

---

# Breaking Physical and Linguistic Borders: Privacy-Preserving Multilingual Prompt Tuning for Low-Resource Languages

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Pretrained large language models (LLMs) have emerged as a cornerstone in modern  
2 natural language processing, with their utility expanding to various applications  
3 and languages. However, the fine-tuning of multilingual LLMs, particularly for  
4 low-resource languages, is fraught with challenges stemming from data-sharing  
5 restrictions (the physical border) and from the inherent linguistic differences (the  
6 linguistic border). These barriers hinder users of various languages, especially  
7 those in low-resource regions, from fully benefiting from the advantages of LLMs.  
8 To address these challenges, we propose the Federated Prompt Tuning Paradigm  
9 for multilingual scenarios, which utilizes parameter-efficient fine-tuning while  
10 adhering to privacy restrictions. We have designed a comprehensive set of experi-  
11 ments and analyzed them using a novel notion of language distance to underscore  
12 the strengths of this paradigm: Even under computational constraints, our method  
13 not only bolsters data efficiency but also facilitates mutual enhancements across  
14 languages, particularly benefiting low-resource ones. Compared to traditional local  
15 cross-lingual transfer tuning methods, our approach achieves 6.9% higher accuracy,  
16 reduces the training parameters by over 99%, and demonstrates stronger cross-  
17 lingual generalization. Such findings underscore the potential of our approach to  
18 promote social equality, ensure user privacy, and champion linguistic diversity.

## 19 1 Introduction

20 Pretrained large language models (LLMs) have been driving the recent progress in natural language  
21 processing [11, 14, 3, 56, 57]. These large models, built on extensive corpora, offer valuable insights  
22 and impressive results across a range of applications. At the meantime, in order to provide universally  
23 accessible knowledge with LLMs, extending them to multiple languages has become a particularly  
24 relevant research target [17, 16, 5, 45].

25 However, finetuning and deploying multilingual LLMs in practical downstream tasks are not as easy  
26 as its monolingual counterpart. First of all, sharing data across different regions can be difficult  
27 or even impossible. Regulations like the General Data Protection Regulation (GDPR) [32] limit  
28 cross-region data-sharing. Moreover, languages in various regions can be radically different, e.g.  
29 Sino-Tibetan and Indo-European, posing a Non-Independent and Identically Distributed (non-IID)  
30 challenge when learning a global multilingual model. This situation accentuates privacy concerns,  
31 and highlights the need for effective privacy-preserving techniques when using multilingual LLMs.  
32 To this end, recent works attempt to address privacy-constrained fine-tuning for multilingual tasks  
33 and explore how different languages impact the federated process [60]. However, they primarily  
34 target high-resources languages; research on low-resource languages remains largely unexplored.

35 Addressing low-resource languages is essential to promoting technological fairness and protecting the  
36 linguistic diversity. Unlike their high-resource counterparts, low-resources languages pose intriguing

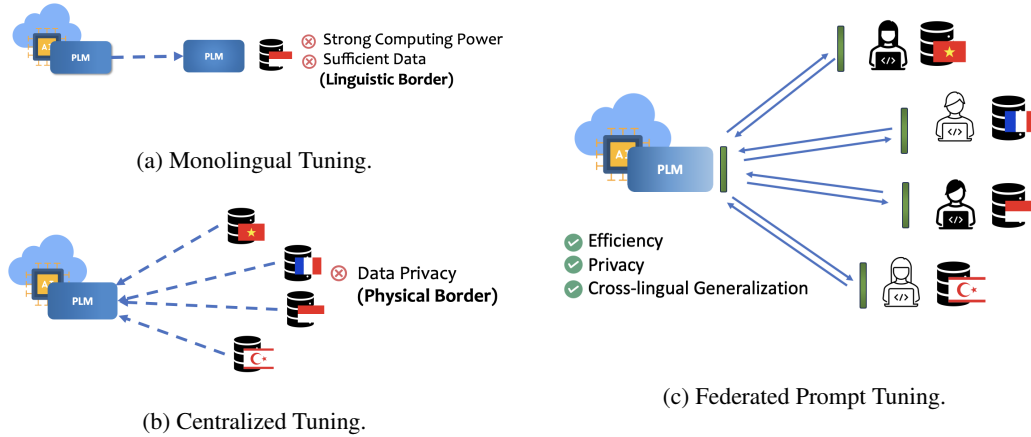


Figure 1: Comparison of three different fine-tuning paradigms for multilingual tasks.

37 research challenges: i) **Limited computational resources.** Regions of low-resources languages are  
 38 often economically developing areas, with little access to huge computational resources required to  
 39 either train language models from scratch or fully fine-tune pre-trained large language models [38, 1].  
 40 ii) **Limited data in the target language.** Due to a small speaking population or the spoken nature  
 41 of the language, data is often scarce [2, 43, 20]. As depicted in Figure 4, the pretraining data for  
 42 LLMs is predominantly in English, with little coverage of low-resource languages. Under such  
 43 circumstances, the performance of low-resources languages is often unsatisfactory during fine-tuning  
 44 because of their under-representation. iii) **Memorization risk.** Recent studies find that as pre-trained  
 45 models scale up, their ability to memorize training data increases [55]. This implies that, when  
 46 fine-tuning these models with limited data, the risk of overfitting and potential privacy issues arises.

47 To counteract the above challenges, we turn to federated learning (FL), where the model training is  
 48 done across multiple decentralized devices or servers while the data is always kept localized [42,  
 49 29, 65]. In a multilingual setting, FL becomes particularly natural, as data from diverse linguistic  
 50 backgrounds can be sourced without compromising user privacy, and due to the geographical spread  
 51 and inherent linguistic diversity of devices, data on each node is likely to exhibit non-IID distribution.

52 In this paper, in order to mitigate the physical border and the linguistic border of multilinguality,  
 53 we propose a new paradigm grounded in FL, Multilingual Federated Prompt Tuning, focusing on  
 54 parameter-efficient fine-tuning for multilingual tasks across various regions or devices. Specifically,  
 55 our global encoder can discern language patterns and cluster languages via federated prompt aver-  
 56 aging, which allows each client to benefit from others' data without direct access. This strategy  
 57 requires minimal computational resources and significantly improves performance, particularly for  
 58 low-resource languages. We demonstrate the effectiveness of our method on standard NLP tasks  
 59 including New Classification and XNLI. The performance of our paradigm achieves 6.9% accuracy  
 60 improvement while protecting the privacy of multilingual source data. Compared with other Fed-  
 61 erated Learning approaches, our paradigm reduces computational cost and communication cost by  
 62 more than 99%. Our approach paves the way for fine-tuning multilingual large language models on  
 63 resource-constraint devices in a privacy-preserving way, and holds the potential to promote social  
 64 equality, privacy, and linguistic diversity in the research community.

## 65 2 A New Paradigm for Multilinguality: Federated Prompt Tuning

### 66 2.1 Notation and Preliminaries

67 In our federated learning setup, we have  $K$  clients. Each client  $k$  has a private dataset, either  
 68 monolingual or multilingual, defined as:

$$\mathcal{D}_k = \{(x_{k,i}, y_{k,i}) \mid i = 1, \dots, n_k\}$$

69 where  $x_{k,i}$  denotes the textual content, and  $y_{k,i}$  is its corresponding label. The server sets up and  
 70 maintains a global prompt encoder  $h_g$ . Conversely, each client  $k$  has its version,  $h_k$ , adjusted based  
 71 on its dataset. Each prompt encoder, whether global or local, is composed of a series of trainable  
 72 prompts:  $h_0, h_1, h_2, \dots$ . These prompts are adjusted during training to better aid the model.

## 73 2.2 Virtual Prompt Encoder

74 Prompt Learning is a parameter-efficient alternative to fine-tuning pretrained language models (PLM).  
75 Instead of selecting discrete text prompts in a manual or automated fashion, prompt tuning utilize  
76 virtual prompt embeddings that can be optimized via gradient descent.

77 Given the utilization of a prompt encoder, for instance the version  $h_k$  for client  $k$ , a textual prompt  $p_k$   
78 tailored for a specific task can be generated. This prompt is subsequently concatenated or combined  
79 in another manner with the original input  $x$ , resulting in a modified input  $x'$ :

$$x'_k = p_k \oplus x \quad (1)$$

80 Here,  $p_k$  represents the prompt generated by  $h_k$ .

81 The modified input  $x'_k$  is then processed by the encoder  $E$  of the pre-trained language model:

$$h'_k = E(x'_k) \quad (2)$$

82 The primary objective of each prompt encoder is to generate an effective prompt, such as  $p_k$  for client  
83  $k$ , to guide the pre-trained model in producing the desired outputs. During the fine-tuning phase,  
84 based on a task-specific loss, the parameters of the prompt encoder  $h_k$  are often adjusted:

$$\mathcal{L}(x, y; h_k) = \text{Loss}(D(E(p_k \oplus x)), y) \quad (3)$$

85 Where  $D$  is a decoder that maps the internal representation to task outputs, and Loss is an appropriate  
86 loss function, like cross-entropy loss. Throughout the fine-tuning, both the model’s parameters and  
87 the prompt encoder  $h_k$ ’s parameters are updated in accordance with this loss function.

## 88 2.3 Federated Prompt Averaging

89 Derived from FedAvg [42] mentioned in Appendix B, we propose the following federated prompt  
90 averaging algorithm (also shown in Algorithm 1):

91 **Initialization:** The server initializes the global prompt encoder with its prompts  $h_0, h_1, \dots$ . Each  
92 client sets its local version tailored to its dataset.

93 **Client Selection:** Every communication round selects a fraction  $C$  of the total  $K$  clients for training.  
94 This subset is  $m = \max(C \times K, 1)$ . A subset  $S$  of  $m$  clients is chosen.

95 **Local Model Training and Tuning:** Each client  $k$  in  $S$  goes through: The client fetches the current  
96 global prompt encoder. It assembles a model using its local prompt encoder with prompts  $h_0, h_1, \dots$   
97 and PLM parameters. Training on  $\mathcal{D}_k$  fine-tunes the local prompt encoder and its prompts, while  
98 most of the PLM remains unchanged. After training, each client computes model updates, especially  
99 the refined local prompt encoder and its prompts.

100 **Aggregation:** The server aggregates updates from all clients. The global prompt encoder and its  
101 prompts  $h_0, h_1, \dots$  are updated using:

$$h_g = \frac{1}{m} \sum_{k=1}^m h_k$$

102

## 103 3 Evaluation: How Federated Prompt Tuning helps Multilinguality

### 104 3.1 Experimental Setup

105 **Tasks and Datasets** We evaluate our model using the popular XGLUE benchmark [31], a cross-  
106 lingual evaluation benchmark for our multilingual evaluation. We conduct our experiments on two  
107 classification tasks: News Classification (NC) and XNLI [18], covering both high-resource languages  
108 and low-resource languages, with details provided in Appendix D. Accuracy (ACC) of the multi-  
109 class classification is used as the metric for both of the tasks. Our base model for both tasks is the  
110 XLM-RoBERTa base-sized model (270M parameters) [16].

Method	en	es	fr	de	ru	Avg
Monolingual	92.4	84.7	79.5	88.3	89.0	86.8
Centralized	93.9	86.7	<b>82.9</b>	<b>89.5</b>	88.6	88.3
FL (IID)	<b>94.1</b>	<b>86.9</b>	82.7	89.4	<b>88.8</b>	<b>88.4</b>
FL (Non-IID)	92.4	86.3	81.2	88.9	84.7	86.7
PE_Monolingual	82.9	59.7	47.3	71.4	60.0	64.3
PE_Centralized	89.1	76.2	67.4	78.8	75.9	77.5
PE_FL (IID) (Ours)	<b>91.2</b>	<b>82.2</b>	<b>76.5</b>	<b>86.4</b>	<b>81.6</b>	<b>83.6</b>
PE_FL (Non-IID) (Ours)	87.8	79.2	73.7	83.1	79.5	80.7

Table 1: Results for FL experiments on the NC task. Bold scores indicate the best in the column.

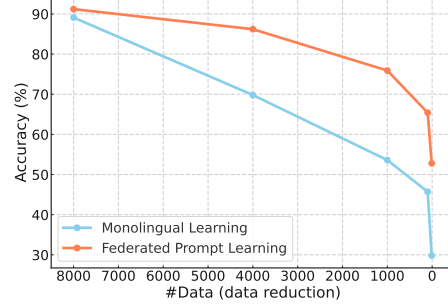


Figure 2: Performance comparison on NC with decreasing dataset size.

Table 2: Results for FL experiments on the XNLI task. Bold scores indicate the best in the column. The PE\_FL is evaluated under the Non-IID setting.

Method	en	fr	es	de	el	bg	ru	tr	ar	vi	th	zh	hi	sw	ur	Avg
Monolingual	39.1	35.1	36.6	35.7	35.3	35.9	35.5	26.2	32.1	31.7	31.5	33.7	31.6	26.0	28.1	32.94
Centralized	35.3	36.9	33.3	35.3	30.5	36.5	33.7	35.7	33.3	<b>40.1</b>	36.1	30.5	37.3	<b>38.6</b>	29.3	34.86
PE_FL (Ours)	<b>43.2</b>	<b>40.6</b>	<b>42.9</b>	<b>40.2</b>	<b>39.7</b>	<b>40.8</b>	<b>41.1</b>	<b>37.6</b>	<b>39.1</b>	39.9	<b>39.4</b>	<b>39.8</b>	<b>38.2</b>	37.1	<b>37.8</b>	<b>39.83</b>

111 **Baselines** 1) *Monolingual Fine-tuning (No FL)*: Traditional local fine-tuning where a separate  
 112 model is finetuned using the corresponding dataset for each single language; 2) *Centralized Fine-*  
 113 *tuning (No FL)*: Standard Fine-tuning using a combined dataset of all languages centralized in one  
 114 location; 3) *Full Fine-tuning with FL*: Directly fine-tuning the whole pre-trained language model in a  
 115 federated manner, with a full pre-trained model on each client; 4) *Prompt Fine-tuning with FL*: Only  
 116 training the prompt encoder in a federated manner, with a prompt encoder on each client.

117 For FL experiments, we adjust the parameter  $\alpha$  that controls the mixture of languages in the dataset.  
 118 An  $\alpha$  value of 1.0 signifies a uniform mixture of all languages, while values closer to 0 indicate a  
 119 dominant representation of individual languages or a more separated mixture.

### 120 3.2 Main Results

121 Table 1 presents the outcomes of experiments focused on news classification. When employing  
 122 Prompt Tuning in comparison to Full Finetuning, there is an acceptable decline in accuracy. Despite  
 123 this decrease, the overall performance remains consistent and stable. A significant gain in accuracy is  
 124 observed when adopting the FL approach. It is worth noting that the fine-tuning time is considerably  
 125 reduced when employing the Prompt Tuning method as opposed to without it. For a comprehensive  
 126 analysis of this, refer to the section 3.5.

127 Table 2 summarises the results of our FL experiments on the XNLI task. To accentuate the potency  
 128 of our Federated Prompt Tuning approach, a juxtaposition was made with traditional monolingual  
 129 training. As the data portrays, our Federated Prompt Tuning, particularly on Non-IID dataset,  
 130 consistently outperformed the monolingual method across all languages. Remarkably, this superior  
 131 performance was maintained even for languages with limited available data. The average score  
 132 further substantiates the prowess of Federated Prompt Tuning, marking a noticeable improvement  
 133 from 32.94% in the monolingual approach to 39.83% with Non-IID Federated Prompt Tuning.

134 From our results in section 3.2, we observe that some languages demonstrate superior accuracy with  
 135 the FL method compared to the centralized approach. This enhanced performance might be attributed  
 136 to the Federated Prompt Averaging in FL, which could introduce similar implicit regularization  
 137 effects. Additionally, the prompt encoder serves as a parameter-efficient alternative. By freezing the  
 138 core language model parameters, we prevent the model from altering its foundational understanding  
 139 of language. As a result, the model’s propensity to overfit to a dataset is reduced, minimizing the risk  
 140 of memorizing specific lexical cues and spurious correlations.

### 141 3.3 Ablation Study I: Data Efficiency

142 As previous sections mentioned, one characteristic of low-resource languages is their limited available  
 143 data. Hence, enhancing data sample efficiency is crucial when fine-tuning pre-trained models for

144 downstream tasks. To better validate and simulate the advantages of our approach in real-world  
 145 scenarios, we reduced the data volume for one language and observed the performance under  
 146 traditional local fine-tuning as well as our Federated prompt fine-tuning method. We conducted  
 147 experiments on German News Classification. German was chosen because it represents the language  
 148 with the fewest resources among the five languages included in this task.

149 As shown in the Figure 2, our Federated Prompt Tuning method consistently outperforms the  
 150 traditional monolingual approach. As we reduce the dataset size from 8,000 to near-zero, the accuracy  
 151 of the traditional method drops significantly. On the other hand, the Federated Prompt Tuning method  
 152 retains its performance, demonstrating its robustness even with limited data. This clearly indicates  
 153 that our Federated Prompt Tuning approach is better suited for scenarios with limited data availability.

### 154 3.4 Ablation Study II: Language Distance

155 As previously mentioned, another characteristic of low-resource languages is that their linguistic  
 156 features differ from those of high-resource languages, particularly in aspects including syntax,  
 157 phonology, and inventory. Consequently, direct fine-tuning on models pre-trained with highly  
 158 dissimilar languages often yields unsatisfying results. Therefore, we conducted an ablation study  
 159 to examine the impact of language similarity on performance, comparing our Federated Prompt  
 160 fine-tuning method to the traditional local fine-tuning approach.

161 We define the *pretrained language* as a representative composite language formed by blending each  
 162 language in the multilingual corpus used for pre-training, in proportion to their amount. This serves as  
 163 a formal representation for the mixed dataset composition. We define distance for a specific language  
 164 in the downstream tasks, in terms of the negative logarithm of its cosine similarity to the pre-trained  
 165 language. More details are shown in Appendix F. Leveraging the distance metric, we compared model  
 166 performance of languages with varying degrees of distance to the pre-trained language. We present  
 167 our results from two key experiments on the NC and XNLI tasks. From Figure 3, a conspicuous trend  
 168 is observed: As the language similarity to the pre-trained language decreases, the model’s accuracy  
 169 tends to drop. However, when we apply our Federated Prompt method, this decline is notably less  
 170 steep. This means that even when we are dealing with languages that are quite different from the  
 171 pre-trained one, our method manages to retain a decent level of accuracy. The difference between our  
 172 method and the traditional local finetuning becomes even more obvious for languages with less data,  
 173 indicating that our Federated Prompt Tuning method offers significant advantages, particularly in  
 174 low-resource scenarios.

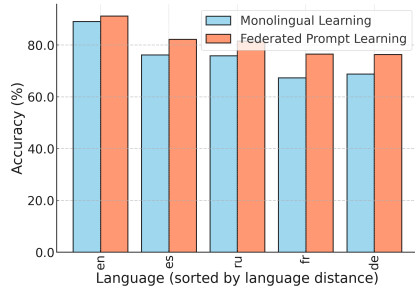
### 175 3.5 Ablation Study III: Parameter Efficiency

176 **Computational Cost** From the perspective of trainable parameters, this significant reduction in  
 177 parameters demonstrates exceptional parameter efficiency. In both of the tasks, despite the total  
 178 number of parameters exceeding 278 million, the trainable parameters are only around 1.2 million,  
 179 accounting for less than 0.5% of the total. Such a design can substantially reduce training time and  
 180 computational resources, while also mitigating the risk of overfitting. In the context of LLMs, this high  
 181 parameter efficiency offers potential for model deployment in resource-constrained environments.

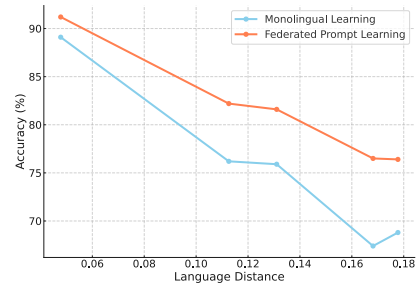
Table 3: Comparison of parameter efficiency and communication overhead in NC task.

	Federated Full Finetuning	Federated Prompt Tuning (Ours)	Optimization Scale
Trainable Params	278,655,764	1,202,708	<b>231.69</b>
Communication Cost	108GB	478.93MB	

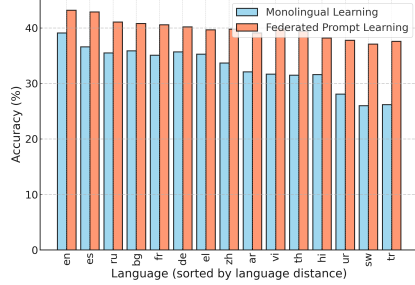
182 **Communication Cost** XLM-Roberta-Base’s data transmission in Federated Learning with 5 clients  
 183 and 10 communication rounds was 108 GB. After our optimization, using a prompt encoder with a  
 184 2x768 structure, the transmission size reduced to 478.93 MB, a 99% reduction shown in Table 3.  
 185 This optimization enhances efficiency in federated prompt tuning and expands its applicability to  
 186 bandwidth-constrained environments like edge devices and mobile networks.



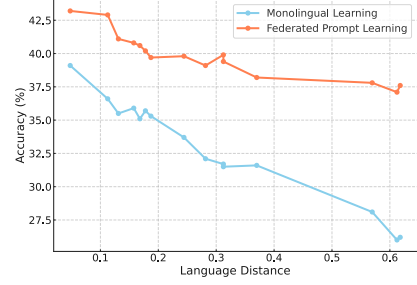
(a) Finetuning accuracy across different languages on the NC task.



(b) Finetuning accuracy across languages with varying similarity to the pre-trained language on the NC task.



(c) Finetuning accuracy across different languages on the XNLI task.



(d) Finetuning accuracy across languages with varying similarity to the pre-trained language on the XNLI task.

Figure 3: Comparative performance of traditional local finetuning and our Federated Prompt Tuning method across languages with varying similarity to the pre-trained language for XNLI and NC.

187 **4 Conclusion**

188 **Future work** Privacy attacks have been discussed in (author?) [21] on how gradient inversion  
 189 can be used to attack language models and break the privacy protection that FL naturally adds.  
 190 Future work may include privacy experiments and additional privacy protection with various secure  
 191 aggregation (SA) [10, 7, 49] and differential privacy (DP) techniques [59, 50]. Future work will also  
 192 explore the impact on the Multilingual Federated Prompt Tuning method as models scale up.

193 **Social Impacts** Addressing the physical and linguistic challenges of multilingual LLMs, especially  
 194 for low-resource languages, requires innovative approaches that can balance efficiency, privacy  
 195 concerns, and performance. Our Multilingual Federated Prompt Tuning paradigm provides a solution  
 196 to these challenges. By aggregating lightweight multilingual prompts, this approach offers enhanced  
 197 fine-tuning capabilities with minimal computational demand. The robustness of our method is  
 198 especially pronounced for low-resource languages with sparse data and rare linguistic features. Its  
 199 potential to democratize access to technology, preserve linguistic diversity, and ensure user privacy  
 200 can have profound implications for the future of technology and society.

201 **References**

202 [1] Ife Adebara and Muhammad Abdul-Mageed. Towards afrocentric NLP for African languages:  
 203 Where we are and where we can go. In *Proceedings of the 60th Annual Meeting of the*  
 204 *Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3814–3841, Dublin,  
 205 Ireland, May 2022. Association for Computational Linguistics.

206 [2] David Ifeoluwa Adelani, Jade Abbott, Graham Neubig, Daniel D’souza, Julia Kreutzer, Con-  
 207 stantine Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, et al.  
 208 Masakhaner: Named entity recognition for african languages. *Transactions of the Association*  
 209 *for Computational Linguistics*, 9:1116–1131, 2021.

210 [3] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,  
 211 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report.  
 212 *arXiv preprint arXiv:2305.10403*, 2023.

- 213 [4] Alan Ansell, Edoardo Ponti, Anna Korhonen, and Ivan Vulić. Composable sparse fine-tuning  
214 for cross-lingual transfer. In *Proceedings of the 60th Annual Meeting of the Association for*  
215 *Computational Linguistics (Volume 1: Long Papers)*, pages 1778–1796, 2022.
- 216 [5] Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of  
217 monolingual representations. In *Proceedings of the 58th Annual Meeting of the Association for*  
218 *Computational Linguistics*, pages 4623–4637, Online, July 2020. Association for Computational  
219 Linguistics.
- 220 [6] Syreen Banabilah, Moayad Aloqaily, Eitaa Alsayed, Nida Malik, and Yaser Jararweh. Federated  
221 learning review: Fundamentals, enabling technologies, and future applications. *Information*  
222 *Processing Management*, 59(6):103061, 2022.
- 223 [7] James Henry Bell, Kallista A Bonawitz, Adrià Gascón, Tancrede Lepoint, and Mariana Raykova.  
224 Secure single-server aggregation with (poly) logarithmic overhead. In *Proceedings of the 2020*  
225 *ACM SIGSAC Conference on Computer and Communications Security*, pages 1253–1269, 2020.
- 226 [8] Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel. BitFit: Simple parameter-efficient  
227 fine-tuning for transformer-based masked language-models. In *Proceedings of the 60th Annual*  
228 *Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1–9,  
229 Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 230 [9] Daniel J. Beutel, Taner Topal, Akhil Mathur, Xinchu Qiu, Javier Fernandez-Marques, Yan  
231 Gao, Lorenzo Sani, Kwing Hei Li, Titouan Parcollet, Pedro Porto Buarque de Gusmão, and  
232 Nicholas D. Lane. Flower: A friendly federated learning research framework, 2020.
- 233 [10] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan,  
234 Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for  
235 federated learning on user-held data. *arXiv preprint arXiv:1611.04482*, 2016.
- 236 [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
237 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
238 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 239 [12] Ethan C. Chau, Lucy H. Lin, and Noah A. Smith. Parsing with multilingual BERT, a small  
240 corpus, and a small treebank. In *Findings of the Association for Computational Linguistics:*  
241 *EMNLP 2020*, pages 1324–1334, Online, November 2020. Association for Computational  
242 Linguistics.
- 243 [13] Yihong Chen, Kelly Marchisio, Roberta Raileanu, David Ifeoluwa Adelani, Pontus Stenetor,  
244 Sebastian Riedel, and Mikel Artetx. Improving language plasticity via pretraining with active  
245 forgetting. *NeurIPS 2023*, 2023.
- 246 [14] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
247 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:  
248 Scaling language modeling with pathways. 2022.
- 249 [15] Hyung Won Chung, Dan Garrette, Kiat Chuan Tan, and Jason Riesa. Improving multilingual  
250 models with language-clustered vocabularies. In *Proceedings of the 2020 Conference on*  
251 *Empirical Methods in Natural Language Processing (EMNLP)*, pages 4536–4546, Online,  
252 November 2020. Association for Computational Linguistics.
- 253 [16] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek,  
254 Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Un-  
255 supervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual*  
256 *Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July 2020.  
257 Association for Computational Linguistics.
- 258 [17] Alexis Conneau and Guillaume Lample. *Cross-Lingual Language Model Pretraining*. Curran  
259 Associates Inc., Red Hook, NY, USA, 2019.
- 260 [18] Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R. Bowman, Holger  
261 Schwenk, and Veselin Stoyanov. Xnli: Evaluating cross-lingual sentence representations. In  
262 *EMNLP*, 2018.

- 263 [19] Sumanth Doddapaneni, Gowtham Ramesh, Anoop Kunchukuttan, Pratyush Kumar, and  
264 Mitesh M. Khapra. A primer on pretrained multilingual language models. *CoRR*,  
265 abs/2107.00676, 2021.
- 266 [20] Abteen Ebrahimi, Manuel Mager, Arturo Oncevay, Vishrav Chaudhary, Luis Chiruzzo, Angela  
267 Fan, John Ortega, Ricardo Ramos, Annette Rios, Ivan Vladimir Meza Ruiz, Gustavo Giménez-  
268 Lugo, Elisabeth Mager, Graham Neubig, Alexis Palmer, Rolando Coto-Solano, Thang Vu,  
269 and Katharina Kann. AmericasNLI: Evaluating zero-shot natural language understanding of  
270 pretrained multilingual models in truly low-resource languages. In *Proceedings of the 60th  
271 Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
272 pages 6279–6299, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 273 [21] Liam H. Fowl, Jonas Geiping, Steven Reich, Yuxin Wen, Wojtek Czaja, Micah Goldblum, and  
274 Tom Goldstein. Decepticons: Corrupted transformers breach privacy in federated learning for  
275 language models. *ArXiv*, abs/2201.12675, 2022.
- 276 [22] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe,  
277 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning  
278 for NLP. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th  
279 International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning  
280 Research*, pages 2790–2799. PMLR, 09–15 Jun 2019.
- 281 [23] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and  
282 Weizhu Chen. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685,  
283 2021.
- 284 [24] Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin John-  
285 son. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual  
286 generalization. *CoRR*, abs/2003.11080, 2020.
- 287 [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua  
288 Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations,  
289 ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- 290 [26] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh,  
291 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv  
292 preprint arXiv:1610.05492*, 2016.
- 293 [27] Anne Lauscher, Vinit Ravishanker, Ivan Vulić, and Goran Glavaš. From zero to hero: On the  
294 limitations of zero-shot language transfer with multilingual Transformers. In *Proceedings of  
295 the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages  
296 4483–4499, Online, November 2020. Association for Computational Linguistics.
- 297 [28] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient  
298 prompt tuning. *CoRR*, abs/2104.08691, 2021.
- 299 [29] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia  
300 Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning  
301 and Systems*, 2:429–450, 2020.
- 302 [30] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.  
303 In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics  
304 and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long  
305 Papers)*, pages 4582–4597, Online, August 2021. Association for Computational Linguistics.
- 306 [31] Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun  
307 Shou, Daxin Jiang, Guihong Cao, Xiaodong Fan, Ruofei Zhang, Rahul Agrawal, Edward Cui,  
308 Sining Wei, Taroon Bharti, Ying Qiao, Jiun-Hung Chen, Winnie Wu, Shuguang Liu, Fan Yang,  
309 Daniel Campos, Rangan Majumder, and Ming Zhou. Xglue: A new benchmark dataset for  
310 cross-lingual pre-training, understanding and generation. *arXiv*, abs/2004.01401, 2020.
- 311 [32] Wei Yang Bryan Lim, Nguyen Cong Luong, Dinh Thai Hoang, Yutao Jiao, Ying-Chang Liang,  
312 Qiang Yang, Dusit Niyato, and Chunyan Miao. Federated learning in mobile edge networks: A  
313 comprehensive survey. *IEEE Communications Surveys & Tutorials*, 22(3):2031–2063, 2020.



- 314 [33] Bill Yuchen Lin, Chaoyang He, Zihang Ze, Hulin Wang, Yufen Hua, Christophe Dupuy, Rahul  
315 Gupta, Mahdi Soltanolkotabi, Xiang Ren, and Salman Avestimehr. FedNLP: Benchmarking  
316 federated learning methods for natural language processing tasks. In *Findings of the Association  
317 for Computational Linguistics: NAACL 2022*, pages 157–175, Seattle, United States, July 2022.  
318 Association for Computational Linguistics.
- 319 [34] Patrick Littell, David R Mortensen, Ke Lin, Katherine Kairis, Carlisle Turner, and Lori Levin.  
320 Uriel and lang2vec: Representing languages as typological, geographical, and phylogenetic  
321 vectors. In *Proceedings of the 15th Conference of the European Chapter of the Association for  
322 Computational Linguistics: Volume 2, Short Papers*, volume 2, pages 8–14, 2017.
- 323 [35] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig.  
324 Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language  
325 processing. *CoRR*, abs/2107.13586, 2021.
- 326 [36] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt  
327 understands, too. *arXiv:2103.10385*, 2021.
- 328 [37] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International  
329 Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.  
330 OpenReview.net, 2019.
- 331 [38] Manuel Mager, Arturo Oncevay, Abteen Ebrahimi, John Ortega, Annette Rios Gonzales,  
332 Angela Fan, Ximena Gutierrez-Vasques, Luis Chiruzzo, Gustavo Giménez-Lugo, Ricardo  
333 Ramos, et al. Findings of the americasnlp 2021 shared task on open machine translation  
334 for indigenous languages of the americas. In *Proceedings of the First Workshop on Natural  
335 Language Processing for Indigenous Languages of the Americas*, pages 202–217, 2021.
- 336 [39] Chaitanya Malaviya, Graham Neubig, and Patrick Littell. Learning language representations  
337 for typology prediction. In *Conference on Empirical Methods in Natural Language Processing  
338 (EMNLP)*, Copenhagen, Denmark, September 2017.
- 339 [40] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, and Sayak Paul.  
340 Peft: State-of-the-art parameter-efficient fine-tuning methods. [https://github.com/  
341 huggingface/peft](https://github.com/huggingface/peft), 2022.
- 342 [41] Kelly Marchisio, Patrick Lewis, Yihong Chen, and Mikel Artetxe. Mini-model adaptation:  
343 Efficiently extending pretrained models to new languages via aligned shallow training. *ACL  
344 2023, Findings of the Association for Computational Linguistics*, 2022.
- 345 [42] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas.  
346 Communication-efficient learning of deep networks from decentralized data. In *Artificial  
347 intelligence and statistics*. PMLR, 2017.
- 348 [43] Shamsuddeen Hassan Muhammad, David Ifeoluwa Adelani, Sebastian Ruder, Ibrahim Said  
349 Ahmad, Idris Abdulmumin, Bello Shehu Bello, Monojit Choudhury, Chris Chinenye Emezue,  
350 Saheed Salahudeen Abdullahi, Anuluwapo Aremu, et al. Naijasenti: A nigerian twitter  
351 sentiment corpus for multilingual sentiment analysis. *arXiv preprint arXiv:2201.08277*, 2022.
- 352 [44] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
353 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
354 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
355 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-  
356 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-  
357 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,  
358 volume 32. Curran Associates, Inc., 2019.
- 359 [45] Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. MAD-X: An Adapter-Based  
360 Framework for Multi-Task Cross-Lingual Transfer. In *Proceedings of the 2020 Conference  
361 on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7654–7673, Online,  
362 November 2020. Association for Computational Linguistics.

- 363 [46] Telmo Pires, Eva Schlinger, and Dan Garrette. How multilingual is multilingual BERT? In  
364 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages  
365 4996–5001, Florence, Italy, July 2019. Association for Computational Linguistics.
- 366 [47] Telmo Pires, Eva Schlinger, and Dan Garrette. How multilingual is multilingual BERT? In  
367 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages  
368 4996–5001, Florence, Italy, July 2019. Association for Computational Linguistics.
- 369 [48] Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna  
370 Korhonen. XCOPA: A multilingual dataset for causal commonsense reasoning. In *Proceedings*  
371 *of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,  
372 pages 2362–2376, Online, November 2020. Association for Computational Linguistics.
- 373 [49] Xinchi Qiu, Heng Pan, Wanru Zhao, Chenyang Ma, Pedro PB Gusmao, and Nicholas D Lane.  
374 vfedsec: Efficient secure aggregation for vertical federated learning via secure layer. *arXiv*  
375 *preprint arXiv:2305.16794*, 2023.
- 376 [50] Thilina Ranbaduge and Ming Ding. Differentially private vertical federated learning, 2022.
- 377 [51] Timo Schick and Hinrich Schütze. True few-shot learning with Prompts—A real-world per-  
378 spective. *Transactions of the Association for Computational Linguistics*, 10:716–731, 2022.
- 379 [52] Yu-An Chung Mariano Cora Meglioli David Dale Ning Dong Paul-Ambroise Duquenne Hady  
380 Elsahar Hongyu Gong Kevin Heffernan John Hoffman Christopher Klaiber Pengwei Li Daniel  
381 Licht Jean Maillard Alice Rakotoarison Kaushik Ram Sadagopan Guillaume Wenzek Ethan Ye  
382 Bapi Akula Peng-Jen Chen Naji El Hachem Brian Ellis Gabriel Mejia Gonzalez Justin Haaheim  
383 Prangthip Hansanti Russ Howes Bernie Huang Min-Jae Hwang Hirofumi Inaguma Somya Jain  
384 Elahe Kalbassi Amanda Kallet Iliia Kulikov Janice Lam Daniel Li Xutai Ma Ruslan Mavlyutov  
385 Benjamin Peloquin Mohamed Ramadan Abinesh Ramakrishnan Anna Sun Kevin Tran Tuan  
386 Tran Igor Tufanov Vish Vogeti Carleigh Wood Yilin Yang Bokai Yu Pierre Andrews Can  
387 Balioglu Marta R. Costa-jussà<sup>3</sup> Onur Celebi Maha Elbayad Cynthia Gao Francisco Guzmán  
388 Justine Kao Ann Lee Alexandre Mourachko Juan Pino Sravya Popuri Christophe Ropers  
389 Safiyah Saleem Holger Schwenk Paden Tomasello Changhan Wang Jeff Wang Skyler Wang  
390 Seamless Communication, Loïc Barrault. Seamlessm4t—massively multilingual & multimodal  
391 machine translation. *ArXiv*, 2023.
- 392 [53] Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. Auto-  
393 Prompt: Eliciting knowledge from language models with automatically generated prompts. In  
394 *Empirical Methods in Natural Language Processing (EMNLP)*, 2020.
- 395 [54] Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the*  
396 *22nd ACM SIGSAC conference on computer and communications security*, pages 1310–1321,  
397 2015.
- 398 [55] Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization  
399 without overfitting: Analyzing the training dynamics of large language models. *Advances in*  
400 *Neural Information Processing Systems*, 35:38274–38290, 2022.
- 401 [56] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-  
402 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open  
403 and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 404 [57] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,  
405 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open  
406 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 407 [58] Haoyu Wang, Handong Zhao, Yaqing Wang, Tong Yu, Jiuxiang Gu, and Jing Gao. Fedkc: Fed-  
408 erated knowledge composition for multilingual natural language understanding. In *Proceedings*  
409 *of the ACM Web Conference 2022, WWW ’22*, page 1839–1850, New York, NY, USA, 2022.  
410 Association for Computing Machinery.

- 411 [59] Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H Yang, Farhad Farokhi, Shi Jin, Tony QS  
412 Quek, and H Vincent Poor. Federated learning with differential privacy: Algorithms and  
413 performance analysis. *IEEE Transactions on Information Forensics and Security*, 15:3454–  
414 3469, 2020.
- 415 [60] Orion Weller, Marc Marone, Vladimir Braverman, Dawn J Lawrie, and Benjamin Van Durme.  
416 Pretrained models for multilingual federated learning. In *Proceedings of the 2022 Conference  
417 of the North American Chapter of the Association for Computational Linguistics: Human  
418 Language Technologies (NAACL-HLT)*, 2022.
- 419 [61] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony  
420 Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer,  
421 Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain  
422 Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art  
423 natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in  
424 Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020.  
425 Association for Computational Linguistics.
- 426 [62] Shijie Wu and Mark Dredze. Beto, bentz, becas: The surprising cross-lingual effectiveness  
427 of BERT. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Lan-  
428 guage Processing and the 9th International Joint Conference on Natural Language Processing  
429 (EMNLP-IJCNLP)*, pages 833–844, Hong Kong, China, November 2019. Association for  
430 Computational Linguistics.
- 431 [63] Shijie Wu and Mark Dredze. Are all languages created equal in multilingual BERT? In  
432 *Proceedings of the 5th Workshop on Representation Learning for NLP*, pages 120–130, Online,  
433 July 2020. Association for Computational Linguistics.
- 434 [64] Hu Xu, Lei Shu, Philip Yu, and Bing Liu. Understanding pre-trained BERT for aspect-based  
435 sentiment analysis. In *Proceedings of the 28th International Conference on Computational Lin-  
436 guistics*, pages 244–250, Barcelona, Spain (Online), December 2020. International Committee  
437 on Computational Linguistics.
- 438 [65] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept  
439 and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19,  
440 2019.

## 441 A Background

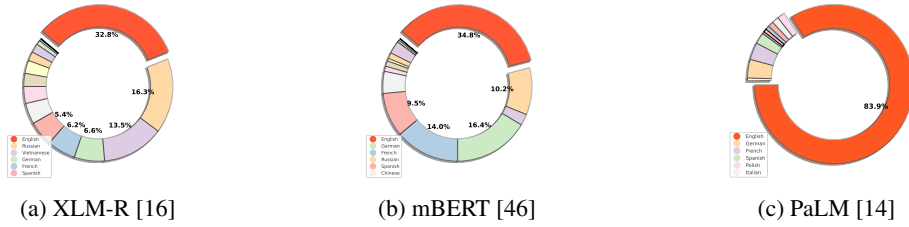


Figure 4: Linguistic coverage of different large language models.

## 442 B Related Work

443

444 **Federated Learning.** Federated Learning has garnered significant attention in the academic realm.  
445 A notable contribution from (author?) [26] underscores the potential of this methodology. One of its  
446 primary benefits is the execution of deep learning algorithms while maintaining an emphasis on user  
447 privacy, a premise originally posited by (author?) [54]. The overarching objective is to amalgamate  
448 insights from diverse data repositories without compromising sensitive particulars. In this spectrum,  
449 the FedAvg algorithm [42] stands out. This algorithm operates by independently training models on  
450 distinct client devices and subsequently aggregating their updates centrally for an averaged outcome.  
451 However, it’s imperative to acknowledge that FedAvg, while powerful, is principally structured for  
452 identically and independently distributed (IID) data, while its application on Non-IID datasets may  
453 lead to potential discrepancies in results or even model instability [42, 26]. Despite abundant research  
454 made on problems at hospitals, legal firms, and financial institutions, extending language models  
455 for multilingual usages effectively and efficiently, especially for low-resource languages remains  
456 under-explored.

457 In the general NLP domain, FL has been instrumental in tasks such as language modeling, sentiment  
458 analysis, and machine translation, showcasing its potential to revolutionize the way models are  
459 trained and deployed [6]. FedNLP introduces a benchmarking framework for evaluating various FL  
460 methods across NLP tasks, providing a universal interface between Transformer-based models and FL  
461 methods [33]. FedKC [58] is a federated approach designed for multilingual Natural Language Un-  
462 derstanding (NLU) that integrates knowledge from multiple data sources through federated learning  
463 techniques to enhance the efficacy and accuracy of multilingual text processing. However, considera-  
464 tions regarding computational and communication efficiency in resource-constrained environments  
465 have not been adequately addressed.

466 **Multilingual Language Models.** Multilingual Pretrained Language Models such as mBERT [46],  
467 XLM-R [16], and SeamlessM4T [52] have emerged as a viable option for bringing the power of  
468 pretraining to a large number of languages [19]. Many studies analyzed mBERT’s and XLM-R’s  
469 capabilities and limitations, finding that the multilingual models work surprisingly well for cross-  
470 lingual tasks, despite the fact that they do not rely on direct cross-lingual supervision (e.g., parallel or  
471 comparable data, translation dictionaries [47, 62, 5, 64]).

472 However, these multilingual LMs are not without limitations. Particularly, (author?) [16] observed  
473 the *curse of multilinguality* phenomenon: given a fixed model capacity, adding more languages  
474 does not necessarily improve the multilingual performance but can deteriorate the performance  
475 after a certain point, especially for underrepresented languages [63, 24, 27]. Prior work tried to  
476 address this issue by increasing the model capacity [5, 45, 12] or through additional training for  
477 particular language pairs [45, 48] or by clustering and merging the vocabularies of similar languages,  
478 before defining a joint vocabulary across all languages [15]. Despite these efforts, the multilingual  
479 LMs still struggle with balancing their capacity across many languages in an sample-efficient and  
480 parameter-efficient way [4, 41, 13].

481 **Prompt Learning and Parameter-Efficient Fine-Tuning.** The size of pre-trained language models  
 482 has been increasing significantly [11], presenting challenges to traditional task transfer based on  
 483 full-parameter finetuning. Recent research has shifted its attention to Parameter-Efficient Fine-Tuning  
 484 techniques, such as prompt tuning [28, 30, 36], adapters [22], and combined approaches [23, 8].  
 485 These methods utilize a minimal number of tuning parameters, yet they offer transfer performance  
 486 that is comparable with traditional finetuning.

487 Prompt learning is a burgeoning area in machine learning where models are steered towards desired  
 488 outputs using prompts, typically without exposure to explicit label information [35]. This paradigm  
 489 shows promise in effectively leveraging large pre-trained models in a data-efficient manner by  
 490 reducing the need for extensive labeled datasets [51]. Additionally, prompt learning has exhibited a  
 491 remarkable ability to generalize across a variety of tasks, suggesting a step towards more flexible and  
 492 adaptable machine learning systems [53].

## 493 C Federated Prompt Averaging Algorithm

Algorithm 1: Federated Prompt Averaging

1: **Initialization:**

- 2: Server initializes global prompt encoder  $h_g$
- 3: Each client initializes local prompt encoder  $h_k$

1: **Server executes:**

- 2: **for each** round  $t$  **do**
- 3:   Select subset  $S$  of  $m$  clients
- 4:   **for each** client  $k$  in  $S$  **do**
- 5:     Send  $h_g$  to client  $k$
- 6:   **end for**
- 7:   Aggregate client updates:
- 8:    $h_g = \frac{1}{m} \sum_{k=1}^m h_k$
- 9: **end for**

1: **Client  $k$  executes:**

- 2: Retrieve current  $h_g$
- 3: Assemble full model using  $h_k$  and PLM parameters
- 4: Train model on  $\mathcal{D}_k$
- 5: Update local prompt encoder  $h_k$
- 6: Send updated  $h_k$  to server

## 495 D Dataset

496 **News Classification (NC)** is a classification problem with 10 classes across 5 languages: English,  
 497 Spanish, French, German, and Russian. This task aims to predict the category given a news article.  
 498 Since only 10k annotated examples are available for each language (excluding the official test set),  
 499 we sample 8k instances for training and 1k for evaluation sets.

500 **Cross-lingual Natural Language Inference (XNLI)** is a cross-lingual sentence understanding  
 501 problem which covers 15 languages, including high-resource languages (English, French, Spanish,  
 502 German, Russian and Chinese), medium-resource languages (Arabic, Turkish, Vietnamese and  
 503 Bulgarian), and low-resource languages (Greek, Thai, Hindi, Swahili and Urdu). The task involves  
 504 determining the relationship between a premise and a hypothesis sentence, and this relationship can  
 505 be categorized into one of three classes: entailment, contradiction, or neutral. We sample 2k instances  
 506 for training and 250 for evaluation sets for each language. NLI serves as an effective benchmark for  
 507 assessing cross-lingual sentence representations, and better approaches for XNLI will lead to better  
 508 general Cross-Lingual Understanding (XLU) techniques.

## 509 **E Implementation**

510 We use Hugging Face’s transformers library [61] and PEFT library [40] for loading pre-trained  
511 models and prompt tuning configurations. For our federated training and evaluation, we use the  
512 Flower framework [9] and PyTorch as the underlying auto-differentiation framework [44]. We use the  
513 AdamW optimizer [37, 25] for all experiments. All experiments are conducted using NVIDIA A40.

## 514 **F Multilingual Distance Measurement**

515 We leverage the database from (author?) [34, 39] to extract feature vectors for each language. These  
516 vectors are then weighted according to the token count of each language in the pre-trained corpus  
517 to calculate the feature vector of the pretrained language. Given the feature vector  $V_i$  for the  $i$ -th  
518 language, token count  $T_i$ , and total tokens  $T_{\text{total}}$ , the weight  $w_i$  is given by  $w_i = \frac{T_i}{T_{\text{total}}}$  and the feature  
519 vector  $V_p$  for the pre-trained model is computed as  $V_p = \sum_{i=1}^n w_i \cdot V_i$ .

520 We define distance for a specific language in the downstream tasks, in terms of the negative logarithm  
521 of its cosine similarity to the pre-trained language. Let  $v$  represent the feature vector of a specific  
522 language in the downstream task. The diversity measure  $\phi$  between this language and the average  
523 language of the pre-trained model is defined as  $\phi(v_i) = -\log(\cos(v_i, V_p))$ .