ± 2.44	50.75 ± 1.12	45.55 ± 4.02	43.68 ± 1.09	54.90 ± 0.00	46.22 ± 4.86	49.17 ± 2.17
	Typo	47.98 ± 2.29	54.82 ± 1.36	48.12 ± 4.30	45.75 ± 1.56	42.62 ± 1.52	54.83 ± 0.19	50.07 ± 5.83	49.17 ± 2.06
	All	49.97 ± 1.18	54.43 ± 0.80	47.97 ± 0.14	49.08 ± 0.50	43.53 ± 1.07	54.93 ± 0.07	51.07 ± 1.12	50.14 ± 0.43
	-	48.2	55.0	52.45	47.45	43.7	54.95	50.35	50.3
mT5-xl	Fam	45.72 ± 0.31	52.92 ± 2.62	53.50 ± 1.18	53.50 ± 0.77	53.47 ± 0.68	54.68 ± 0.40	52.13 ± 2.00	52.28 ± 0.76
	Geo	45.90 ± 0.93	52.85 ± 2.89	52.65 ± 0.81	53.23 ± 1.95	52.38 ± 0.50	54.78 ± 0.19	52.98 ± 2.26	52.11 ± 0.91
	Learn	46.13 ± 0.52	53.52 ± 1.36	53.10 ± 3.12	50.82 ± 1.97	52.35 ± 1.97	54.62 ± 0.26	51.27 ± 4.62	51.69 ± 1.05
	Rnd	45.85 ± 1.06	53.18 ± 1.33	52.50 ± 2.07	50.82 ± 3.59	53.23 ± 2.24	54.70 ± 0.50	51.92 ± 1.93	51.74 ± 1.74
	Sem	45.60 ± 0.94	51.17 ± 3.70	54.02 ± 1.55	51.82 ± 0.56	52.17 ± 1.18	54.52 ± 0.83	52.10 ± 2.57	51.63 ± 0.66
	Typo	46.00 ± 1.06	52.57 ± 1.48	51.05 ± 1.40	49.90 ± 2.27	53.65 ± 0.25	54.58 ± 0.19	52.05 ± 3.17	51.40 ± 0.35
	All	45.18 ± 0.14	50.43 ± 2.16	52.43 ± 1.87	48.65 ± 3.09	52.33 ± 1.01	54.73 ± 0.52	51.58 ± 3.06	50.76 ± 1.15
	-	45.7	47.9	55.2	55.25	45.3	46.2	55.4	50.14

7.2 Effect of Model Sizes

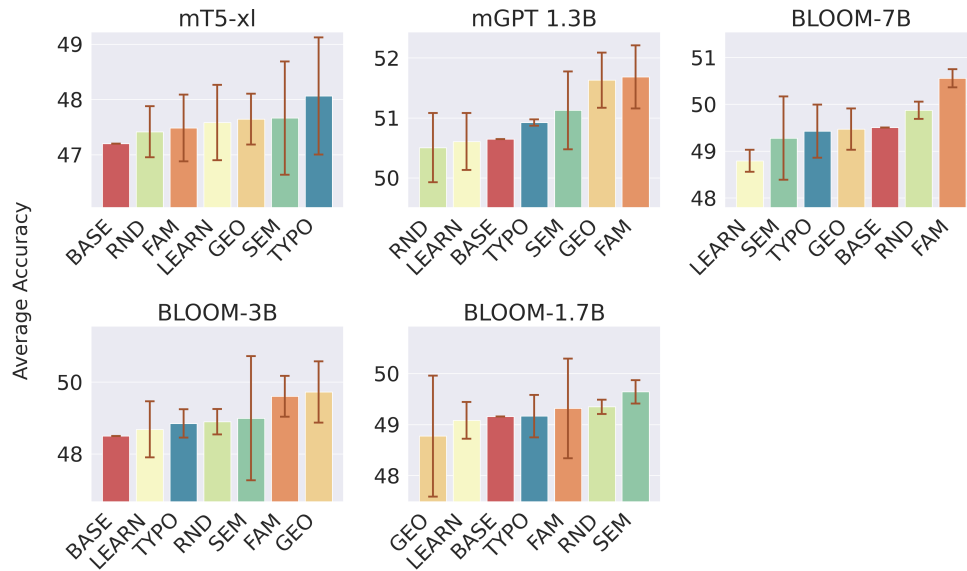
We extend our analysis to the impact of model size on average zero-shot performance for five tasks encompassing natural language understanding and commonsense reasoning. The comparative results are visualized in Figure 2.

Consistent with the scaling law, we observe that larger models generally lead to better performance. Specifically, all instruction-tuned models surpass the base model’s performance. After the base model, the lowest performance is given by the model trained with the randomly selected language subset. In the 3B model, the language family subset and the subset based on typological feature vectors perform similarly to the random selection.

Table 12

Composition of language subsets for instruction tuning analysis outlining the specific languages included in each subset, categorized based on selection criteria: language family (FAM), typological features (TYPO), learned language vectors (LEARN), geographical feature vector (GEO), Language Embeddings from Semantic Typology (SEM) and a randomly selected set (RND).

Subset	Languages
FAM	az,en,fi,he,ja,ka,ko,ml,mn,my,th,tl,vi,xh
GEO	af,bn,et,fr,he,hr,id,ja,ka,kk,mn,ta,ur,vi
LEARN	cs,fi,hr,km,ko,lt,lv,my,nl,pt,sl,ta,uk,vi
RND	cs,gu,hi,id,ko,lv,mk,ml,ps,pt,si,ta,vi,zh
SEM	bn,gl,gu,ka,ko,ml,my,pl,ps,sl,sv,tl,uk,vi
TYPO	az,bn,de,et,fa,hi,it,ja,mk,sw,ta,tl,vi,zh

**Figure 1**

Average performance of models instruction-tuned with various language subsets on natural language understanding and commonsense reasoning tasks, based on eight languages not present in any language subset (unseen).

The geographical language subset (GEO) consistently excels, (see Appendix A for visualization of languages by geographical feature vectors and a more detailed discussion), ranking first at 3B and 7B1 sizes. Conversely, the family-based selection (FAM) starts strong in the smallest model but diminishes in larger models. The typological vector (TYPO) showcases a unique progression, initially lagging in the 1.7B and 3B sizes but surging to second place in the 7B1 model. Lastly, subsets defined by informative

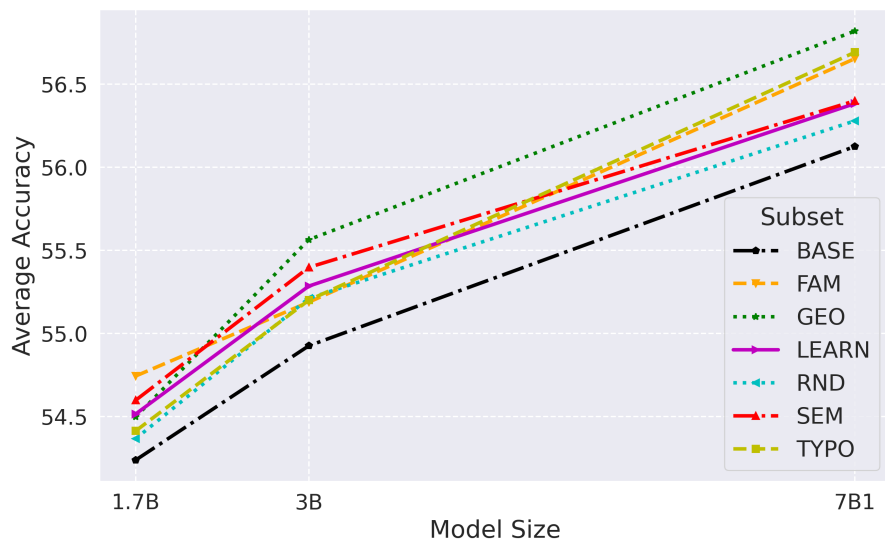


Figure 2
Comparison of average zero-shot performance across model sizes, demonstrating the scaling law and performance trends by language subset selection.

language features generally outperform random selection, particularly as the model size scales to 7B, where distinct performance clusters emerge, with base, random, learned and semantic typology language vector selections forming the lower part.

7.3 Effect of Varying Number of Languages

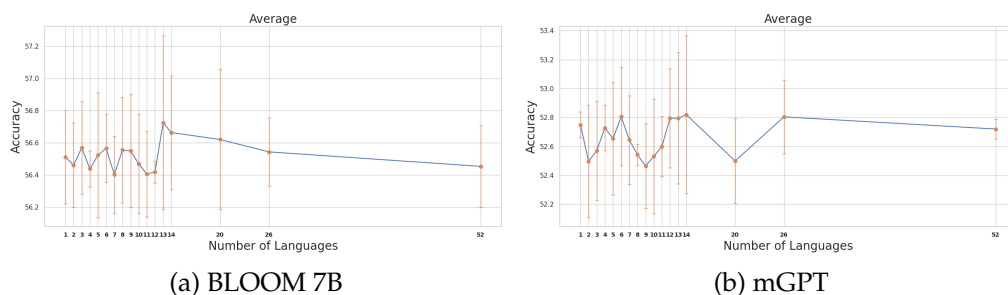


Figure 3
Average performance of BLOOM-7B and mGPT models with confidence intervals trained with varying numbers of languages, based on a geographical feature vector using our language selection algorithm, across natural language understanding and commonsense reasoning tasks.

Initially, we selected 14 languages to examine the effect of combining languages on maximizing cross-lingual and cross-task generalization. This was a deliberate choice to avoid the curse of multilinguality (Conneau et al. 2019), a phenomenon where the performance of a multilingual model degrades as the number of languages increases, due

to the model’s limited capacity to learn effectively from too many languages. However, an interesting research question remains: What are the effects of varying the number of clusters, and consequently, the number of languages, on multilingual instruction tuning using our clustering approach based on linguistically informed language features? To investigate this, we employ a k-means clustering algorithm to find centroids of language feature vectors and then select the closest language to each centroid. We vary the number of clusters from 1 to 14 (the number used in previous sections) and include additional settings of 20, 26, and 52 clusters, with 52 corresponding to the ALL setting. Regarding the selected language subset, our previous analysis and main results indicate that the GEO subset generally shows better and more consistent performance compared to other subsets using different linguistic features. Therefore, for this analysis, we focus on subsets selected using the geographical feature vector (GEO). Lastly, this analysis is conducted using the mGPT (Shliazhko et al. 2022) and BLOOM-7B (Scao et al. 2022) models to optimize resource usage.

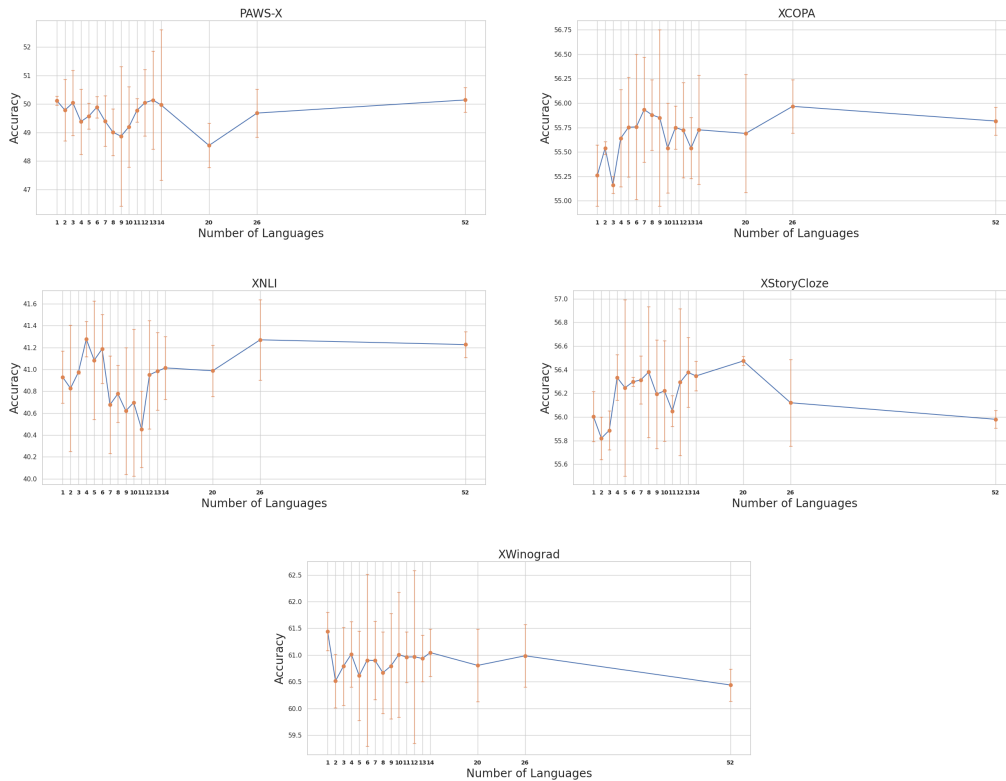


Figure 4
Effect of varying number of languages on different benchmarks for the mGPT model.

Figure 3 presents the averaged task performances with confidence intervals for the mGPT and BLOOM-7B models. Analyzing the results of the mGPT model in Figure 3.b, we observe that the results are not highly consistent. However, several insights can still be drawn. For instance, the model trained in all languages (52 languages) does not yield the best performance, in parallel with expectations. Despite the wide confidence intervals, when examining both the top of the confidence intervals and the mean results,

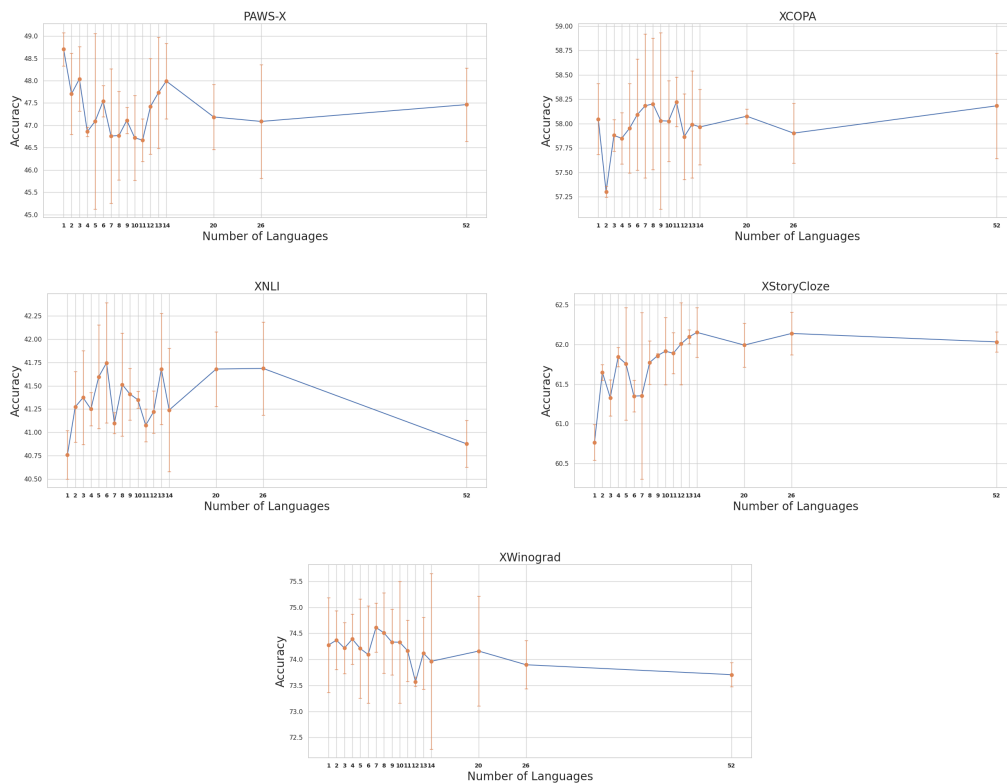


Figure 5
Effect of varying number of languages on different benchmarks for the BLOOM 7B model.

it can be stated that the default selection of 14 languages performs marginally better compared to other configurations.

Examining the results of the BLOOM-7B model in Figure 3.a, we observe more interpretable outcomes. The performance remains relatively stable up to 13 languages, with 13 and 14 languages yielding the best mean results. This further validates our choice of using 14 languages in our default setting. Beyond 14 languages, a consistent decline in performance is observed. Figures 4 and 5 present the average results per task. Considering the overall average and task-specific results, it can be considered that the impact of varying the number of languages used in multilingual instruction tuning is both task and model dependent. This observation opens an interesting research direction to investigate why certain tasks for specific models do not degrade in performance even with a high number of languages. It is important to note that the total computational budget for all experiments remains constant.

We can attribute these results to several factors: The inconsistency observed with fewer languages may be influenced by whether the selected languages are present in the evaluation tasks. Since we applied our language selection algorithm with random seeds, different results could emerge even for the same number of languages, potentially leading to relatively high confidence intervals. The performance decrease when moving from 26 to 52 languages cannot be explained by increased language coverage for tasks. We hypothesize that this decline, consistent with the findings of (Conneau et al. 2019), may be due to potential interference between more distantly related languages and

reduced exposure to each individual language as the number of languages in training increases.

7.4 A Case for Monolingual vs Multilingual Instruction Tuning: Vietnamese

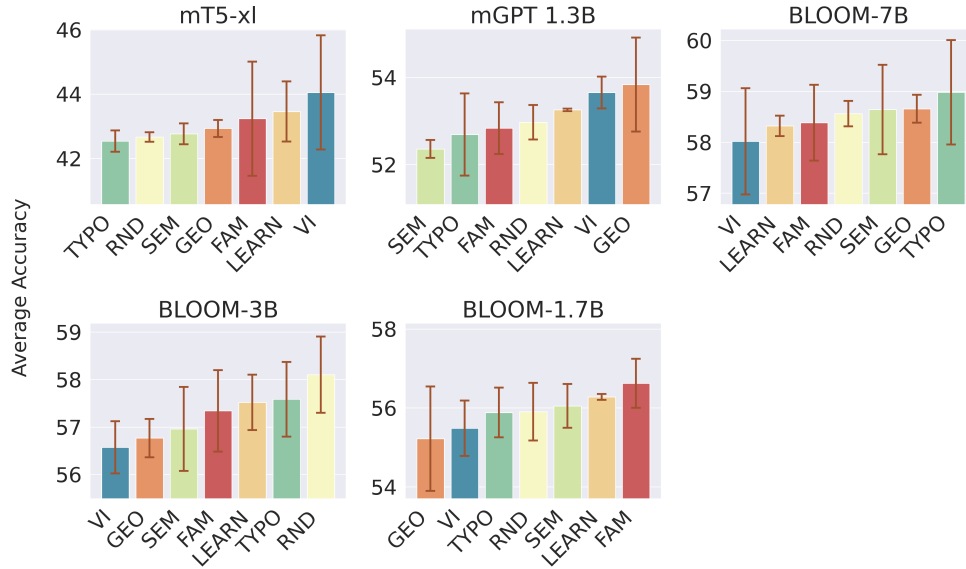


Figure 6

Average performance comparison of monolingual (VI) and multilingual instruction tuning with confidence intervals for Vietnamese across various models on XNLI and XCOPA tasks.

Inspired from [Chen et al. \(2023\)](#), we perform a case study on Vietnamese to analyze how monolingual instruction tuning compares with the proposed linguistically-informed multilingual instruction tuning techniques. We choose Vietnamese—identified as VI in our analyses—since it offers a unique case due to its inclusion based on language family, typological features, geographical positioning, and semantic typology, making the language a consistent participant across selected language subsets (see Analysis and Discussion § 7 for language selection). Exceptionally, we also deliberately place Vietnamese in the random subset to facilitate this targeted case.

As depicted in Figure 6, we compare the average accuracy of Vietnamese on NLU and commonsense reasoning tasks, XNLI ([Conneau et al. 2018](#)) and XCOPA ([Ponti et al. 2020](#)). It is worth noting that our experimental setup maintains computational parity; such that the number of iterations during training is consistent across all trials. In other words, for language subsets with N samples per language, the Vietnamese case used $N \times 14$ samples.

Interpreting the results, the monolingual Vietnamese case is the leading performer for mT5, and for mGPT, it ranks second in average score, with GEO performing slightly better. However, paired-t tests show no statistically significant difference between VI and GEO. In this case, there are multilingual training setups that perform statistically

worse than the monolingual case, which might be due to the increased training data size for the target language, Vietnamese. In contrast, within the BLOOM family, the monolingual Vietnamese case consistently exhibits a lower performance across all tested model sizes in terms of average scores. These results partially suggest that, depending on the model, a well-chosen mix of languages consisting of fewer samples per language might yield better outcomes than a monolingual approach, aligning with the findings of [Chen et al. \(2023\)](#).

8. Conclusion

In this work, we study the challenge of cross-lingual and cross-task generalization in multilingual instruction tuning by employing a linguistically guided language selection strategy and a simple clustering method. Our exhaustive experiments across various model types and parameter sizes, languages and evaluation benchmarks demonstrate that linguistic guidance to selected language subsets often perform better than random selection—especially as the model size grows. However, there is no one best selection technique that consistently outperforms the others across tasks and models—results are task and model dependent. We believe our approach can improve the development of better multilingual models through better instruction tuning. Furthermore, our results can guide researchers for data processing by helping to choose/curate the training data in a more informed way.

9. Limitations

The study’s methodology has limitations, particularly due to the exclusive use of LoRA ([Hu et al. 2021](#)) as a PEFT method and the reliance on the Bactrian-X ([Li et al. 2023](#)) dataset, which consists of 52 languages. While this dataset is extensive, it includes automatically translated content, which may introduce inaccuracies. The same issue is present in the evaluation benchmarks, as they also contain automatic translations. Additionally, these evaluations are inherently limited in representing real-world language model usage ([Lyu, Wu, and Aji 2024](#); [Biderman et al. 2024](#)). However, we chose this dataset and evaluation benchmarks because, despite these limitations, they remain among the most comprehensive and best available options for multilingual research. The outcomes might differ with a more diverse language dataset, as the language selection process could yield different results. The research was conducted under resource constraints, which limited the extent of training, the impact of prolonged training durations and the applicability of findings to models beyond the 7B scale might differ. Due to these constraints, we conducted our experiments three times (with three different seeds), which affects the statistical robustness of our results and leaves room for improvement in the calculation of confidence intervals and p-values for paired t-tests. Furthermore, our work relies on a relatively straightforward algorithmic approach, hinting at the possibility of achieving different or improved results with more complex methodologies targeting the nature of specific models and their exposition of pretraining corpus languages.

Ethics Statement

Our research, grounded in a commitment to ethical principles such as inclusivity, fairness, and the responsible use of data, acknowledges the profound impact that multilingual language models can have on global communication. Striving to promote linguistic

diversity and equity, we also recognize the environmental implications of computational research and have endeavored to balance the need for thorough experimentation with the imperative for sustainability. Finally, we are firmly against using our model weights in any harmful way and are dedicated to promoting responsible usage.

Appendix A: Visualization of Geographical Features

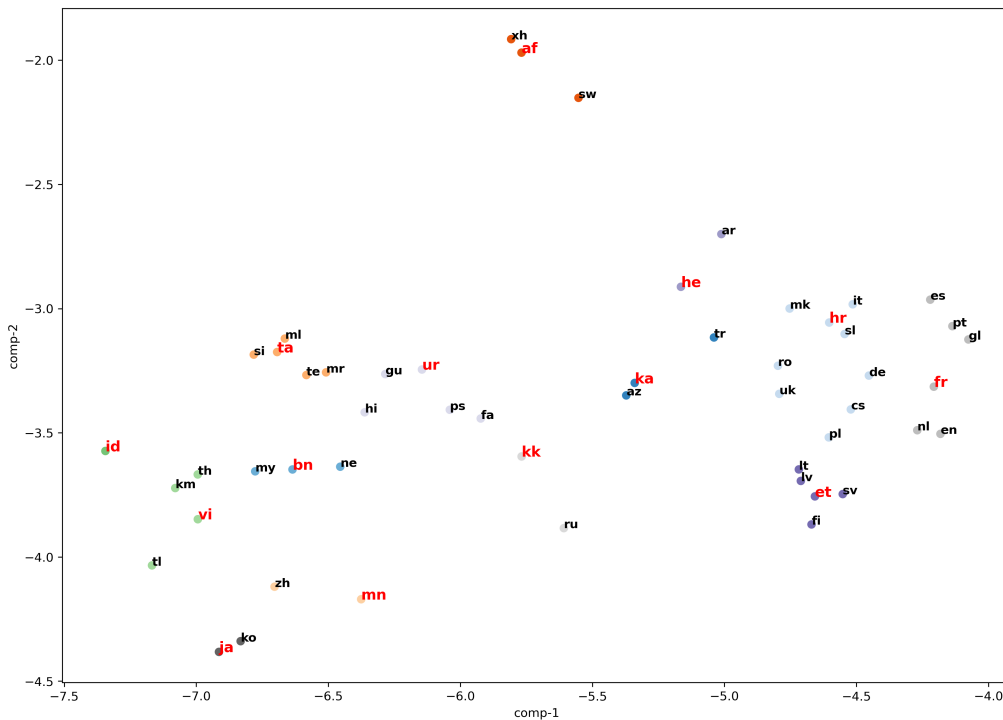


Figure A.1
t-SNE visualization of 512-dimensional geographical feature vectors for 52 languages, highlighting the selected languages in red text.

Geographical feature vectors, originally in 512 dimensions for the 52 languages featured in the Bactrian-X dataset (Li et al. 2023), are visualized in a two-dimensional space using t-SNE in Figure A.1. These geographical features have consistently performed well across various benchmarks and models, including NLU and commonsense reasoning tasks, performance on unseen languages, and the Vietnamese case study.

We attribute the success of the geographically-informed language subset to our algorithm’s capacity to select a diverse array of languages, as evidenced by the red labels in the figure. This suggests that, beyond typological or semantic features, there may be latent alternative representations in geographical features that are significant for language modeling.

Appendix B: Statistical Analysis Methods

We calculated 95% confidence intervals to estimate the precision of our results and used paired t-tests to assess the statistical significance of differences in performance between language selection strategies and baselines, such as random selection.

Appendix C: Confidence Interval

The margin of error for the confidence intervals, denoted by the ‘±’ symbol in the tables, was computed using the following formula for the t -distribution:

$$\text{Margin of Error} = t_{\alpha/2, \text{df}} \times \frac{s}{\sqrt{n}}$$

where:

- $t_{\alpha/2, \text{df}}$ is the critical value from the t -distribution for a 95% confidence level, where $\alpha = 0.05$ and df (degrees of freedom) is $n - 1$.
- s is the standard deviation of the sample.
- n is the sample size.

Acknowledgments

This work has been supported by the Scientific and Technological Research Council of Türkiye (TÜBİTAK) as part of the project “Automatic Learning of Procedural Language from Natural Language Instructions for Intelligent Assistance” with the number 121C132. We also gratefully acknowledge KUIS AI Lab for providing computational support.

References

- Biderman, Stella, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber Abbasi, Alham Fikri Aji, Pawan Sasanka Ammanamanchi, Sid Black, Jordan Clive, Anthony DiPofi, Julen Etxaniz, Benjamin Fattori, Jessica Zosa Forde, Charles Foster, Mimansa Jaiswal, Wilson Y. Lee, Haonan Li, Charles Lovering, Niklas Muennighoff, Ellie Pavlick, Jason Phang, Aviya Skowron, Samson Tan, Xiangru Tang, Kevin A. Wang, Genta Indra Winata, Francois Yvon, and Andy Zou. 2024. Lessons from the trenches on reproducible evaluation of language models. *ArXiv*, abs/2405.14782.
- Chang, Tyler A., Catherine Arnett, Zhuowen Tu, and Benjamin K. Bergen. 2023. When is multilinguality a curse? Language modeling for 250 high- and low-resource languages. *arXiv preprint*.
- Chen, Pinzhen, Shaoxiong Ji, Nikolay Bogoychev, Barry Haddow, and Kenneth Heafield. 2023. Monolingual or multilingual instruction tuning: Which makes a better alpaca. *arXiv preprint arXiv:2309.08958*.
- Chen, Yang and Alan Ritter. 2020. Model selection for cross-lingual transfer. In *Conference on Empirical Methods in Natural Language Processing*.
- Chen, Yiyi, Russa Biswas, and Johannes Bjerva. 2023. Colex2Lang: Language embeddings from semantic typology. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 673–684, University of Tartu Library, Tórshavn, Faroe Islands.
- Chronopoulou, Alexandra, Dario Stojanovski, and Alexander Fraser. 2023. Language-family adapters for low-resource multilingual neural machine translation. In *Proceedings of the The Sixth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2023)*, pages 59–72, Association for Computational Linguistics, Dubrovnik, Croatia.
- Collins, Chris and Richard Kayne. 2011. *Syntactic Structures of the World’s Languages*. New York University, New York.

- Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. In *Annual Meeting of the Association for Computational Linguistics*.
- Conneau, Alexis, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485, Association for Computational Linguistics, Brussels, Belgium.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Association for Computational Linguistics, Minneapolis, Minnesota.
- Dryer, Matthew S. and Martin Haspelmath, editors. 2013. *WALS Online (v2020.3)*. Zenodo.
- Fan, Yimin, Yaobo Liang, Alexandre Muzio, Hany Hassan, Houqiang Li, Ming Zhou, and Nan Duan. 2021. Discovering representation sprachbund for multilingual pre-training. In *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, pages 881–894, Association for Computational Linguistics.
- FitzGerald, Jack, Christopher Hench, Charith Peris, Scott Mackie, Kay Rottmann, Ana Sanchez, Aaron Nash, Liam Urbach, Vishesh Kakarala, Richa Singh, Swetha Ranganath, Laurie Crist, Misha Britan, Wouter Leeuwis, Gokhan Tur, and Prem Natarajan. 2023. MASSIVE: A 1M-example multilingual natural language understanding dataset with 51 typologically-diverse languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4277–4302, Association for Computational Linguistics, Toronto, Canada.
- Fujinuma, Yoshinari, Jordan L. Boyd-Graber, and Katharina Kann. 2022. Match the script, adapt if multilingual: Analyzing the effect of multilingual pretraining on cross-lingual transferability. In *Annual Meeting of the Association for Computational Linguistics*.
- Gao, Leo, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. A framework for few-shot language model evaluation.
- Glavaš, Goran and Ivan Vulić. 2021. Climbing the tower of treebanks: Improving low-resource dependency parsing via hierarchical source selection. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4878–4888, Association for Computational Linguistics, Online.
- Hammarström, Harald, Robert Forkel, Martin Haspelmath, and Sebastian Bank. 2023. Glottolog database 4.8. Available online at <https://glottolog.org>.
- Hu, Edward J, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Jancso, Anna, Steven Moran, and Sabine Stoll. 2020. The ACQDIV corpus database and aggregation pipeline. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 156–165, European Language Resources Association, Marseille, France.
- Joshi, Pratik, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Association for Computational Linguistics, Online.
- Kew, Tannon, Florian Schottmann, and Rico Sennrich. 2023. Turning english-centric llms into polyglots: How much multilinguality is needed?
- Laurençon, Hugo, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, Jörg Froberg, Mario Šaško,

- Quentin Lhoest, Angelina Mcmillan-Major, Gérard Dupont, Stella Biderman, Anna Rogers, Loubna Ben Allal, Francesco de Toni, Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Romero Muñoz, Jian Zhu, Daniel van Strien, Zaid Alyafeai, Khalid Almubarak, Vu Minh Chien, Itziar Gonzalez-Dios, Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Ifeoluwa Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilić, Margaret Mitchell, Sasha Luccioni, and Yacine Jernite. 2022. The BigScience ROOTS Corpus: A 1.6TB Composite Multilingual Dataset. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, New Orleans, United States.
- Li, Haonan, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. 2023. Bactrian-x : A multilingual replicable instruction-following model with low-rank adaptation.
- Lin, Peiqin, Chengzhi Hu, Zheyu Zhang, André F. T. Martins, and Hinrich Schütze. 2023. mglm-sim: Unveiling better cross-lingual similarity and transfer in multilingual pretrained language models. *CoRR*, abs/2305.13684.
- Lin, Xi Victoria, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shrutu Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona T. Diab, Veselin Stoyanov, and Xian Li. 2021. Few-shot learning with multilingual language models. *CoRR*, abs/2112.10668.
- Lin, Yu-Hsiang, Chian-Yu Chen, Jean Lee, Zirui Li, Yuyan Zhang, Mengzhou Xia, Shrutu Rijhwani, Junxian He, Zhisong Zhang, Xuezhe Ma, Antonios Anastasopoulos, Patrick Littell, and Graham Neubig. 2019. Choosing transfer languages for cross-lingual learning. In *The 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, Florence, Italy.
- Littell, Patrick, David R Mortensen, Ke Lin, Katherine Kairis, Carlisle Turner, and Lori Levin. 2017. Uriel and lang2vec: Representing languages as typological, geographical, and phylogenetic vectors. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, volume 2, pages 8–14.
- Lyu, Chenyang, Minghao Wu, and Alham Aji. 2024. Beyond probabilities: Unveiling the misalignment in evaluating large language models. In *Proceedings of the 1st Workshop on Towards Knowledgeable Language Models (KnowLLM 2024)*, pages 109–131, Association for Computational Linguistics, Bangkok, Thailand.
- Malaviya, Chaitanya, Graham Neubig, and Patrick Littell. 2017. Learning language representations for typology prediction. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Copenhagen, Denmark.
- Moran, Steven and Daniel McCloy, editors. 2019. *PHOIBLE 2.0*. Max Planck Institute for the Science of Human History, Jena.
- Muennighoff, Niklas, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2022. Crosslingual generalization through multitask finetuning.
- Ògúnrémí, Tolúlopé, Dan Jurafsky, and Christopher Manning. 2023. Mini but mighty: Efficient multilingual pretraining with linguistically-informed data selection. In *Findings*.
- Pfeiffer, Jonas, Naman Goyal, Xi Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. 2022. Lifting the curse of multilinguality by pre-training modular transformers. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3479–3495, Association for Computational Linguistics, Seattle, United States.
- Ponti, Edoardo Maria, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. XCOPA: A multilingual dataset for causal commonsense reasoning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2362–2376, Association for Computational Linguistics, Online.
- Scao, Teven Le, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilic, Daniel

- Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, and et al. 2022. BLOOM: A 176b-parameter open-access multilingual language model. *CoRR*, abs/2211.05100.
- Shaham, Uri, Jonathan Herzig, Roei Aharoni, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. Multilingual instruction tuning with just a pinch of multilinguality. *ArXiv*, abs/2401.01854.
- Shliazhko, Oleh, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. 2022. mgpt: Few-shot learners go multilingual.
- Taori, Rohan, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Tikhonov, Alexey and Max Ryabinin. 2021. It’s all in the heads: Using attention heads as a baseline for cross-lingual transfer in commonsense reasoning.
- Üstün, Ahmet, Arianna Bisazza, Gosse Bouma, and Gertjan van Noord. 2020. UDapter: Language adaptation for truly Universal Dependency parsing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2302–2315, Association for Computational Linguistics, Online.
- Üstün, Ahmet, Arianna Bisazza, Gosse Bouma, Gertjan van Noord, and Sebastian Ruder. 2022. Hyper-X: A unified hypernetwork for multi-task multilingual transfer. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7934–7949, Association for Computational Linguistics, Abu Dhabi, United Arab Emirates.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wang, Mingyang, Heike Adel, Lukas Lange, Jannik Strötgen, and Hinrich Schütze. 2023. Gradsim: Gradient-based language grouping for effective multilingual training. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 4631–4646, Association for Computational Linguistics.
- Xue, Linting, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Association for Computational Linguistics, Online.
- Yang, Yinfei, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. PAWS-X: A cross-lingual adversarial dataset for paraphrase identification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3687–3692, Association for Computational Linguistics, Hong Kong, China.
- Yong, Zheng-Xin, Hailey Schoelkopf, Niklas Muennighoff, Alham Fikri Aji, David Ifeoluwa Adelani, Khalid Almubarak, M Saiful Bari, Lintang Sutawika, Jungo Kasai, Ahmed Baruwa, Genta Indra Winata, Stella Biderman, Dragomir R. Radev, and Vassilina Nikoulina. 2022. Bloom+1: Adding language support to bloom for zero-shot prompting. *ArXiv*, abs/2212.09535.
- Zhou, Li, Antonia Karamolegkou, Wenyu Chen, and Daniel Hershcovich. 2023. Cultural compass: Predicting transfer learning success in offensive language detection with cultural features. *ArXiv*, abs/2310.06458.