

1. Introduction

Multilingual pretrained language models (MPLMs), pretrained on multilingual corpora with >100 languages, exhibit strong multilinguality on downstream tasks. Low-resource languages (LRLs), for which little text data is available for pretraining monolingual pretrained language models (PLMs), benefit from MPLMs. However, LRLs are under-represented in the pretraining corpora, resulting in suboptimal performance for these languages. On the other hand, the scarcity of task-specific annotated data makes it difficult for LRLs to employ the pretraining-finetuning paradigm. To mitigate the high demand for labeled data, another line of research – prompt-based learning – emerges, focusing on exploiting large PLMs by reformulating the input. The prompt is designed to help PLMs “understand” the task better and “recall” what has been learned during pretraining. Prompt-based methods provide a new form of zero-shot or few-shot learning in multilingual NLP studies. It involves performing a specific task using prompts, without labeled data in the target language, and is an effective method for LRLs lacking annotated data.

2. Motivation

Our work aims to improve the zero-shot transfer learning performance of LRLs on natural language understanding tasks by taking advantage of cross-lingual information retrieval and the multilinguality of MPLMs. Specifically, we retrieve semantically similar cross-lingual sentences from high-resource languages (HRLs) as prompts and use the cross-lingual retrieval information to benefit the LRLs from the multilinguality of MPLMs. To this end, we made the following contributions in this work:

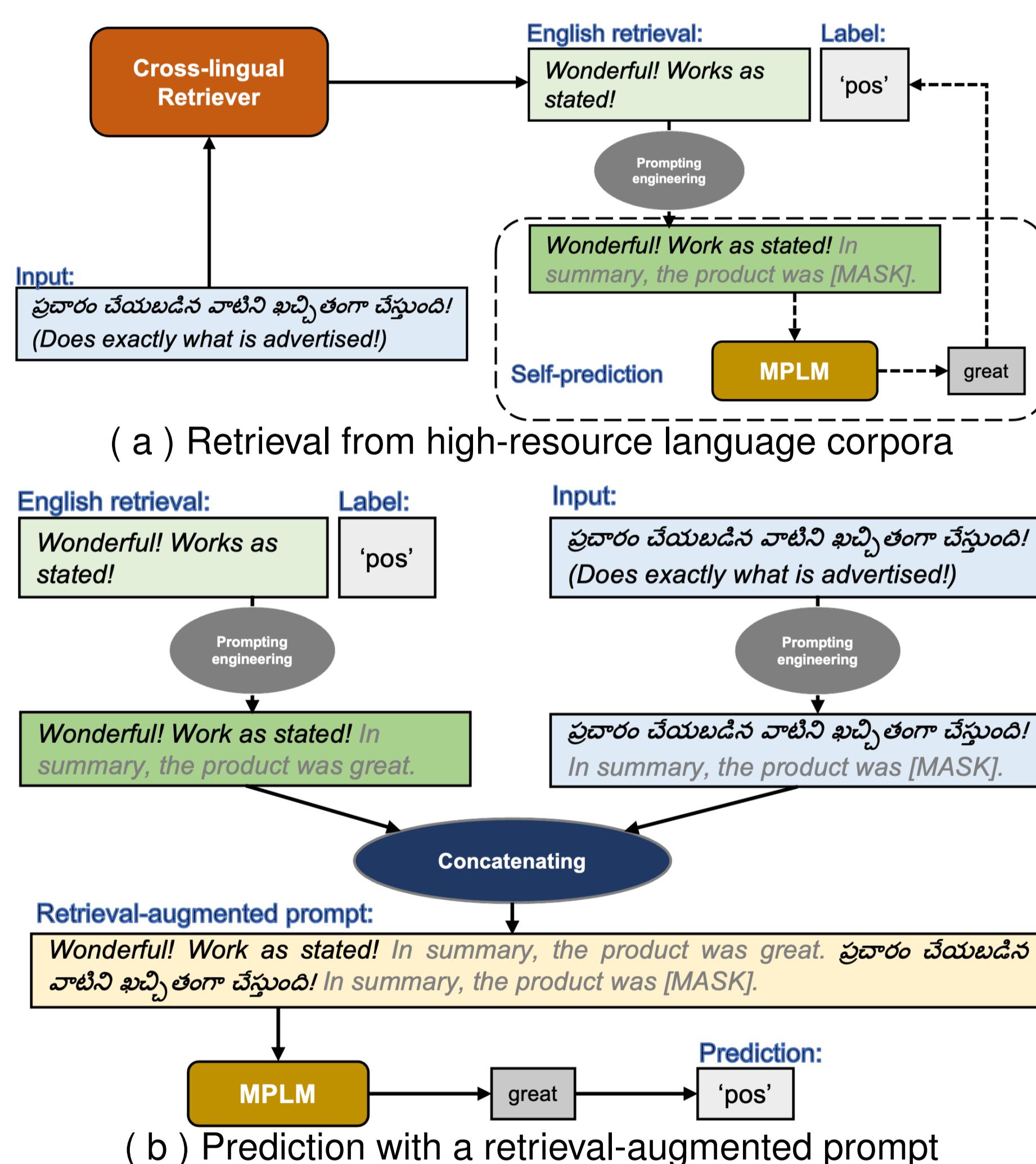


Figure 1: The pipeline of our proposed PARC method

1. We propose **Parc**, a pipeline for integrating retrieved cross-lingual information into the prompts for improved zero-shot learning (Figure 1).
2. We conduct experiments on three different multilingual classification tasks: binary sentiment analysis of product reviews, news topic classification, and natural language inference task. The results show that Parc improves the zero-shot performance on LRLs by retrieving examples from both labeled and unlabeled HRL corpora.
3. To find an optimal configuration of our Parc pipeline, we conduct a comprehensive study on the variables that affect the zero-shot performance: the number of prompts, the choice of HRL, and the robustness w.r.t. other retrieval methods and MPLMs.

3. PARC: Prompts Augmented by Retrieval Crosslingually

We enhance zero-shot learning for low-resource languages (LRLs) by cross-lingual retrieval from **labeled/unlabeled** HRLs. As Figure 1 shows, the PARC pipeline consists of two steps: (a) Cross-lingual retrieval from high-resource language corpora, and (b) prediction with a retrieval-augmented prompt.

3.1 Retrieval from HRL corpora

The cross-lingual retriever takes an LRL input sample as query and retrieves the semantically most similar HRL sample from the HRL corpus. The label of the retrieved HRL sample is obtained either from the corpus (**labeled** setting) or by self-prediction (**unlabeled** setting).

3.2 Prediction with a retrieval-augmented prompt

The retrieved HRL sample, its label and the input sample are transformed into a prompt. For that, we need

- (i) a *pattern P(.)* which converts the input sentence into a cloze-style question with a mask token, e.g.:

$$P(X) = X \circ \text{“In summary, the product was [MASK].”}$$

where X is the input sentence, \circ refers to the string concatenation operator, and

- (ii) a representative word for each possible class. e.g.:

$$\{\text{pos} \rightarrow \text{“great”}, \text{neg} \rightarrow \text{“terrible”}\}$$

- (iii) a mapping (called verbalizer) from the class labels to representative HRL words. e.g.:

$$\{\text{pos} \rightarrow \text{“great”}, \text{neg} \rightarrow \text{“terrible”}\}$$

The prompt pattern $P(.)$ transforms the retrieved HRL sample into the cross-lingual context.

$$C_k^i = P(X_k^{R_i}, v(y_k^{R_i})) \quad (1)$$

where C_k^i is the cross-lingual context for the input X_i^L with the k -th most similar HRL sample $X_k^{R_i}$ and $y_k^{R_i}$ is the label for $X_k^{R_i}$.

Next, the cross-lingual retrieval-augmented prompt I_i is created by concatenating the cross-lingual context and the cloze-style question.

$$I_i = C_k^i \circ P(X_i^L) \quad (2)$$

The augmented prompt I_i is fed to the MPLM M . M performs masked token prediction and returns a probability distribution $p = M(I_i)$ over all HRL words. We predict the class \hat{y} whose verbalizer $v(\hat{y})$ received the highest probability from model M :

$$\hat{y} = \arg \max_{y \in Y} p(v(y)) \quad (3)$$

4. Experimental Results and Analysis

4.1 Main results

| | Amazon | AGNews | XNLI | Avg. |
|----------------|-------------|-------------|-------------|-------------|
| MAJ | 50.0 | 25.0 | 33.3 | 36.1 |
| Random | 48.2 | 25.6 | 32.4 | 35.4 |
| Direct | 53.8 | 36.3 | 33.1 | 41.1 |
| Finetune | 68.6 | 57.9 | 34.5 | 53.7 |
| PARC-unlabeled | 58.4 | 46.7 | 33.5 | 46.2 |
| PARC-labeled | 68.9 | 67.6 | 35.8 | 57.4 |

Table 1: Overview of results on three classification tasks. The reported numbers are averaged across 10 evaluation LRLs. The number of prompts $k=1$ in relevant baselines and our methods for all three tasks.

4.2 Effect of languages

| Unlabeled | Sim. | source size | target size | Labeled | Sim. | source size | target size | |
|-----------|------|-------------|-------------|---------|------|-------------|-------------|-------|
| | corr | p | corr | p | corr | p | corr | p |
| Spearman | 0.28 | 0.05 | 0.20 | 0.16* | 0.31 | 0.03 | 0.42 | 2e-03 |
| Pearson | 0.27 | 0.06 | 0.22 | 0.12* | 0.38 | 6e-03 | 0.41 | 3e-03 |

Table 2: Correlations between Amazon review performance and three features. Sim.: language similarity between an LRL and an HRL; source (target) size: the log of the data size (MB) of source (target). *: insignificant result with a p value larger than 0.05.

| | | | | |
|-----------------------------|------|------|------|------|
| zh(11) det(12) en(14) | 55.3 | 54.9 | 60.1 | 63.4 |
| hi(7) zh(11) det(12) en(14) | 50.9 | 55.1 | 55.7 | 58.5 |
| hi(7) | 54.3 | 54.9 | 60.3 | 62.8 |
| jv(5) | 48.7 | 50.8 | 58.4 | 57.8 |
| uz(6) | 28.0 | 42.6 | 33.7 | 48.9 |
| te(7) | 39.9 | 46.3 | 34.1 | 44.3 |
| ur(7) | 22.9 | 53.1 | 63.3 | 92.2 |

(a) Performance (Unlabeled)

| | | | | |
|-----------------------------|------|------|------|------|
| zh(11) det(12) en(14) | 59.3 | 64.7 | 75.8 | 80.6 |
| hi(7) zh(11) det(12) en(14) | 60.0 | 62.8 | 71.2 | 77.5 |
| hi(7) | 59.8 | 62.5 | 72.4 | 75.1 |
| jv(5) | 59.4 | 64.0 | 72.7 | 80.0 |
| uz(6) | 20.0 | 49.8 | 27.3 | 50.6 |
| te(7) | 33.7 | 44.3 | 42.9 | 51.1 |
| ur(7) | 22.9 | 53.1 | 63.3 | 92.2 |

(b) Language Similarity

| | | | | |
|-----------------------------|------|------|------|------|
| zh(11) det(12) en(14) | 59.3 | 64.7 | 75.8 | 80.6 |
| hi(7) zh(11) det(12) en(14) | 60.0 | 62.8 | 71.2 | 77.5 |
| hi(7) | 59.8 | 62.5 | 72.4 | 75.1 |
| jv(5) | 59.4 | 64.0 | 72.7 | 80.0 |
| uz(6) | 20.0 | 49.8 | 27.3 | 50.6 |
| te(7) | 33.7 | 44.3 | 42.9 | 51.1 |
| ur(7) | 22.9 | 53.1 | 63.3 | 92.2 |

(c) Performance (Labeled)

Figure 2: Visualization of the correlation between zero-shot performance and language similarity, pretraining data size of source and target language.

4.3 Effect of k

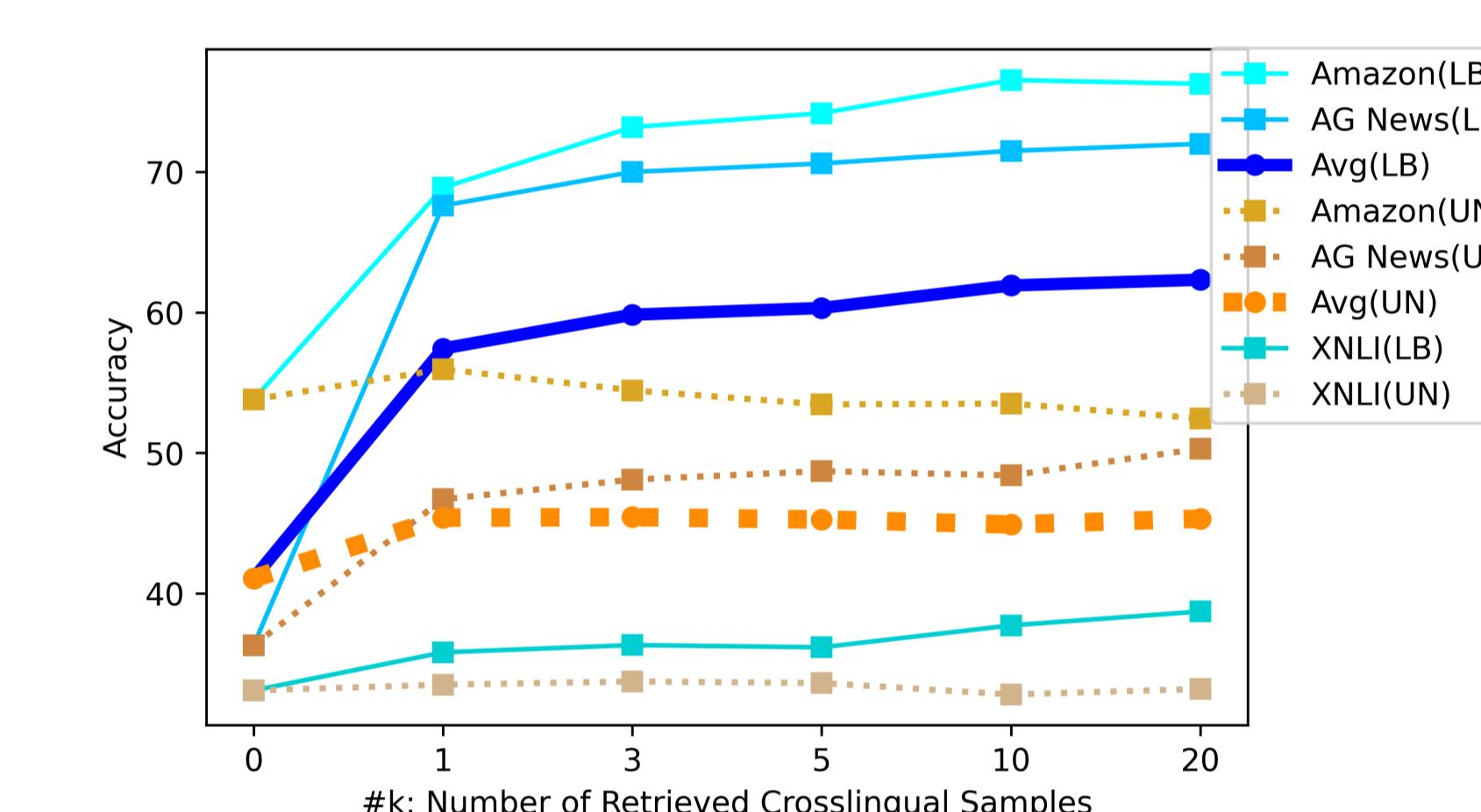


Figure 3: Accuracy on three tasks with different k in the labeled (LB) and unlabeled (UN) setup.

4.4 Generalization to other retrievers and MPLMs

| | Amazon | AGNews | XNLI | Avg. |
|-------------------|--------|--------|------|------|
| Direct | 53.8 | 36.2 | 33.1 | 41.0 |
| mbERT+pooling | 53.1 | 36.9 | 33.6 | 41.2 |
| mbBERT+distiluse | 54.7 | 38.4 | 34.0 | 42.3 |
| mbBERT+paraphrase | 54.6 | 46.7 | 33.7 | 46.7 |
| XLM-R+paraphrase | 70.1 | 57.4 | 34.7 | 54.1 |
| mBERT+LaBSE | 59.4 | 43.8 | 33.1 | 46.1 |
| LB | 53.8 | 58.0 | 33.8 | 48.5 |
| mbBERT+pooling | 53.6 | 58.0 | 32.8 | 48.5 |
| mbBERT+distiluse | 62.8 | 63.8 | 34.6 | 53.7 |
| mbBERT+paraphrase | 72.9 | 67.6 | 36.8 | 59.1 |
| XLM-R+paraphrase | 73.0 | 76.0 | 35.7 | 61.6 |
| mBERT+LaBSE | 72.2 | 80.0 | 37.5 | 63.2 |

Table 3: Accuracy with different models used in our approach. pooling: cosine similarity of the fast hidden states from the MPLM; distiluse: distiluse

