# Clustering Monolingual Vocabularies to Improve Cross-Lingual Generalization

*Riccardo Bassani*

6866207

| *First Supervisor* | *Daily Supervisor* | *Second Examiner* |
| Tejaswini DEOSKAR | Anders SØGAARD | Dong NGUYEN |
| *Utrecht University* | *University of Copenhagen* | *Utrecht University* |

*14 September 2021*

# Abstract

Multi-lingual language models based on multi-lingual BERT (mBERT) have made it possible to perform linguistic tasks on virtually any language for which some raw text data is available. Through a particular type of transfer learning, these models exploit the large amount of data available for resource-rich languages like English and then apply what they learn when handling low-resource languages. This happens both in the pre-training phase, where low-resource languages' data would not be enough for the model to train, and in the fine-tuning phase, where labeled data is often available only for a few languages. Multi-lingual language models, however, exhibit better performance for some languages than for others, and many languages do not seem to benefit from multi-lingual sharing at all. The broad goal of this thesis is to improve the performance of multi-lingual models on low-resource languages. We focus on two issues: First, to create a better multi-lingual representational space where text tokens are represented independently of the language to which they belong. Second, to improve multi-lingual segmentation. In current work, indeed, low-resource languages are segmented with a tokenizer mainly trained on resource-rich languages text, thus yielding sub-optimal segmentation for low-resource languages. Instead, we use a dedicated tokenizer and a large subword vocabulary for better segmentation of each language. We then use a clustering algorithm to discover sensible multi-lingual groupings of segments across languages. A multilingual Bert model is then trained on the obtained clusters. Not only does this allows to use the aforementioned large vocabularies, without increasing the multi-lingual model capacity, but it also yields a truly interlingual model, which learns how to represent language-neutral semantic clusters, instead of language-specific text tokens as in traditional Bert-like models. We called this model ICEBERT, standing for Interlingual-Clusters Enhanced BERT. We show significant improvements over standard multi-lingual segmentation and training in a question answering task covering nine languages, both in a small model regime and in a BERT-base training regime, thus demonstrating the effectiveness of our clustering. The proposed approach could be easily expanded to more languages, or applied to model architectures different from mBERT.

# Contents

# Chapter 1

# Introduction

Pre-trained language models (Devlin et al. 2019, Conneau and Lample 2019, Conneau, Khandelwal, et al. 2020) provide general-purpose sentence representations by assigning to each token in a sentence a contextual vector that encodes information about the token and its relationship with other tokens in the sentence. This is achieved by learning a vector representation - an embedding - for each token in a vocabulary, and the weights of a multiple layers model. Starting from the token embeddings, the model outputs contextual vectors that take into account the full sentence when encoding the meaning of single tokens. Pre-trained language models can be used for several different linguistic tasks such as question answering, sentiment analysis or part-of-speech tagging, by fine-tuning them on task-specific datasets (Dai and Quoc V Le 2015). Pre-training constitutes a crucial phase in the training of a language model, since the quality of the pre-trained model determines its performance on downstream tasks.

In state-of-the-art language models, the vocabulary not only comprises words, but also fractions of words and even single characters, the elements of the vocabulary being therefore referred to as subwords. The process of splitting text into tokens is called tokenization, or segmentation. Various subword tokenization algorithms have been proposed - as WordPiece (Schuster and Nakajima 2012), BPE (Sennrich et al. 2016), and Unigram LM (Kudo 2018) - yielding differences in the token vocabulary which have a large impact on the model performance.

The growing interest in multi-lingual NLP has led to the development of multi-lingual versions of successful models, such as mBERT and XLM (Devlin et al. 2019, Conneau and Lample 2019), where the subword vocabulary is supposed to cover multiple languages (e.g. 104 languages for mBERT). Multi-lingual language models are not only useful in translation-related tasks, but also allow to transfer knowledge from one language to another, drastically improving the availability of NLP tools for low-resource languages (Ruder et al. 2019). Multi-lingual language models, however, exhibit better performance for some languages than for others (Singh et al. 2019). This is mostly due to the imbalance of training data across languages, which leads to an under-representation of certain scripts and languages (Conneau, Khandelwal, et al. 2020). On top of feeding less data to the model, low resource languages also get penalized by being tokenized according to resource-rich languages' statistics (Chung et al. 2020). As a consequence, for some languages, multi-lingual models like mBERT perform worse than their monolingual versions, in which a language-specific tokenizer is used, despite relying on a much larger amount of training data overall (Pyysalo et al. 2020).

This thesis aims to pre-train a multi-lingual BERT model (Devlin et al. 2019) minimizing the negative impact that the low amount of training data has on the model's performance on low-resource languages. The method proposed in this study has two main goals:

- Increase the degree of multi-linguality of language models by making their vector space a true interlingua.

- Increase low-resource languages performance by adopting a tokenization method which does not penalize them.

We try to achieve these goals firstly by using a dedicated tokenizer - and therefore a large dedicated subword vocabulary - for each language, and secondly by semantically clustering subwords across all the languages' vocabularies. By fixing the number of clusters and training a model on cluster-IDs, it is possible to optimize tokenization without resorting to a large subword vocabulary, therefore controlling for the model capacity. Moreover, the model is not trained on language-specific tokens, but on language-neutral cluster IDs representing a concept which includes multiple subwords from different languages. The core of this study is therefore the identification of an effective clustering process, followed by the training and the evaluation of a multi-lingual BERT language model. We called the model trained following this method ICEBERT, for Interlingual-Clusters Enhanced BERT. In this work, we focused on the 9 languages included in the Gold Passage task of the TyDiQA datasets (Clark et al. 2020a): Arabic, Bengali, English, Finnish, Indonesian, Korean, Russian, Swahili, Telugu. These languages were chosen by the dataset creators to be typologically diverse, i.e. to express different linguistic features. Using such a set of languages, we can reasonably expect that a model performing well on this set will generalize across a large number of the languages in the world.

We evaluate the model in a zero-shot scenario: The model, pre-trained on data from all languages, is fine-tuned only on English labeled data, and is then evaluated on test data in other languages, for which no data was seen during fine-tuning. We chose this evaluation method since this is the most realistic situation for the majority of low resource languages, who lack labeled data for most downstream tasks. We chose the GoldP-TyDiQA task (Clark et al. 2020b), a question answering task, due to the high typological diversity of the language included. We expected to see improvements across all the 9 languages apart from English (whose performance were expected to drop due to the replacement of actual English words with cluster IDs in the fine-tuning data, and to the scarce benefits deriving from training on other languages' data), but our interest was mainly focused on low-resource languages. This means we aimed to improve the lowest score, and that we would have considered a success, for instance, an increase in performance for Swahili and Telugu, even if it had occurred along with a drop in performance for Russian.

When comparing our ICEBERT model to a baseline using standard multi-lingual segmentation and training, we improved F1 scores for most of the languages in the experiment, as well as on average. The most relevant improvements are the ones obtained on Telugu, the language with lowest baseline score among the ones included in the experiment. Our approach significantly improved performance on Telugu both according to the Exact Match metric (7.12 *vs.* 4.24) and to the F1 metric (14.03 *vs.* 10.38).

The rest of the thesis is structured as follows: Chapter 2 gives the reader a theoretical background, and presents some key works motivating the thesis. Sections 2.1, 2.2, 2.3 and 2.4 respectively introduce word embeddings, neural word embeddings, pre-trained word embeddings and cross-lingual word embeddings.

Section 2.5 deals with text tokenization in multi-lingual language models, and section 2.6 underlines the need for a novel approach to multi-lingual text tokenization, outlining the motivations and the goals of the thesis.

Chapter 3 describes the creation of subword clusters, starting with the building of a joint multi-lingual vocabulary (section 3.2) and then moving to the clustering itself (section 3.3).

Chapter 4 discusses the training of a Bert model on the learnt clusters. The chapter includes a description of the developed baseline model, the training data, and two Bert models, in a small-scale (section 4.5) and in a large-scale (section 4.6) scenario.

Chapter 5 reports and discusses ICEBERT's performance on the GoldP-TyDiQA task, while chapter 6 identifies the main questions raised by our results, and the possible next steps to take to build on the work presented in the thesis.

Finally Chapter 7 briefly summarizes the work done and the achieved results.

The code developed during this thesis is available at https://github.com/BassaniRiccardo/transformers/tree/master/examples/pytorch/language-modeling.

# Chapter 2

# Theoretical Background and Related Work

This chapter provides some theoretical background and presents works that are relevant to this thesis. It starts by introducing the idea of word vectors (2.1), and it then moves to describe neural word embeddings (section 2.2) and pre-trained word embeddings (2.3). Section 2.4 introduces cross-lingual embeddings, giving an overview of different methods of cross-lingual embeddings creation: word-level approaches (2.4.1), sentence-level approaches (2.4.2), document-level approaches (2.4.3), and neural approaches (2.4.4). The focus of the chapter then moves to tokenization in multi-lingual models (2.5), since the main goal of this thesis is optimizing tokenization in such models. The possible levels of tokenization are outlined in section 2.5.1. Four popular subword tokenization algorithms are described: Byte Pair Encoding (section 2.5.2), WordPiece (section 2.5.3), Unigram Language Modeling (section 2.5.4), and SentencePiece (section 2.5.5). Finally, section 2.6 points out the need for a fairer segmentation approach, motivating the thesis and defining its goals.

## 2.1 Distributional Semantics: the Introduction of Word Vectors

At the dawn of Natural Language Processing, words were represented with simple strings or IDs to lookup in a dictionary, following the so-called symbolic approach. In this way, no information about the meaning of words were encoded in their representations, which could not be exploited to derive relationships (e.g. similarity) between words. The inspiration about how to create meaningful representations for words came from the field of Distributional Semantics, particularly by the idea that "You shall know a word by the company it keeps" (Firth 1957): The context in which words appear can be used a starting point to build vector representations of the words. In the late 1980s, also thanks to the increase in computational power, a shift of paradigm took place, and data-driven, statistical methods started to become dominant: Researchers begun to feed machines with very large text corpora, in order to extract statistical information and to use this information to create vector representations of words.

Early works followed a so-called count-based approach: the Latent Semantic Analysis (LSA) model (Deerwester et al. 1990), for instance, creates a word-documents co-occurrence matrix from which it extracts vector representations for words and documents. Other approaches, as distributional semantic models, use other words as contexts instead of documents, but still rely on large co-occurrence matrices,

which are usually reduced through techniques such as Single Value Decomposition (SVD). Some count-based models, such as Glove (2014), managed to remain competitive until today, but, in general, the focus in research is currently on a different type of methods to learn word representations: the so-called Neural Word Embeddings.

## 2.2   Neural NLP and Neural Word Embeddings

Starting from the 1990s, using vectors to represent words became the de-facto standard in NLP, and the differences between distinct approaches in word representations were given by how these vectors were created. The already cited statistical methods, indeed, were not the only option available: Several events in the 1980s, such as the rediscovery of backpropagation (Rumelhart et al. 1986) and the introduction of recurrent neural networks (Hopfield 1982), not only directly provided tools that would have turned out to be extremely useful to NLP, but also increased the general interest in neural networks, further boosting progress in the field (Macukow 2016).

One of the first breakthrough application of neural networks for NLP is considered to be the work described in the paper "A Neural Probabilistic Language Model" (Bengio et al. 2000), where the authors adopt a neural approach in a Language Modeling context. Language Modeling (LM), the task of assigning probabilities to sequences of words, already had a 20 years long history at the time. However, language models up to 2000 typically exploited statistical information using "very little knowledge of what language really is", as reported in Rosenfeld (2000). The model proposed by Banjo et al., on the other hand, used vector representations of words as input of a feed-forward neural network, and trained them jointly with the model parameters. The neural approach, therefore, provided a method to create word vectors without extracting explicit statistical information from a corpus, but training them on a particular task, in Banjo's case LM.

Following Banjo et al.'s article, several models were proposed with different architectures and training objectives, and the word vectors learned via this approach started to be referred to with a new term: "word embeddings". Despite this term having often been credited to Banjo et al., the latter simply referred to their methods as "learning a distributed representation for words", and this is what "word embedding" actually amounts to. Significant improvements in LM (and general) performance were obtained through the use of more adequate architectures, such as Recurrent Neural Networks (RNNs; Mikolov, Karafiát, et al. 2010) and, later, long short-term memory networks (LSTMs; Graves 2013). These types of networks were able to capture and reflect the temporal, sequential network of natural language, and led to better word representations. Moreover, the introduction of easier and more efficient methods to train word embeddings (Mikolov, Quoc V. Le, et al. 2013, Pennington et al. 2014) contributed to make their use even more popular. Current state-of-the-art models almost universally adopt the Transformer architecture (Vaswani et al. 2017), based on the enhancement of sequence-to-sequence models (Sutskever et al. 2014) with an attention mechanism (Bahdanau et al. 2015).

## 2.3   Pre-trained Word Embeddings

On top of improved model architectures, a key idea concerning word representation is to credit for the incredible success that many NLP models have obtained in recent years: pre-trained models (and word embeddings). Works like the one of Dai and Quoc V Le (2015) explored the idea of training a model on an unsupervised task (e.g. LM), and to then use the so-trained model (and embeddings) as a starting point, fine-tuning it on a specific task requiring supervised data. The good results achieved in these

studies showed how such a semi-supervised approach allows models to generalize better and makes the training more stable.

Fine-tuning is one of two types of transfer learning, i.e. a technique aiming at using knowledge gained from solving one problem to solve a different but related one (Bozinovski and Fulgosi 1976). The second approach, adopted by some successive works, goes under the name of feature extraction. Feature extraction exploits pre-trained embeddings differently, creating task-specific embeddings by learning functions of the internal states of a pre-trained model (ELMo; Peters et al. 2018). However, fine-tuning remained the most used transfer learning technique, notably being adopted by BERT (Devlin et al. 2019), arguably the most breakthrough model of recent years. A Transformer-based model, BERT introduced a new learning objective, Masked Language Modeling (MLM), consisting in predicting randomly masked words in a sentence. This inherently bidirectional task allowed the model to overcome LM's limitations and led to much better performance. Due to its astounding achievements, in the last two years BERT has often been regarded as the starting point to build on for the creation of improved models. A key difference between BERT-like neural approaches and not neural ones is that the former not only provides static embeddings, but also a multiple layer network in which each layer represents the words in an increasingly more contextualized fashion. The input layer simply assigns a learned static embedding to each word, while the last one offers a contextualized representation of each word, i.e. each word is represented taking into account also the other words in the sentence.

## 2.4    Cross-Lingual Embeddings

The English language has had a special place among other languages in the NLP field. Besides English being the language internationally adopted in research, another crucial factor for its predominance in NLP is the large amount of data available in the language in comparison to other languages, especially for what concerns annotated data. The development of less computationally expensive methods has made it possible to build models for languages different from English (e.g. French, German, Dutch), but low-resource languages remain strongly penalized. On top of that, monolingual word embeddings belong to language-specific vector spaces and are not comparable across languages.

Recently, the research community has become increasingly aware of the necessity of overcoming this limitation, building word embeddings that inhabit a common multi-lingual space: Such embeddings are referred to as "cross-lingual embeddings". Efforts towards the creation of high-quality cross-lingual word embeddings are motivated by two reasons (Ruder et al. 2019). First, being able to capture word similarity across languages, they are crucial in cross-lingual NLP applications as machine translation or bilingual/multi-lingual lexicon induction. Second, they constitute the starting point to narrow the gap in NLP tools availability across languages, by allowing cross-lingual transfer from resource-rich to low-resource languages (Tela et al. 2020).

The main differences between the various proposed methods to create cross-lingual embeddings lie in the nature of the data they use. Early works in the creation of cross-lingual embeddings typically relied on a significant amount of parallel data (Klementiev et al. 2012). Three types of parallel data can be identified, depending on the type of alignment: this can be at word level (2.4.1), at sentence level (2.4.2), or document level (2.4.3).

### 2.4.1   Word-Level Approaches

Word-level approaches are arguably the most popular, and, among them, mapping-based approaches are the most common ones since the introduction of Mikolov et al.'s linear transformation method (2013). Mapping-based methods first build language-specific embedding spaces, and then learn a mapping from an embedding space to the other, in the form of a transformation matrix. In particular, Mikolov et al. proposed to learn a linear projection maximizing the similarity between the embedding space of the target language and the embedding space of the mapping of the source language's space obtained with such a projection. This approach has later been optimized by constraining the transformation matrix to be orthogonal (Xing et al. 2015). Methods imposing such a constraint preserve length normalization and have also been shown to ensure monolingual invariance (Artetxe et al. 2016). Other mapping-based approaches use Canonical Correlation Analysis (CCA) to map the embedding space of both the source and target language to a single shared space, by maximizing the correlation between the projections of translation pairs (Faruqui and Dyer 2014). Finally, margin methods learn the mapping in such a way that the similarity between translations is high, while the similarity between random word pairs is low (Lazaridou et al. 2015).

Despite most of the cited studies being concerned with only bilingual embeddings, all these methods can be extended to a multi-lingual setting, an example being the work of Ammar et al. (2016) extending the CCA approach to multiple languages. However, requiring a large seed lexicon is often unrealistic for low-resource languages: In order to make the approach more suitable to a largely multi-lingual context, some studies proposed to use easy-to-obtain weak supervision signals (Smith et al. 2017, Artetxe et al. 2017, Søgaard et al. 2018), or even not to use any supervision at all (Conneau, Lample, et al. 2017, Artetxe et al. 2018, Hoshen and L. Wolf 2018, Alvarez-Melis and Jaakkola 2018).

Two other classes of methods exploiting word-aligned data must be mentioned: The first one consists of methods based on pseudo-bilingual corpora (Xiao and Guo 2014), i.e. automatically constructed corpora where words from both the source and the target language are present. The second one is made of methods that exploit a cross-lingual regularization term to jointly minimize the monolingual losses of the source and the target languages (Klementiev et al. 2012). All these three approaches can be extended to a multi-lingual context and are comparable in performance. The mapping-based approach is the most popular one, being easy to use and computationally less expensive than the alternatives (Ruder et al. 2019).

### 2.4.2   Sentence-level Approaches

Despite word-level supervision being easy to obtain, the need for cross-lingual sentence representations has recently started to motivate the study and the usage of sentence-level supervision. Sentence level approaches exploit sentence-aligned parallel corpora as the Bible corpus (Christodoulopoulos and Steedman 2015). Levy et al. (2017) observed that most methods use parallel corpora in one of two ways. The first class of methods represents each word using all the words appearing in the same sentences in which it appears or in sentences aligned to the ones in which it appears, counting the co-occurrences frequency. The second class of sentence-aligned methods represents each word by the set of sentences in which it appears, regardless of its frequency in each one. Regardless of the used feature set, sentence-level approaches are appealing since they provide truly inter-lingual representation, and, differently from word-level approaches, they do not rely on the assumption that the monolingual spaces are isomorphic, an

assumption which results to be problematic for languages typologically distant from each other (Søgaard et al. 2018).

### 2.4.3    Document-level Approaches

Sentence-aligned parallel corpora are not easy to obtain: Similar approaches with easier requirements exploit document-aligned data, where the documents do not have to be parallel, but only comparable. It is the case of INVERTED (Søgaard et al. 2015), which uses inverted indexes (document-IDs) to represent words by the Wikipedia concepts they are used to describe. This method was actually proposed before Levy, Søgaard and Goldberg's work and, at the time of publication, it outperformed state-of-the-art approaches to neural net word embeddings (Klementiev et al. 2012, Chandar et al. 2014). Even though document-aligned data has been shown to be less informative than sentence-aligned data (Levy et al. 2017), it is an easier requirement, and methods leveraging this kind of data can therefore be extended more easily to truly low-resource languages. Comparable data aligned at word or sentence level can also be used, but parallel data are preferred, especially in word-level approaches where comparable data yields significantly worse results (Ruder et al. 2019).

### 2.4.4    Neural Approaches

The success of Transformer-based LMs in the monolingual scenario was shortly followed by the extensions of such models to the multi-lingual context, an iconic case being the multi-lingual version of BERT (mBERT). Large multi-lingual language models such as mBERT revolutionized the approach to the creation of word embeddings: training a model on a large corpus including data from different languages yields a multi-lingual model and automatically creates cross-lingual embeddings. Differently from the previously discussed methods, these models are able to learn contextualized embeddings. Moreover, even though improvements leveraging sentence-aligned data have been proposed (Conneau and Lample 2019), in principle they do not need any parallel data. However, the computational effort required to perform the training process is a significant issue, not only for what concerns the training time, but also from a financial and an environmental perspective (Strubell et al. 2019). This makes it prohibitive for most researchers to train a full-scale model, and the standard strategy adopted when developing new methods consists in focusing on the far less expensive fine-tuning side or, if necessary, experimenting with small-size models to assess the validity of the proposed improvements.

Since its publication in 2018, mBERT has been the object of several studies and proven to have some degree of multi-linguality. Pires et al. (2019) showed its good performance at zero-shot cross-lingual transfer: by fine-tuning on a specific task using training data in a certain language and evaluating on the same task in another language, they assessed the degree to which the model generalizes across languages. Despite results being encouraging, transfer performance resulted to be significantly better for similar languages, and the analysis mainly focused on syntactic tasks. More importantly, Pires et al.'s study showed how there is a language-neutral component in mBERT, but did not prove the model made no use of language-specific information.

Following this work, the question of how language-neutral mBERT is was addressed by Libovický et al. (2019). They found that mBERT representations are not really language-neutral, but are instead

constituted of a language-neutral and a language-specific component. While the former is enough for easy semantic tasks like sentence retrieval or word-alignment, the latter is necessary to perform harder semantic tasks such as Machine Translation quality estimation. The language-specific component, however, was shown to mainly consist in a constant, language-specific shift in the embedding space. Removing this shift resulted in a drastic decrease in language identification performance (confirming it encoded language-specific information), and in a significant improvement in parallel sentence retrieval performance. By comparing the language-specific shift vectors, Libovický et al. (2019) also concluded that mBERT clusters languages by their families.

Similar results were achieved by Singh et al. (2019), who showed that mBERT's vector space is not really an interlingua - i.e. a common representational space for semantically similar text across languages - by analyzing mBERT's representation of similar sentences across languages. Their results revealed how these representations diverge in Projection Weighted Canonical Correlation Analysis (PWCCA) similarity as moving deeper through mBERT's layers (a similar trend was also observed by Pires et al. (2019)). In particular, their study exposes how the model does not really exploit a single multi-lingual space, but partitions the space to better reflect the linguistic and evolutionary relationships between languages. These findings are coherent with the ones of Libovický et al. (2019), and underline the need for new techniques yielding more interlingual models. Singh et al. (2019) also showed that the phenomena they observed are stronger when using subword tokenization compared to character or word level tokenization, highlighting the issues of the subword approach and the importance of tokenization in general. The next section will introduce tokenization, discussing the various methods already proposed and their limitations in a multi-lingual context.

## 2.5   Tokenization Approaches in Multi-Lingual Models

Text tokenization, or text segmentation, is the task of splitting text in "tokens", that is to say in meaningful units (Berasategi 2020). Tokenization therefore determines the vocabulary used by the model, hence the word types that the model can recognize, and has a strong impact on the model performance. This section introduces tokenization (2.5.1), describes the most common tokenization algorithms (sections 2.5.2, 2.5.3, 2.5.4, 2.5.5), and then underlines the issues with current approaches to multi-lingual tokenization. Finally, section 2.6 motivates the work done in this thesis and it outlines its main goals.

### 2.5.1   Word, Character, and Subword Tokenization

The naive approach to text tokenization, consisting of identifying each word with a token, is called word-level tokenization. The main issue with this approach is the large vocabulary size it creates, a problem that is even more relevant in the multi-lingual context, where the shared vocabulary would easily contain more than one million word types. Such a large vocabulary would prohibitively increase the model capacity, making the training too expensive. A related issue would arise with out-of-vocabulary words: When a new word is encountered, a model relying on word representation would fail to assign a representation to the new word. At the other end of the spectrum there is the possibility of using single characters as tokens (Karpathy et al. 2015). This drastically reduces the vocabulary size down to a few hundred elements (including special characters on top of letters). Character-level tokenization has been shown to be surprisingly effective in various NLP tasks such as sentiment analysis (Radford, Józefowicz, et al. 2017) and machine translation (Kalchbrenner et al. 2016, Lee et al. 2017). However,

using characters as tokens, while the vocabulary size decreases, the sequence length (i.e. the number of tokens necessary to represent a sequence) increases significantly, leading to higher training times for equal amounts of information.

The interim solution most commonly adopted is to use fractions of words as tokens: the so-called subword-level tokenization. This is the approach used by the majority of modern language models, mBERT included. It allows to reduce the vocabulary size and to deal with out of vocabulary words by including characters in the vocabulary and by resorting to character-level tokenization when larger subwords are not matched. At the same time, sequence length decreases significantly in comparison to character-level tokenization. Differently from the word and the character case, however, there is not a unique way to split text into subwords: various algorithms have been proposed, of which the four more popular will be introduced in the following sections.

### 2.5.2   Byte Pair Encoding (BPE)

BPE was first introduced in 1994 as "a technique for data compression that works by replacing common pairs of consecutive bytes with a byte that does not appear in that data." (Gage 1994). More than 20 years later, it was adapted to text tokenization: The vocabulary is initialized with the character vocabulary and the most frequently occurring subword pairs are merged together (instead of being replaced by another byte/symbol as in original BPE) until the desired vocabulary is size is reached (Sennrich et al. 2016). A famous and successful application of BPE can be found in the OpenAI GPT-2 model (Radford, J. Wu, et al. 2019).

### 2.5.3   WordPiece

WordPiece is a subword tokenization algorithm introduced by Schuster and Nakajima (2012) when solving Japanese and Korea voice problem. It is very similar to BPE, the main difference being that new subwords are not added by frequency, but depending on the likelihood on the training data which is achieved when the new token is added to the model. Y. Wu et al. (2016) slightly modified the algorithm a few years later, and this version was exploited by BERT (and consequently by mBERT). It is important to stress the distinction made between subwords that are or are not at the beginning of a word: BERT tokens contain a '##' suffix when they are not at the beginning of a word.

### 2.5.4   Unigram Language Modeling (Unigram LM)

The Unigram LM algorithm, proposed by Kudo (2018), takes a different approach to subword creation: instead of starting from a small vocabulary, they start from a very large one (e.g. obtained through BPE with a large number of merge operations) and then it excludes the least useful subwords until the desired size is reached. To quantify the utility of each subword, the algorithm exploits a unigram language model and a training corpus, and it defines the loss of each subword as the reduction of the likelihood on the corpus when the subword is removed from the vocabulary. At each iteration, the top 80% of subwords with the highest loss are kept, while the remaining 20% are discarded. This process is repeated until the

desired vocabulary size is reached.

### 2.5.5    SentencePiece

The three algorithms described above cannot be applied to raw text, but they require the text to be pre-processed using a toolkit such as Moses[1] for word alignment. Therefore, these algorithms do not provide an end-to-end system able to deal with any type of text, but they rely on language-dependent pre-processors, which can be problematic when dealing with a large number of different languages, especially with non-segmented languages (e.g. Chinese). Motivated by the necessity of a language-independent standard for text tokenization, Kudo and Richardson (2018) developed SentencePiece, a purely end-to-end system available as a C++, Python, and Tensorflow library. SentencePiece implements both BPE and Unigram LM, and it offers the option to customize character normalization. A great example of its effectiveness in the multi-lingual context is its exploitation in the work of Pyysalo et al. (2020), where it was used to create monolingual BERT models for 42 languages.

## 2.6    The Need for a Fairer Segmentation Approach

The limitations of current tokenization approaches that will be discussed in the next paragraph derive mainly by the fact that models like mBERT use a single subword vocabulary for all languages (they exploit subword sharing). When handling subwords belonging to different languages, indeed, mBERT does not use any marker to denote the input language, the motivation being that this makes zero-shot training work[2]. If no markers are used, indeed, any text in a script already known by the model can be segmented into known subwords. If subwords were marked, instead, the tokenizer would not know which subwords to use when encountering a new language, not having any subword marked with the unknown language's marker in the vocabulary. Another reason to support the omission of language-specific markers is that in this way subwords present in multiple languages can be exploited to make the model more multi-lingual. This hypothesis - that word piece overlap is to credit for mBERT performance - has been tested in several works. Pires et al. (2019) showed that zero-shot NER performance does depend on word piece overlap, but the latter is not the only reason for mBERT good performance, showing that the model is able to create representations that go beyond the exploitation of vocabulary memorization. In a subsequent study, S. Wu and Dredze (2019) found a strong correlation between word piece overlap and the model performance on zero-shot MLDoc, NER, POS, and Dependency Parsing. A third study, however, gave contrasting results, reporting word piece overlap not to be significant and therefore reducing the importance attributable to shared subwords (K et al. 2020).

Possible downsides of subword sharing were stressed by Chung et al. (2020) in a recent paper. The authors claimed that the joint approach taken in the identification of the best overall subword vocabulary overemphasizes the importance of cross-lingual subword sharing. According to their view, this especially penalizes low resource languages, whose words end up being segmented according to statistics of resource-rich languages, and therefore in suboptimal subwords that do not have a meaning in the low-resource language they belong to. Moreover, the subword vocabulary creation is biased towards common scripts, leading to poor performance for languages with rare scripts. Suboptimal subwords tokenization could also explain the results obtained by Pyysalo et al. (2020), which showed how training monolingual

---

[1]http://www.statmt.org/moses/
[2]https://github.com/google-research/bert/blob/master/multi-lingual.md

BERT models can improve performance even if the amount of training data decreases. Finally, Hu et al. (2020), when evaluating various multi-lingual language models such as mBERT, XLM, and XLM-R in their XTREME paper, observe lower performance on non-Indo-European languages. They identify in the under-representation of ideograms in the joint SentencePiece vocabulary a possible cause of this phenomenon, this thesis being supported by the findings of Conneau, Khandelwal, et al. (2020) on the importance of vocabulary size in a cross-lingual model's performance.

The described shortcomings of current approaches to multi-lingual subword tokenization motivate the search for a better and fairer segmentation method, which takes into account the necessity for all languages, also low-resource ones, to be tokenized according to a specific subword vocabulary, instead of relying on a general one that favours languages that are most represented in the training corpus. An important challenge to address when developing such a method is to avoid an excessive increase of the model capacity: in multi-lingual models, indeed, the embedding matrix (the matrix whose rows are embeddings of the word types in the vocabulary) accounts for a large portion of the model parameters (e.g. 47% in XLM-R, Conneau, Khandelwal, et al. 2020).

The method proposed in this study, therefore, has two main goals:

- Increase the degree of multi-linguality of language models by making their vector space a true interlingua.

- Increase low-resource languages performance by adopting a tokenization method which does not penalize them.

# Chapter 3

# Subword Clusters Creation

## 3.1 Project Overview

In order to overcome the limitations of current segmentation methods, we introduce a novel approach to subword segmentation. The method also allows training a model which represents text from different languages in a common representational space (i.e. an interlingua). The key idea was to build a single multi-lingual model by using distinct monolingual tokenizers - therefore language specific subword vocabularies - and then clustering the obtained subwords by semantic similarity. A single model was then trained on the desired number of clusters, learning an embedding for each cluster. This allowed us to use of a very large vocabulary without increasing the model capacity, leading to a better segmentation for each language.

The project can be seen as articulated in three main phases:

1. Creation of a large, joint, multi-lingual subword vocabulary (section 3.2).

2. Implementation of a subword clustering algorithm for reducing the vocabulary size and creating an interlingual representational space (section 3.3).

3. Model training and evaluation (section 4).

An overview of the project is given in figure 3.1, which outlines the three phases listed above: join vocabulary creation (in green), subword clustering (in blue), and model training (in fucsia).

## 3.2 Joint Vocabulary Creation

The joint vocabulary was build by merging monolingual vocabularies. The monolingual vocabularies were created starting from the BPEmb vocabularies offered by Heinzerling and Strube (2018). BPEmb[1] is a collection of BPE-based subword[2] embeddings in 275 languages. Monolingual vocabularies are available in different sizes. We extracted vocabularies of size 30k[3] for the 9 languages present in the TyDiQA-GoldP task (Clark et al. 2020a), a question answering task included in the XTREME benchmark (Hu et

---

[1] //bpemb.h-its.org/

[2] From now on, with the term subwords and tokens we will refer to tokens in general, regardless of the them being full words or portion of words. When talking about non-full words specifically, we will use the term 'strict subwords'.

[3] This vocabulary size was chosen to be similar to the size of the monolingual BERT's vocabulary.
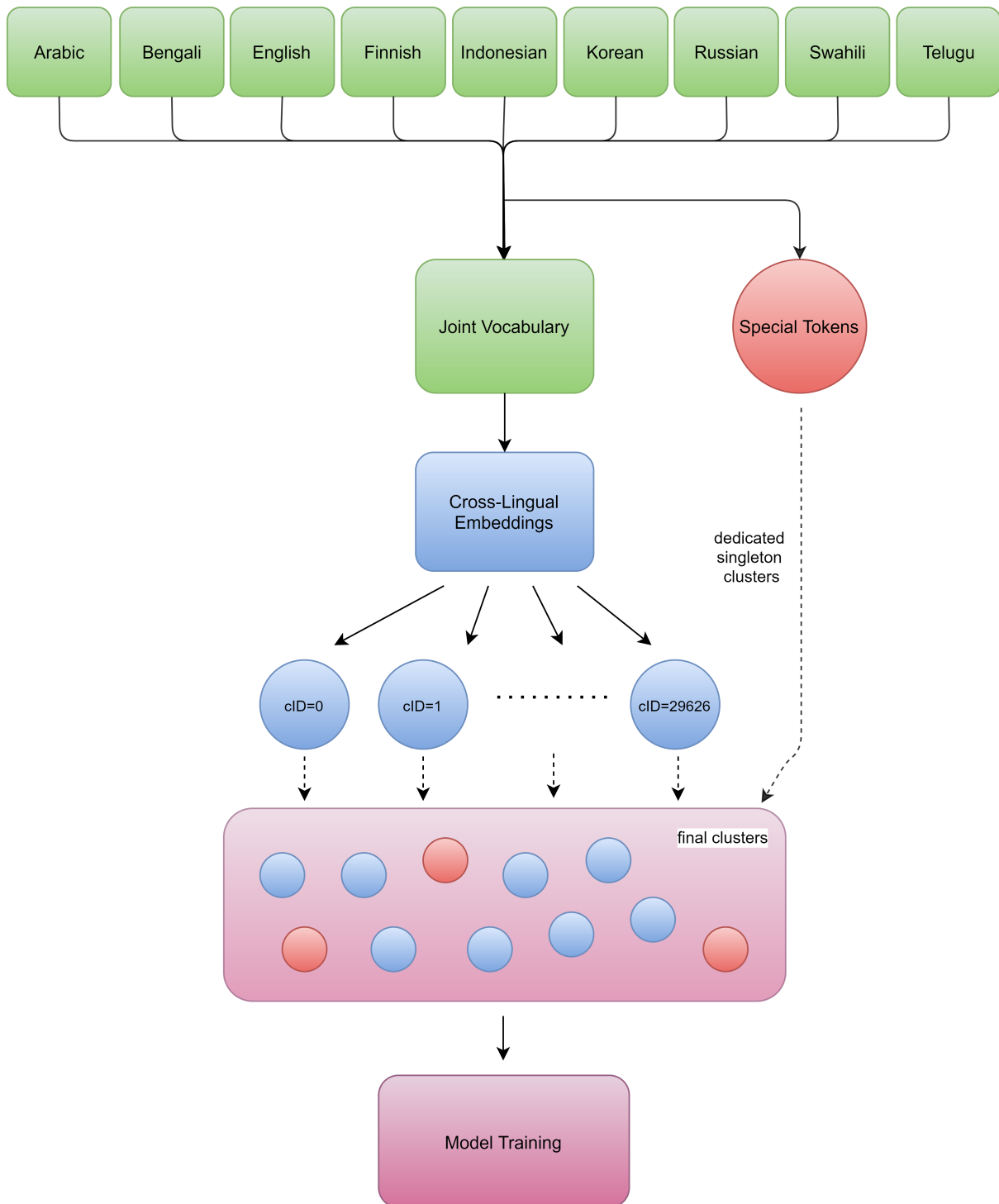
Figure 3.1: Overview of the three phases of the project. From top to bottom: joint vocabulary creation, subword clustering, and model training. The clusters are identified by clusters IDs (cIDs), and are represented as circles in the figure. Blue circles represent clusters obtained by partitioning the translation graph, while red circles represent singleton clusters dedicated to special tokens.

al. 2020). The languages considered in this study are thus Arabic, Bengali, English, Finnish, Indonesian, Korean, Russian, Swahili, and Telugu. These languages were chosen for two reasons. First, they are typologically very different and constitute an adequate challenge for a multi-lingual model. Second, they allows us to evaluate on a downstream task (in general, most low-resource languages lack test data for

many downstream tasks), and therefore to assess the effectiveness of our approach.

The vocabularies were transformed from a SentencePiece format to a WordPiece format: the '_' marker was removed from subwords at the beginning of a word, while a '#' marker was added to subwords starting in the middle of a word. Following the approach of Dufter and Schütze (2020), we mapped all digits to zeros. For each language, two different vocabularies were created:

1. A vocabulary containing only the subwords to cluster. Digits, punctuation and other special symbols which could be considered language-agnostic were excluded from the clustering process. Each one of these special tokens was assigned to a dedicated, singleton cluster. These vocabularies were saved together with the relative lists of embeddings, extracted from the models provided by Heinzerling and Strube (2018). The embeddings were then mapped into a shared space by projecting all non-English vectors into the English vector space (section 3.3.1), and used in the clustering process.

2. A complete vocabulary, containing not only the subwords to cluster, but also special symbols and digits. These vocabularies were necessary to create monolingual tokenizers.

We will refer to the first set of vocabularies when discussing the clustering process. The vocabularies were merged into a unique vocabulary (the "joint vocabulary"). During this process, the tokens were marked with a code ("_ln") representing the language they belonged to, to avoid identical subwords with different meanings across languages being treated as the same token.

### 3.2.1   Subword Marking

The choice to mark tokens was made due to Chung et al.'s observation (2020) that some languages risk being tokenized according to other languages' statistics. This excludes the possibility of doing zero-shot transfer toward languages not included in the training process, which is instead possible in mBERT thanks to the authors' choice not to use markers. However, this decision guarantees that the meaning of language-specific subwords is taken into account and it is crucial to guarantee a better segmentation. Moreover, an important form of zero-shot transfer remains possible: for a specific task, the model can be trained on English (for which annotated data is available) and then used to process any other language included in the clustering process. Indeed, being the model fine-tuned on cluster-IDs, the fine-tuning should be effective not only for English but also for the other languages. Finally, zero-shot transfer towards unseen languages would still be feasible: a possible solution could consist in exploiting a parallel corpus like JW300 (Agić and Vulić 2019) - containing bilingual data for most language pairs - to learn subword mappings from the new language to each of the already supported languages. For each word in the new language, the cluster-ID[4] majority label across the cluster-IDs of the target subwords in the other languages could then be taken. However, this idea was not explored in this study. It is also worth to note that the same approach adopted in this work can be extended, with few simple modification in the mapping process, to up to all the 275 languages for which BPEmb provides tokenizers and embeddings.

At the end of the vocabulary creation process, the joint vocabulary comprised 269649 tokens, and the number of special tokens excluded from the clustering process amounted to 665[5].

---

[4]The clustering process will be described in section 3.3.

[5]The sum of these two numbers is higher than 9 x 30k = 270k, since some special tokens were included both in hashed

## 3.3    Subword Clustering

This section discusses the clustering process. The goal of this process is to obtain clusters made of tokens which are semantically similar, grouping together tokens from different language regardless of their language. The ideal clusters will therefore contain tokens which are translations or synonyms of each other. Moreover, we want clusters to be generally small and of similar size: if all languages' representational spaces were isomorphic, each abstract token would have a concrete instance in each language, and we would like to group these instances together, obtaining clusters of size equal to the number of languages. Since vector spaces are not really isomorphic, we want to be more elastic about clusters sizes, and to allow for larger and smaller clusters. However, very large clusters are undesirable, since they risk to group many tokens which are not actually semantically close. Therefore, while allowing some variance in cluster size, we still want to limit the maximum cluster dimension.

In order to cluster tokens from different languages, it is necessary that the tokens' embeddings inhabit the same representational space, so that semantically similar tokens are represented by similar vectors independently of their language. The section therefore starts by describing the steps followed in the creation of cross-lingual embeddings and analysing the quality of the resulting embeddings (section 3.3.1). It then moves to illustrating the creation of a translation graph, an undirected acyclic graph whose edges connect putative translations and synonyms (section 3.3.2). The way this translation graph is partitioned in order to obtain the desired clusters is described in section 3.3.3.

### 3.3.1    Cross-Lingual Embeddings Creation

Given the large size of the joint vocabulary, in order to avoid an excessive increase of the model capacity, the successive issue to face was how to reduce the number of embeddings to train. The proposed method relies on clustering the joint vocabulary's tokens into a number of clusters corresponding to the size of the desired vocabulary (i.e approximately 30k, as in monolingual BERT). The subwords in the joint vocabulary were clustered by semantic similarity exploiting the already available BPEmb-embeddings. To compute semantic similarity (namely cosine similarity) we needed all the embeddings to lay in the same representational space: We started from monolingual BPEmb embeddings of size $d = 300$, and projected them into a shared semantic space by mapping all vectors but English ones into the English vector space.

The mapping was performed using the scripts and seed dictionaries[6] published by Glavas et al. (2019). In particular, the embeddings were mapped using the PROC algorithm, hence using the linear map solving the Procrustes problem (Schönemann 1966). Seed dictionaries were available for English-Russian and English-Finnish, respectively containing 5000 pairs for the training and 2000 for the testing. Dictionaries for the remaining languages were derived with the help of Google Translate. In a more extreme low-resource scenario, a smaller seed dictionary would also be sufficient. Glavas et al. 2019, indeed, showed how by applying a bootstrapping algorithm it is possible to obtain similar performances with seed dictionaries containing only 1000 pairs. Moreover, a completely unsupervised approach is also possible (Artetxe et al. 2018), therefore this approach is applicable also to very low-resource languages for which Google Translate is not available. This method for multi-lingual embeddings generation, therefore, has the great advantage of being easily scalable to a scenario comprising all 275 languages supported by

---

and in non-hashed form

[6]https://github.com/codogogo/xling-eval

Heinzerling and Strube (2018), and potentially any language for which enough raw text is available for the training of monolingual embeddings. The mapping provided us with 300-sized multi-lingual embeddings. In Table 3.1 we report Bilingual Lexicon Induction (BLI) performance: For each source/target language pair, and for each word pair in the relative test dictionary, the top k translations of the source word are extracted from the target language vocabulary as its k nearest neighbours according to cosine similarity. Precision at k (p@k) scores refer to the fraction of times in which the gold translation (the one in the test dictionary) is present among the top k extracted translations. p@1 is reduced to standard precision, i.e. the fraction of correctly extracted translations. The mean reciprocal rank (MRR) gives an idea of the average position of the gold translation among the extracted translations, as expressed by the formula:

$$\frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i}$$

where Q is the set of all the tokens in the target vocabulary. BLI scores show that the mapping was effective, but far from perfect. On average, only one time out of five the ideal translations is chosen, and less than one time out of two it figures in the top-10 candidates. Despite this, if used wisely, BLI predictions can be very useful: The goal of the clustering process is indeed to group together translations, and the BLI predictions are in fact list of translations. The way they are integrated in the clustering method is described in the next sections 3.3.2 and 3.3.3.

| pair | p@1 | p@5 | p@10 | MRR |
|------|-----|-----|------|-----|
| ar→en | 0.169 | 0.388 | 0.443 | 0.268 |
| en→ar | 0.208 | 0.414 | 0.495 | 0.306 |
| bn→en | 0.089 | 0.297 | 0.384 | 0.188 |
| en→bn | 0.201 | 0.375 | 0.447 | 0.285 |
| fi→en | 0.192 | 0.408 | 0.447 | 0.289 |
| en→fi | 0.229 | 0.415 | 0.498 | 0.317 |
| id→en | 0.276 | 0.531 | 0.605 | 0.395 |
| en→id | 0.416 | 0.648 | 0.713 | 0.521 |
| ko→en | 0.103 | 0.340 | 0.425 | 0.217 |
| en→ko | 0.269 | 0.442 | 0.513 | 0.353 |
| ru→en | 0.231 | 0.456 | 0.519 | 0.330 |
| en→ru | 0.284 | 0.483 | 0.552 | 0.378 |
| sw→en | 0.026 | 0.103 | 0.148 | 0.067 |
| en→sw | 0.112 | 0.222 | 0.281 | 0.168 |
| te→en | 0.106 | 0.270 | 0.372 | 0.188 |
| en→te | 0.171 | 0.353 | 0.418 | 0.257 |
| **avg** | **0.193** | **0.384** | **0.454** | **0.283** |

Table 3.1: BLI scores for the multilingual embeddings derived from BPEmb.

### 3.3.2   Translation Graph Creation

Starting from the cross-lingual embeddings obtained in the previous section, we first created a sparse graph containing for each subword the top-k translations in each language. The top-k translations of a subword in a source language were found as its k-nearest neighbours among the subowords in the target language. We therefore eliminated from the sparse translations graphs all non-symmetric edges, so that the obtained graph contained the edge between "subword_1" and "subword_2" if and only if "subword_1" was among the top-k translations of "subword_2" and vice-versa. We refer to this kind of translation couples with the term "strong translations".

To select the most adequate value for k, we created five graphs, one graphs for each value of k between 1 and 5, and we looked at the distribution of cliques sizes in each graph. When looking for the best value of k, we must choose an high enough value to avoid having many subwords which are not part of any clique of size higher than 1 (a clique of size 1 is simply a singleton), since this would make the clustering process impossible. An high value of k, however, can lead to the formation of "wrong" cliques, hence cliques containing subwords which are not reciprocal translations/synonyms. We therefore wanted to find the lowest value of k yielding a reasonably low number of singletons, that is to say cliques containing only one subword. Figure 3.2 shows how, as expected, the number of singletons decreases as k increases. Figure 3.3 shows how the proposed approach also works for strict subwords (i.e. tokens which do not constitute a whole world by themselves, of which the hashed tokens, analysed in Figure 3.3, represent an easy identifiable subset), which follow the same pattern observed for general tokens.
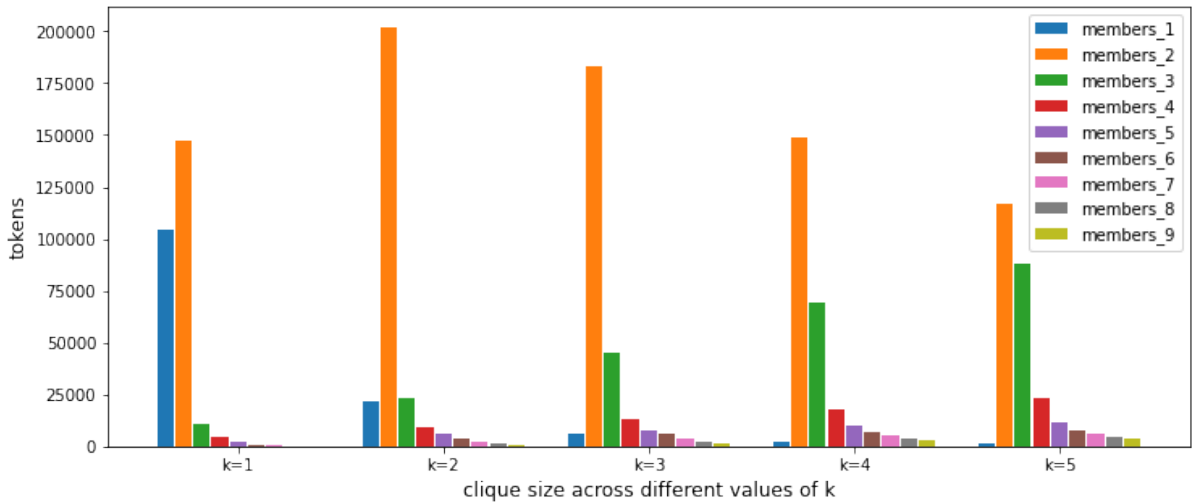


Figure 3.2: Number of tokens belonging to a clique of size at most s ('members_s'), for s in [1,9] and for different values of k.
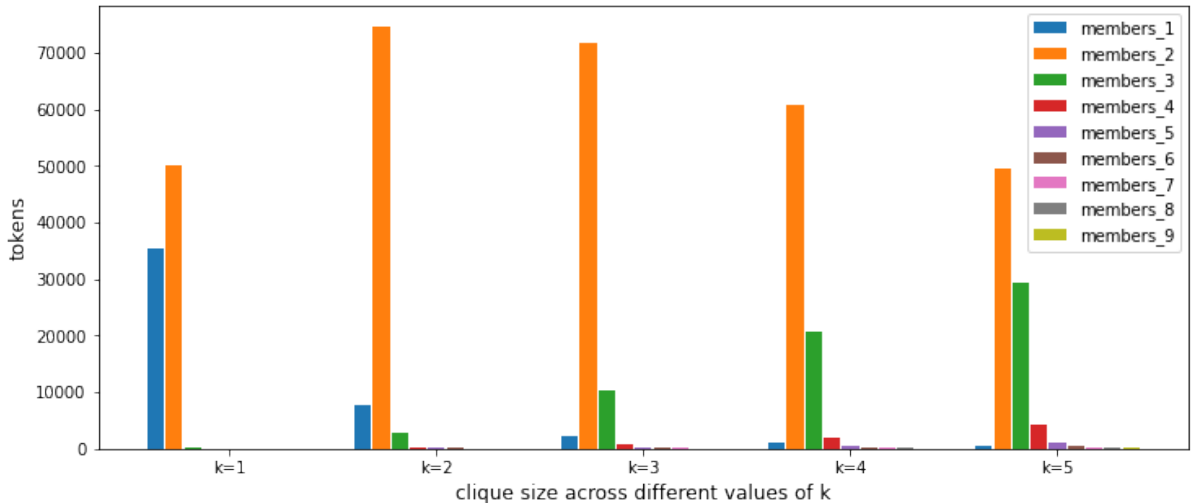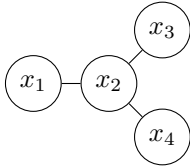


Figure 3.3: Number of hashed tokens belonging to a clique of size at most s ('members_s'), for s in [1,9] and for different values of k.

We selected k=5, it being the smallest value giving a reasonably low number of isolated subwords (subwords only appearing in cliques of size 1, hence with no "strong translations"). In particular, with

k=5, there are 1232 isolated subwords, requiring 1232 dedicated clusters.

Another category of tokens to analyze is the one made of subwords only occurring in cliques of size 2, which could yield very small clusters. To investigate whether the high number of such tokens (43%) could be a problem, we analyzed how many "strong translations" these tokens had. Indeed, even a token not occurring in any clique of size higher than 2 could have several translations, as shown in the example below, where $x_2$ has three translations despite not being in any clique of size higher than 2.

Therefore, more than two tokens can be grouped together even if they do not form a clique, as for nodes $x_1, x_2, x_3, x_4$ in the example above. Figure 3.4 shows how the majority of the subwords occurring only in groups of size 2 have actually more than one "strong translation". Moreover, even subwords with a single strong translation can be grouped with more than one subword, as it happens to $x_1, x_2$ and $x_4$ in the graph reported above. This means that k=5 is high enough to yield a graph partition with only few subgraphs of size 2, since most of them can be merged into larger groups.



Figure 3.4: Distribution of "strong translation" across subwords occurring only in cliques of size 2. Many of these subwords have actually several translations.

Finally, we decided to include nearest neighbours of subwords also in their own language (hence to include synonyms on top of translations), since we observed this did not harm the multi-linguality of the cliques, as shown in Table 3.2.

### 3.3.3    Graph Partitioning

After having created the translation graph, we wanted to partition it to derive subword clusters. In order to obtain clusters of similar size, we applied the METIS balanced graph partitioning algorithm on the sparse translation graph. The METIS algorithm was first introduced in 1995 (Karypis and Kumar

| clique size | avg. number of languages | number of cliques |
|:-----------:|:------------------------:|:-----------------:|
| 1 | 1.00 | 1232 |
| 2 | 1.86 | 957445 |
| 3 | 2.29 | 136731 |
| 4 | 3.16 | 51199 |
| 5 | 3.96 | 29393 |
| 6 | 4.63 | 19350 |
| 7 | 5.27 | 13015 |
| 8 | 5.84 | 8768 |
| 9 | 6.38 | 5340 |
| 10 | 6.77 | 3410 |
| 11 | 7.13 | 2273 |
| 12 | 7.38 | 1161 |
| 13 | 7.49 | 578 |
| 14 | 7.66 | 397 |
| 15 | 7.99 | 175 |
| 16 | 7.96 | 128 |
| 17 | 7.98 | 45 |
| 18 | 8.16 | 25 |
| 19 | 7.50 | 4 |
| 20 | 6.50 | 8 |

Table 3.2: Language coverage per cliques size: for each clique size we report the avg number of languages in cliques of that size, and the total number of cliques of that size.

1995), and is now available both as a standalone software[7] and as a Python package[8]. It solves the multiway graph partitioning problem, which consists in finding a partition of the vertex set into a given number of balanced sets while minimizing cut weight (Hashimoto et al. 2010). The METIS algorithms goes through three main phases:

1. Coarsening Phase: The graph G0 is transformed into a sequence of smaller graphs $G_1, G_2, ..., G_m$ such that $|V_0| > |V_1| > |V_2| > > |V_m|$ .

2. Partitioning Phase: A 2-way partition Pm of the graph $G_m = (V_m, E_m)$ is computed that partitions $V_m$ into two parts, each containing half the vertices of $G_0$.

3. Uncoarsening Phase: The partition $P_m$ of $G_m$ is projected back to $G_0$ by going through intermediate partitions $Pm - 1, Pm - 2, . . . , P_1, P_0$.

What makes this algorithm particularly suitable for our scenario is that it allows to control for the number of parts the graph is partitioned into, and it guarantees these to be balanced. The main parameters of the METIS algorithm are indeed two:

1. **p**: the number of parts (or groups) in which the graph must be divided. Since the partition returned by the algorithms also contains empty groups, and since we wanted to perform some post-processing on the output partition, it was not possible to set the exact value of **p**. To reach a final number of clusters close to the size of the original BERT vocabulary (approximately 30000), we set **p** = 28500, considering that 665 clusters were reserved for the special tokens, and that the post-processing would have increased the number of groups;

[7]http://glaros.dtc.umn.edu/gkhome/metis/metis/overview
[8]https://github.com/networkx/networkx-metis

2. **v**: the load imbalance tolerance. The cluster size must be not higher than $\left\lceil \frac{N}{p} * (1+v) \right\rceil$, where N is the total number of nodes in the graph. This value was set to 1.5, in order to allow for a moderate degree of flexibility with the clusters size, but at the same time avoiding the formation of very large groups of tokens.

With these parameter values, the algorithm produced a partition made of groups of size varying between 1 and 23 (the 81 empty groups were immediately removed). As shown in Figure 3.5, however, the large majority of the groups have size 7, 8 or 9, which, since we are clustering subwords from nine languages, is exactly the desired size.
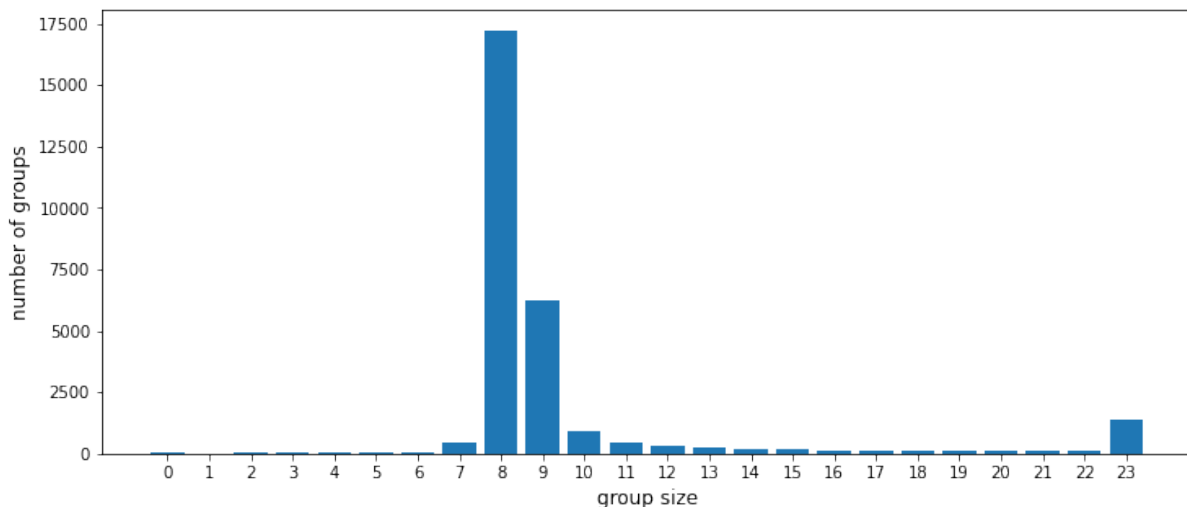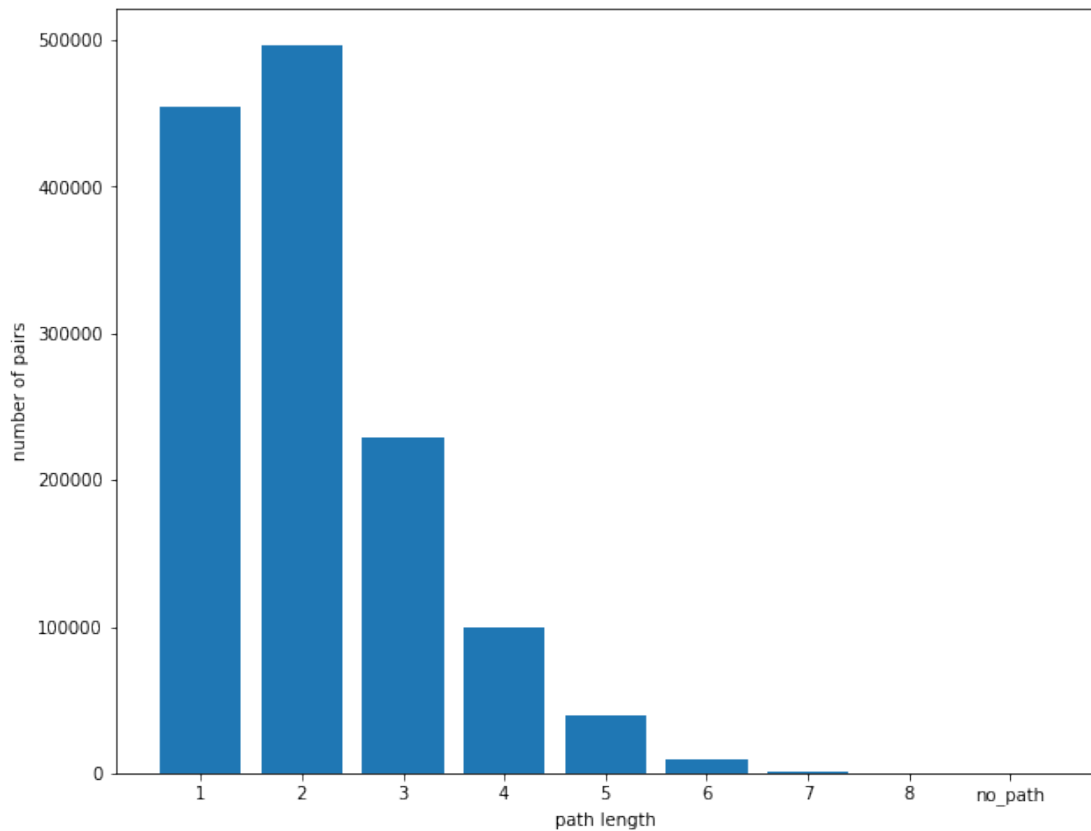


Figure 3.5: Distribution of the sizes of the groups in the partition obtained with the (28500, 1.5) METIS algorithm. Most groups have size 8, close to the number of languages.

The results displayed in Figure 3.5, however, also reveal an issue with the METIS algorithm. If the previous graph analysis detected the presence of 1232 isolated subwords, how can the METIS partition contain only one group of size 1? The answer is that the "cut" objective of the METIS algorithm simply aims to minimize the number of edges cut by the partition. Therefore, it is not necessarily penalized when grouping together disconnected nodes. This holds not only for the METIS algorithm, but for k-way graph partitioning in general, its goal being minimizing the number of edges whose incident vertexes belong to different subsets (Karypis and Kumar 1995). Despite this, however, graph partitioning algorithms prefer grouping together highly connected nodes, therefore the METIS algorithm can prove to be adequate for our task. The problem of isolated subwords can simply be solved by manually removing from the partition the 1232 isolated subwords detected before, and assigning them to dedicated clusters. After doing this, we wondered whether the so-many groups of size 7, 8, and 9 actually reflected the graph structure or they were only an artifact of the METIS algorithm. In order to answer this question, we looked at the distribution of shortest-path lengths across all pairs of subwords belonging to the same group. Figure 3.6a shows that most of the subwords inhabiting the same clusters are immediate neighbours (length=1) or words sharing a neighbour (length=2). As can be seen in Figure 3.6b, only around 100 intra-cluster pairs are not connected. Even if this number is very low, we decided to intervene manually by splitting the 19 groups containing such pairs into 39 groups so that all their members were connected.

Finally, 665 singletons were added for the multi-lingual special tokens. This resulted in a total of 30292 clusters containing 270314 subwords.

**a)** Linear scale.



**b)** Logarithmic scale.



Figure 3.6: Distribution of shortest path lengths across intra-cluster pairs of tokens. Clusters of size 1 where not considered.

## 3.4   Clusters Visualization

A few randomly chosen clusters are reported in Figure 3.7 for illustration. This section does not aim to provide a quantitative analysis, but only a visualization of the clusters. It is nevertheless useful to observe how clusters are made to get an idea of the results of the clustering. Clusters generally contain

subwords from different languages, but not necessarily one token per language. The green cluster, for instance, contains 8 subwords, of which 2 English and 2 Swahili, and it does not contain any Indonesian, Korean, nor Russian subword.



Figure 3.7: Examples of clusters created by the METIS algorithm. Each token is reported both in the original language and in English, except from the cases when providing an English translation is not possible. This happens for many strict subwords in the red and in the yellow clusters. For most clusters is possible to identify an underlying concept: sick (green), pixel/data (cyan), gang/violence (red). This is instead unfeasible for other clusters, especially the ones mainly made of strict subwords (yellow).

In some groups, as the green one and the cyan one, it is easy to identify the concept underlying

the cluster (illness for the former, pixel/data for the latter). In other cases, like the red cluster, more complex connections can be observed. The terms 'gang', 'alley', 'kill him', and somehow 'dragon' can be seen as part of the same context, without being synonyms. The other members of the cluster are harder to interpret, as well as the strict subwords in the yellow cluster. Less interpretable clusters are not necessarily of lower quality, as the model could benefit from them in ways that are not understandable by a human observer. In general, considering the high number of strict subwords in the vocabulary, it is not easy to assess the quality of the clusters. For this reason, we decided to establish whether our approach was effective by training a small BERT model and comparing it with a baseline. The approach taken in the training and in the evaluation is described in the next section.

# Chapter 4

# Training a BERT Model on the Learned Clusters

We decided to train a BERT model mostly following the original BERT paper (Devlin et al. 2019), with some modification motivated by the findings reported by Liu et al. (2019) in their more recent work "RoBERTa: A Robustly Optimized BERT Pretraining Approach":

- **Learning objective:**. Devlin et al. (2019) used two learning objectives during pre-training:

    1. Masked Language Modeling (MLM): A random sample of the input tokens is replaced with the special token [MASK][1]. The MLM objective is defined as a cross-entropy loss on predicting the masked tokens.
    2. Next Sentence Prediction (NSP): A binary classification task consisting in predicting whether two segments are consecutive in the corpus.

    In the footsteps of Liu et al. (2019), we decided to train our model only on the MLM task;

- **Sequence length (T):** The maximum number of tokens in the input. Sequences longer than T are truncated. The authors of the original BERT paper adopted a two-phases training, where the model is trained for 90% of updates with $T = 128$ and for the remaining 10% with $T = 512$. We follow the approach of Liu et al. (2019) and train with a fixed sequence length of $T = 512$;

- **Dynamic Masking:** Devlin et al. (2019) pre-masked the training data in order to avoid that the masking process slowed down the training. This static masking approach, however, cause the same masking to be repeated in different epochs, reducing the variance of the training data. To minimize this phenomenon, they duplicated the training corpus by a factor of 10, so that 10 different masking are actually seen by the model. This solution is not optimal, since it unnecessarily augments the data volume, and since masking redundancy is not completely eliminated (over 40 epochs, each masking is repeated 4 times). We adopt the solution proposed by Liu et al. 2019, called dynamic masking, consisting in generating the masking pattern every time a sequence is fed to the model.

---

[1]15% of the input tokens are selected for possible replacement. Of these, 80% are replaced with [MASK], 10% are replaced by a random token and 10% are left unchanged.

## 4.1   Training Corpus

As training corpus, we concatenated the Wikipedia corpora of all languages. Wikipedia sizes are reported in table 4.1.

| language | Wikipedia size (Mb) | Number of lines (k) |
|:---:|:---:|:---:|
| ar | 1637 | 429 |
| bn | 522 | 74 |
| en | 15284 | 5867 |
| fi | 745 | 287 |
| id | 600 | 224 |
| ko | 703 | 251 |
| ru | 6161 | 149 |
| sw | 35 | 14 |
| te | 498 | 77 |

Table 4.1: Sizes of Wikipedia corpora. For each language we report the raw text size in Mb and the number of lines in thousands.

To compensate for the different Wikipedia sizes across languages, we took an approach similar to the one adopted by Devlin et al. (2019). They performed exponentially smoothed weighting of the data during pre-training data creation. This means that they sampled sentences according to a multinomial distribution $q_{i_{i=1...N}}$, where:

$$q_i = \frac{p_i^\alpha}{\sum_{j=1}^N p_j^\alpha} \qquad \text{being} \quad p_i = \frac{n_i}{\sum_{j=1}^N n_j} \quad ; \tag{4.1}$$

The parameter $\alpha \in (0, 1)$ controls the degree to which low-resource languages are over-sampled: the closer it is to zero, the closer the number of lines across different languages. Devlin et al. (2019) chose a value of $\alpha = 0.7$ in their paper. Successive works, however, used lower values of $\alpha$ (Conneau and Lample 2019, Conneau, Khandelwal, et al. 2020). Xue et al. (2021) recently showed how exponential sampling works best with a value of $\alpha = 0.3$. We therefore used $\alpha = 0.3$ to compute the values of the multinomial distribution $q_{i_{i=1...N}}$. We then created an augmented corpus by keeping fixed the size of the English corpus and duplicating the other languages' data in order to achieve the desired ratios between corpora. During this process we also concatenated/split consecutive lines in the corpus so that each line was made of approximately 512 tokens. This was necessary since the Transformers library handling training on TPUs does not support automatic creation of balanced input segments, but requires the input dataset to be provided as a line by line dataset, i.e. so that each line corresponds to a sequence. If a line is longer than the maximum sequence length T (512 in our case), it is truncated. If it is shorter than T, it is padded with [PAD] tokens. Wikipedia lines exhibit a great variance, ranging from single words to full paragraphs. To avoid waste of data, and that most of the input sequences were made predominantly by [PAD] tokens, we pre-processed the corpus as described above. The ICEBERT cIDs corpus (section 4.2) was created after this pre-processing step, so that the two corpora contained exactly the same lines. The corpora where finally shuffled according to a shared indexes list. This resulted in a training corpus of approximately 80GB for the baseline and 70GB for ICEBERT[2].

---

[2]In UTF-8 encoding foreign characters often takes more than 8 bits, therefore by mapping foreign characters to digits the size of the corpus is reduced. The size of the English corpus, on the other hand, is increased, since many short tokens

## 4.2   ICEBERT pre-training

We called our model ICEBERT, an acronym for Interlingual-Clusters Enhanced BERT. To train our ICEBERT model on cluster IDs instead that actual subwords, we decided to take the following approach:

- At training time, the entire training corpus was mapped to cluster IDs (cIDs). For each sentence in the corpus, this was accomplished in two steps:

  1. The sentence was tokenized using a lowercase BPEmb monolingual tokenizer. Before calling the tokenizer, the sentence was lowercased and digits were replaced by zeros. After the tokenization was completed, all tokens were marked with a "_ln" code;

  2. Each token was mapped to its cID, using the token → cID dictionary created in the clustering phase.

  The model was then trained with a fictitious tokenizer, whose vocabulary contained the string representations of cIDs: "0", "1", ... "30291". In this way, the model's tokenizer simply separates strings of cIDs according to white spaces, while the optimal subword tokenization is made by language-specific tokenizers;

- At inference time, the same actions are performed. This means that the input to the model must be mapped to cIDs before begin fed to the model. This is fairly simple and quick. In most cases (NLI, sentiment analysis, sentence retrieval...) the cID mapping does not represent a problem for the generation of predictions. For some tasks, the model prediction must be modified to take into account the differences between the original input and the mapped one (e.g. span indexes in question answering). The most significant issues would probably occur in text-generative tasks, since a generated cID could correspond to multiple subwords. In this case, the choice of the correct word could be made by restricting the possible output language and by using a statistical language model.

The example below shows how a sentence is first tokenized and then mapped to cIDs:

*"Tyger Tyger, burning bright"*

*['ty_en', '##ger_en', 'ty_en', '##ger_en', ',_en', 'burning_en', 'bright_en']*

*"5159    17398    5159    17398    29669    23261    21619"*

The model's tokenizer then tokenizes the cID-mapped sentence by simply detecting whitespaces:

*['5159', '17398', '5159', '17398', '29669', '23261', '21619']*

## 4.3   Creating a Baseline Model

As a baseline, we trained a BERT model with the same architecture and the same hyperparameters as the ICEBERT model on the same training corpus. The training of a new baseline was necessary since an existing baseline pre-trained in the same condition of the ICEBERT model was not available.

---

(e.g. punctuation) are mapped to 5-digits cluster IDs.

### 4.3.1    Baseline Choice: Balanced Corpus

As a tokenizer for the baseline, we trained a SentencePiece tokenizer on a balanced Wikipedia containing 200M characters for each one of the 9 languages in the experiment. Such a tokenizer already constitutes an improvement over one trained on an unbalanced corpus, since it should tokenize text according to statistics that more fairly reflect the low-resource languages in the corpus. We decided to use a "balanced tokenizer" as a baseline to show that our ICEBERT model improves performance beyond what could have been achieved with such an easy adjustment. Moreover, this allows us to judge the advantages of using monolingual tokenizers *vs.* a single multi-lingual one, controlling for the size of the training data in each language.

### 4.3.2    Baseline Choice: Cased Vocabulary

Early experiments showed how the cased baseline consistently outperforms the uncased one, as also reported in mBERT's documentation[3]. We therefore decided to keep only the cased version as a baseline. As in ICEBERT's case, the SentencePiece vocabulary was converted to a WordPiece format. The vocabulary size of the baseline was set to be the same as the ICEBERT's one[4].

## 4.4    Evaluation

We evaluated our models on the TyDiQA Gold passage task (TyDiQA-GoldP), a question answering task covering the 9 languages we selected for this study. TyDiQA-GoldP is part of the TyDiQA benchmark (Clark et al. 2020b). In this task, the model is given a question and a passage (called "context") that is guaranteed to contain the answer, and it is asked to predict the single contiguous span of characters that answers the question (Clark et al. 2020b). The model is then evaluated with the two metrics proposed in the SQuAD 1.1 paper (Rajpurkar et al. 2016):

1. EM: Exact match, measures the percentage of predictions exactly matching any one of the ground truth answers.

2. F1: (Macro-averaged) F1 score. Computes the F1 score between the prediction and ground truth answers, treating them as bags of tokens. For each prediction, the maximum score over all the ground truth answers is taken. The scores are then averaged over all the questions.

Our models were fine-tuned in a zero-shot scenario, using only the English TyDiQA training dataset. The training data comprised 231 batches of size 16. In order to train the ICEBERT model, the English dataset was mapped to cluster-IDs. The text was mapped following the method described in section 4.2. The input labels (i.e. the span indexes) also had to be mapped to match the indexes of the mapped answers in the mapped contexts. The text in the development dataset was mapped to cluster IDs in the same way as for the training dataset. The span indexes output by the model where then adjusted to extract the correct text span in the original contexts, and the extracted answers were evaluated with EM and F1.

---

[3]https://github.com/google-research/bert/blob/master/multilingual.md
[4]The final vocabulary size resulted to be slighter lower for the baseline (30374 vs 30397) after the removal of tokens made useless by the BertTokenizer preprocessing (e.g. '##,', since punctuation in automatically split and never happens to be hashed by the tokenizer).

## 4.5    Small-Scale Training - ICEBERT$_{\text{SMALL}}$

Before running a large scale experiment, having a significant financial and environmental impact, we tested the efficacy of our method in a smaller-scale scenario.

### 4.5.1    Configuration

In order to support fast experimentation, we adopted the settings proposed by Dufter and Schütze (2020). Our ICEBERT$_{\text{SMALL}}$ model had an hidden size of 64, and intermediate size of 256, 12 hidden layers and a single attention head. We will refer to with this configuration with the term 'small-architecture'.

### 4.5.2    Training Data and Hyperparameters

We pre-trained the small models on a corpus made of 180k Wikipedia sentences (20k per language). The sentences were shuffled to avoid any biased towards a specific language. The models were pre-trained for only 10 epochs, in contrast with the 100 epochs in Dufter and Schütze (2020), to account for the larger amount of training data (around ten times more). We mostly used the same training parameters as Dufter and Schütze (2020):

- optimizer = AdamW (learning rate = 2e-3, weight decay = 0.01, warmup steps = 50)
- max_sequence_length = 128
- batch size = 128[5]

The training required less than ninety minutes per model on a single Tesla T4 GPU on Google Colab. As a baseline, we pre-trained a model of the same size, with the same hyperparameters, and on the same training corpus[6], omitting the mapping to cIDs.

The models were fine-tuned on the English TyDiQA-GoldP dataset for 20 epochs, with an AdamW optimizer with learning rate equal to 5e-5. Fine-tuning took less than an hour per model on a NVidia K80 GPU provided by Kaggle[7]. We pre-trained and fine-tuned both the baseline and the model five different times, and report mean and standard deviation for each language in table 4.2 in section 4.5.3.

### 4.5.3    Small-Scale Training Results

Table 4.2 compares the performance of the baseline and the ICEBERT$_{\text{SMALL}}$ model, showing both language-specific scores and average scores. Average scores refer to all languages except from English. The means and standard deviations of the cross-languages average scores were obtained by first computing the cross-language average score for each run, and then averaging the obtained averages. Table 4.3 reports p-values on a one-way t-test comparing ICEBERT$_{\text{SMALL}}$ to the baseline. In both tables, results with a p-value lower than 0.05 are bolded, and results with a p-value between 0.05 and 0.10 are italicized. A graphical visualization of the small-scale results is given in figure 4.1.

---

[5]Instead of 256, for memory limits reasons.
[6]shuffled using the same random indexes
[7]https://www.kaggle.com/

| ln | EM | | F1 | |
|---|---|---|---|---|
| | baseline | icebert | baseline | icebert |
| ar | 1.41 (0.31) | **5.15 (0.81)** | 19.73 (1.54) | **21.59 (0.77)** |
| bn | 0.89 (0.79) | **3.54 (1.48)** | 8.76 (0.77) | **10.83 (2.13)** |
| en | **19.27 (0.89)** | 14.82 (2.37) | **30.01 (1.43)** | 26.15 (2.62) |
| fi | 2.12 (0.83) | **5.24 (1.47)** | 15.27 (2.14) | **18.11 (1.66)** |
| id | 4.21 (1.05) | **6.44 (0.81)** | 18.70 (1.43) | **21.64 (0.98)** |
| ko | 1.01 (0.36) | **2.97 (0.48)** | 7.44 (0.74) | **9.29 (0.70)** |
| ru | 1.65 (0.69) | **8.08 (1.45)** | 11.84 (1.81) | **16.50 (1.04)** |
| sw | 2.44 (1.71) | **4.85 (1.29)** | 14.42 (2.55) | **17.52 (2.03)** |
| te | 0.66 (0.51) | **1.61 (0.52)** | 6.95 (0.61) | 6.91 (1.07) |
| avg | 1.80 (0.52) | **4.74 (0.89)** | 12.89 (0.77) | **15.30 (0.83)** |

Table 4.2: Small-architecture EM and F1 scores on the TyDiQA-GoldP task. Results are averaged over 5 runs. Standard deviation is reported in brackets.

| ln | p-value | |
|---|---|---|
| | EM | F1 |
| ar | **<0.001** | **0.021** |
| bn | **0.004** | **0.037** |
| en | 0.998 | 0.99 |
| fi | **0.002** | **0.023** |
| id | **0.003** | **0.003** |
| ko | **<0.001** | **0.002** |
| ru | **<0.001** | **0.001** |
| sw | **0.018** | **0.033** |
| te | **0.009** | 0.522 |
| avg | **<0.001** | **0.001** |

Table 4.3: p-values of results reported in table 4.2, on a one-way t-test.
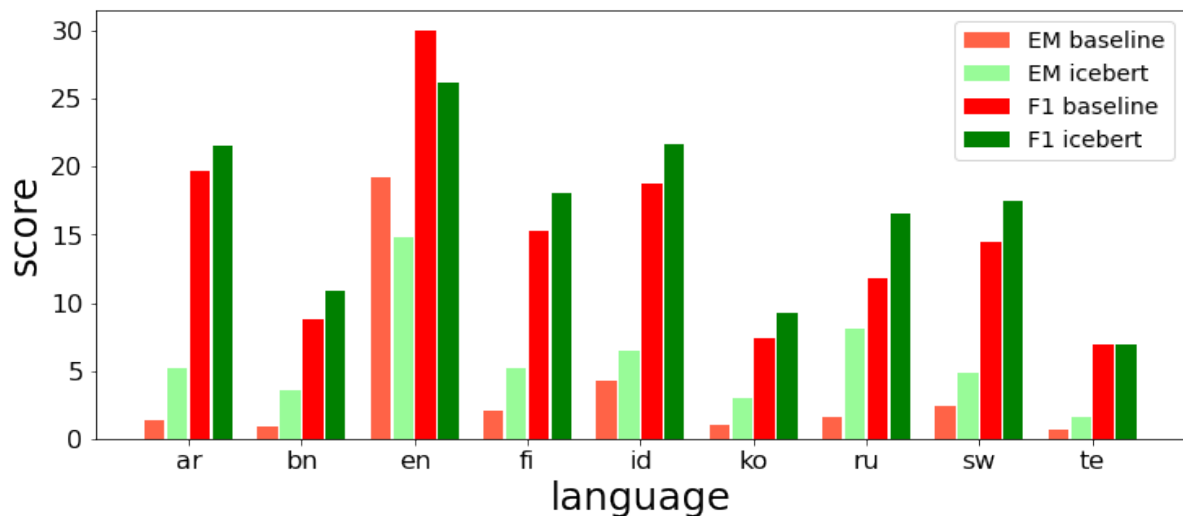


Figure 4.1: EM and F1 scores for the small baseline and ICEBERT model, across 9 languages.

ICEBERT$_{\text{SMALL}}$ appears to perform better than the baseline for every language except from English. The difference is particularly significant for Exact Match scores, which increase significantly for all non-English languages (+163%), but a significant increment can also be seen for F1 (+18.7%).

The drop in English performance was expected: The model benefits from the clustering when evaluated on languages that were not present in the training data. Since all training data is in English, ICEBERT$_{\text{SMALL}}$ does not improve over the baseline, which learns directly from English sentences, without passing through the cID-mapping.

## 4.6   Large-Scale Training - ICEBERT$_{\text{BASE}}$

The positive results obtained in the small-scale experiment motivated the training of a larger ICEBERT model, to investigate whether similar improvements could be achieved also over a better baseline. In the large-scale setting we trained two models of the size of the original BERT model, BERT$_{\text{BASE}}$. Our ICEBERT$_{\text{BASE}}$ model therefore had an hidden size of 768, and intermediate size of 3072, 12 hidden layers and 12 attention heads.

We will refer to this configuration with the term 'base-architecture'. We mostly used the trainer hyper-parameters reported by Devlin et al. 2019 in the BERT paper:

- optimizer = AdamW:
  - ◎ learning rate = 1e-4 (linear decay)
  - ◎ $\beta_1 = 0.9$, $\beta_2 = 0.999$
  - ◎ weight decay = 0.01
  - ◎ epsilon = 1e-6
  - ◎ warmup steps = 10,000
- dropout probability = 0.1
- activation function = gelu
- training steps = 1,000,000
- max_sequence_length = 512[8]
- batch size = 64 (8 per core)[9]

It is important to notice that, despite the total number of training steps being the same as in BERT$_{\text{BASE}}$, the lower batch size implies that the model is actually trained on $\frac{1}{4}$ of the training data. A full-training was not feasible with the given resources (more than 3 weeks per model on a single v3-8 TPU), but we decided to undertrain rather than downscaling, based on the findings of Li et al. (2020), who claimed that larger models tend to achieve better performance than smaller models when trained for the same time.

The training took approximately 5 days per model on a GCP cluster made of a nd2-highmem-8 compute engine (8 vCPUs, 64 GB memory) accelerated by a v3-8 TPU node (8 cores, 128 GiB of TPU memory). The cost of the training on Google Cloud Platform amounted to little more than 60$ per model for the compute engine, while the TPUs where made available for free by Google as part of the TRC program[10]. The costs covered by this offer amount to 300$ per model if using a preemptible v3-8 TPUs, or to 1000$ per model if using an on-demand device.

---

[8]instead of adopting the two-phases approach, as explained in the introduction of chapter 4.

[9]Instead of 256, for memory limits reasons.

[10]https://sites.research.google/trc/

## 4.7  Main Libraries

We trained our models exploiting the Python Transformers[11] library (T. Wolf et al. 2020). This library offers modules and scripts for pre-training and fine-tuning several models. In particular, we used the BertForMaskedLM class, which makes it easy to instantiate and train a Neural Network built according to the BERT architecture.

The transformers library is backed by Jax[12], PyTorch[13], and TensorFlow[14]. We choose PyTorch for our implementation. On top of the standard torch module[15], supporting training on GPUs, we also exploited the torch_xla module[16], allowing PyTorch to run on TPUs.

---

[11]https://github.com/huggingface/transformers
[12]https://jax.readthedocs.io/en/latest/
[13]https://pytorch.org/
[14]https://www.tensorflow.org/
[15]https://pytorch.org/docs/stable/torch.html
[16]https://pytorch.org/xla/release/1.8/index.html

# Chapter 5

# Results and Discussion

Table 5.1 reports EM and F1 scores on the TyDiQA dataset for the base-architecture scenario, comparing the fine-tuned baseline and the ICEBERT$_{\text{BASE}}$ model. The scores of the fully trained BERT$_{\text{BASE}}$ model, as reported in Hu et al. (2020), are also shown for comparison. Table 5.2 reports p-values on a one-way t-test comparing ICEBERT$_{\text{BASE}}$ to the baseline. In both tables, results with a p-value lower than 0.05 are bolded, and results with a p-value between 0.05 and 0.10 are italicized. A graphical visualization of the base-scale results is given in figure 5.1.

| ln | EM | | | F1 | | |
|---|---|---|---|---|---|---|
| | baseline | icebert | xtreme | baseline | icebert | xtreme |
| ar | 9.08 (0.57) | **13.10 (1.24)** | 42.80 | 22.65 (1.32) | **29.32 (0.67)** | 62.20 |
| bn | 8.55 (0.83) | 7.08 (1.91) | 32.70 | 15.35 (0.56) | **17.45 (1.28)** | 49.30 |
| en | 34.32 (1.96) | **37.58 (1.42)** | 63.60 | 43.79 (1.72) | **47.27 (1.39)** | 75.30 |
| fi | 10.57 (1.68) | *12.19 (0.71)* | 45.30 | 24.90 (1.44) | **27.60 (1.44)** | 59.70 |
| id | 21.06 (0.58) | 20.71 (1.15) | 45.80 | 34.52 (0.48) | *35.91 (1.47)* | 64.80 |
| ko | **6.88 (0.89)** | 1.81 (0.30) | 50.00 | **13.82 (0.58)** | 6.56 (0.84) | 58.80 |
| ru | 10.96 (0.92) | *12.44 (1.09)* | 38.80 | 21.85 (0.76) | **26.18 (1.23)** | 60.00 |
| sw | 16.37 (0.93) | 16.43 (0.43) | 37.90 | 28.44 (1.43) | 28.38 (0.82) | 57.50 |
| te | 4.24 (0.43) | **7.12 (0.43)** | 38.40 | 10.38 (0.83) | **14.03 (1.60)** | 49.60 |
| avg | 10.96 (0.69) | 11.36 (0.15) | 41.46 | 21.49 (0.46) | **23.18 (0.64)** | 57.74 |

Table 5.1: Base-architecture EM and F1 scores on the TyDiQA-GoldP task. Results are averaged over 3 fine-tunings. Standard deviation is reported in brackets.

Tables 5.1 and 5.2 show that in the base-architecture scenario ICEBERT$_{\text{BASE}}$ performs better than the baseline in general, even if not across all languages and evaluation metrics. An increment in F1 performance can be observed for Arabic, Bengali, English, Finnish, Russian, and Telugu, while the model achieves comparable performance on Indonesian and Swahili. ICEBERT$_{\text{BASE}}$ performs particularly bad for Korean. The most significant improvements can be observed on Arabic and Telugu. The latter is particularly important, being the language with the lowest scores among the ones included in the study. The average score of non-English languages is also significantly better for ICEBERT$_{\text{BASE}}$ than for the baseline.

| ln | p-value | |
|---|---|---|
| | **EM** | **F1** |
| ar | **0.004** | **0.001** |
| bn | 0.856 | **0.03** |
| en | **0.04** | **0.027** |
| fi | *0.099* | **0.042** |
| id | 0.671 | *0.098* |
| ko | 1.0 | 1.0 |
| ru | *0.073* | **0.003** |
| sw | 0.458 | 0.524 |
| te | **0.001** | **0.012** |
| avg | 0.194 | **0.01** |

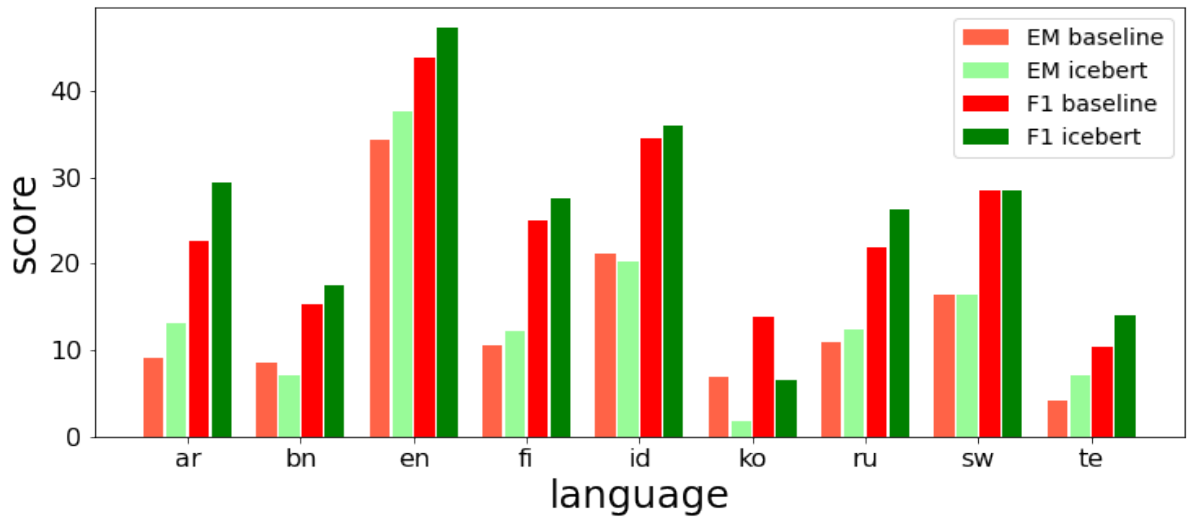Table 5.2: p-values of results reported in table 5.1, on a one-way t-test.



Figure 5.1: EM and F1 scores for the base baseline and ICEBERT model, across 9 languages.

According to the EM metric, our model significantly improves over the baseline only for Arabic, English, and Telugu, and to a minor extent for Finnish and Russian. ICEBERT$_{\text{BASE}}$ scores are comparable with the baseline's ones for the remaining languages, except from Korean, for which a drastic drop in performance can be observed as in the F1 case. The very low Korean score heavily penalizes the average EM score of our model. By Excluding Korean from the average, ICEBERT$_{\text{BASE}}$ gets a better EM score than the baseline (12.73 vs 11.55, p-value=0.023) on the remaining 7 languages.

The most relevant[1] results we have reported, F1 scores, are also the most encouraging ones. Indeed, ICEBERT$_{\text{BASE}}$ improves over the baseline over almost all languages. Excluding Korean, whose performance are clearly harmed by some factors which should be more deeply investigated, the only two languages for which significant improvements could not be seen are also the ones which are more similar to English (Indonesian and Swahili), and for which the proposed approach was less necessary. Languages which are typologically distant from English, like Arabic, Russian, and Telugu, did benefit from the clustering, as shown by the reported increments in performance. Surprisingly, English itself

---

[1]F1 is considered to be more informative than EM, and the base-architecture's results are clearly more important than the small-architecture's ones.

did benefit from the clustering in the base-architecture scenario. The fact that this did not happen in the small-architecture case can be explained by looking at the fine-tuning hyperparameters. The small models needed to be fine-tuned for more epochs and with a higher learning rate to reach reasonably good scores. Since the training data was in English, however, this caused the model to overfit, and to exhibit much higher performance for English that for other languages (10 times the average for EM and 2.5 times the average for F1). ICEBERT$_{\text{SMALL}}$, on the other hand, was not fine-tuned on the original English dataset, but on its cID-mapped version, therefore the overfitting was less severe in that case, leading to higher performance on non-English languages and to poorer performance on English.

Despite the observed improvements, the EM and F1 scores obtained by both the baseline and ICEBERT$_{\text{BASE}}$ remain much lower than the ones of the fully trained BERT$_{\text{BASE}}$ model, as reported by Hu et al. (2020). This was expected, since our model is not fully trained.

Figure 5.2 shows how the MLM loss decreases during the pre-training. ICEBERT's loss decreases faster than the baseline's one, which remains stuck at a value of approximately six for more than 100'000 training steps. The ICEBERT model could therefore be preferred in scenarios with little resources. Neither the baseline nor the ICEBERT model did converge after only 1'000'000 steps, and it can be expected that training for more epochs would improve their performance.
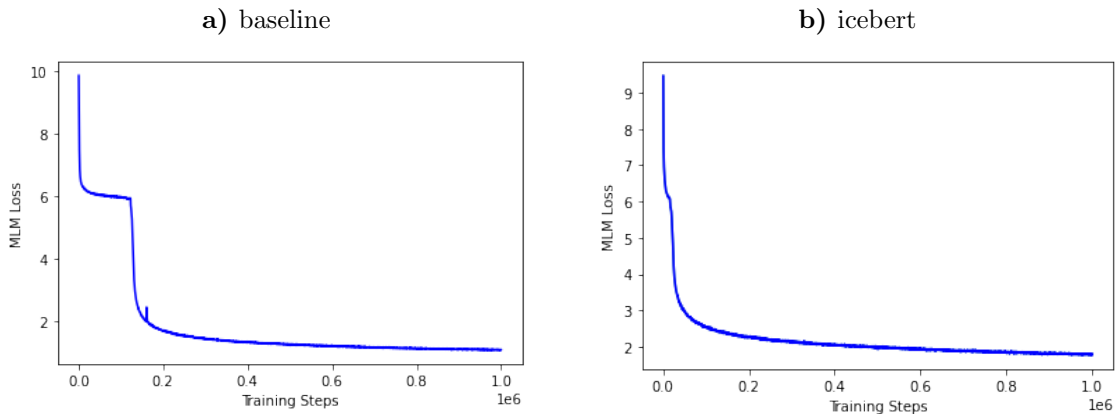
**a)** baseline                                                    **b)** icebert



Figure 5.2: Pretraining MLM loss over one million training steps

An analysis of the baseline's and the model's fine-tuning is reported in Appendix A. Except from some failed runs occurred when fine-tuning the baseline (and not when fine-tuning the ICEBERT model), no significant differences were observed between the two.

# Chapter 6

# Future Work

According to the results reported in the previous section, our approach proved to be effective in increasing the performance of a baseline multi-lingual model. The bad results obtained for Korean, however, must be a warning, and should prompt a careful investigation of the reasons which caused this phenomenon. In the best scenario, the issue does not concern the pre-trained model itself, but the evaluation method. The mapping of the TyDiQA dataset to cIDs, indeed, was not a straightforward process. For many Korean examples (104/448), the mapped answer did not exactly correspond to the text extracted by slicing the mapped context with the derived span indexes. The same phenomenon could be observed for Telugu (in 278/1273 examples), while only few errors occurred in other languages. The almost totality of these errors occurred in examples where the answer started or ended "inside a word", i.e. with a suffix or a prefix and not with a full word. In such cases, if the tokenizer fails to split a word in the specific way that would permit to correctly extract the answer for the context, an exact match cannot be achieved. Therefore, this is an issue not directly related to our mapping method, but to the tokenizer. Paradoxically, a better (richer) tokenizer would be able to recognize longer words, therefore failing to split tokens in the specific way necessary to extract the correct answers more often than a tokenizer with a smaller vocabulary. Besides these phenomenon, other less easily observable issues linked to tokenization could have compromised the dataset mapping, leading to an underestimation of ICEBERT's performance. Further work is necessary to investigate this phenomena. An alternative explanation for Korean's poor performance could be found in the clusters themselves. It is indeed possible that the clustering process was more effective for some languages than for others. Optimizing the clustering process and analyzing its effectiveness across different languages, not necessarily only the ones included in this study, would be a useful goal for possible future works. The fact that the drop in performance for Korean was absent in the preliminary experiments could raise the suspect that some data corruption took place during the creation of the large-scale experiment training corpus. Despite being improbable that this only occurred for Korean, and even though a shallow analysis of the training corpus did not reveal any data corruption, this remains an hypothesis to be taken into account.

The improvement on Indonesian and Swahili in the small-architecture scenario, which could not been observed in the larger scale case, could be explained by the fact that the lower baseline scores made it easier for the model to improve. This leads to wonder whether the improvements observed over our undertrained BERT$_{\text{BASE}}$ baseline could also be achieved when fully training the BERT$_{\text{BASE}}$ models. In future works, it will be important to apply our method to a truly state-of-the-art baseline, to assess

whether the state-of-the-art itself could be improved through our approach. This could mean not only training a full-scale BERT model, but also applying the same method to different model architectures.

Another important step to take in future works is to extend the proposed method to other languages. The good results obtained across the very different languages included in this study are encouraging, suggesting that our approach could be valid for any language, no matter how typologically distant from English. The natural next step would therefore be to extend the clustering process to comprise all 104 languages handled by mBERT. This could be easily done thanks to the high scalability of our method. Monolingual tokenizers and embeddings are made available for 275 languages by Heinzerling and Strube (2018), and can in principle be trained for any language for which a reasonably large amount of raw text data is available. The embeddings mapping necessary for the clustering does not necessarily require bilingual seed dictionaries, since completely unsupervised approaches have been proven to be successful by Conneau, Lample, et al. (2017), Artetxe et al. (2018), Hoshen and L. Wolf (2018), and Alvarez-Melis and Jaakkola (2018). With an higher number of languages, the benefits of the clustering could result to be higher than what observed in this study. With a single multi-lingual vocabulary of fixed size, indeed, increasing the number of languages implies reducing the language specific vocabulary each language can exploit, especially when multiple scripts are present and vocabulary sharing is limited. This explains why mBERT's vocabulary ($\sim$120k tokens) is significantly bigger than BERT's one ($\sim$30k tokens) and, notwithstanding that, some low resource languages have little to none dedicated vocabulary. When using the proposed clustering method, including more languages would be possible without reducing the size of the monolingual vocabularies. Albeit it cannot be denied that the clustering process would be more complex with more languages, in principle the goal of creating clusters of translations would remain reachable. If in our experiments we aimed to create clusters of approximately 9 tokens, with 100 languages the goal would become to create cluster with 100 tokens. The improvements due to a better tokenization, therefore, would be maximized in a scenario with more languages.

Finally, it would be important to evaluate ICEBERT on a wider range of downstream tasks. This would be possible provided that the model supported more language, since different multi-lingual downstream tasks typically cover different languages.

# Chapter 7

# Conclusion

This study aimed to improve multi-lingual language models by training them on clusters of monolingual segments. The proposed approach, focusing on 9 languages, yielded good quality clusters, able to group semantically similar words and subwords across languages. This led to improvements over standard segmentation and training methods on the majority of the included languages, when evaluated on a question answering task. More work is necessary to investigate the performance of our method on other languages and downstream tasks, and to develop more efficient methods for clustering multi-lingual subwords, in order to maximize the advantages of the proposed approach.

# Bibliography

Agić, Ž., & Vulić, I. (2019). JW300: A wide-coverage parallel corpus for low-resource languages. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3204–3210. https://doi.org/10.18653/v1/P19-1310

Alvarez-Melis, D., & Jaakkola, T. (2018). Gromov-Wasserstein alignment of word embedding spaces. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 1881–1890. https://doi.org/10.18653/v1/D18-1214

Ammar, W., Mulcaire, G., Tsvetkov, Y., Lample, G., Dyer, C., & Smith, N. A. (2016). Massively multilingual word embeddings. *CoRR*, *abs/1602.01925*. http://dblp.uni-trier.de/db/journals/corr/corr1602.html#AmmarMTLDS16

Artetxe, M., Labaka, G., & Agirre, E. (2016). Learning principled bilingual mappings of word embeddings while preserving monolingual invariance. *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2289–2294. https://doi.org/10.18653/v1/D16-1250

Artetxe, M., Labaka, G., & Agirre, E. (2017). Learning bilingual word embeddings with (almost) no bilingual data. *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 451–462. https://doi.org/10.18653/v1/P17-1042

Artetxe, M., Labaka, G., & Agirre, E. (2018). A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 789–798. https://doi.org/10.18653/v1/P18-1073

Bahdanau, D., Cho, K., & Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In Y. Bengio & Y. LeCun (Eds.), *International conference on learning representations*.

Bengio, Y., Ducharme, R., & Vincent, P. (2000). A neural probabilistic language model. In T. K. Leen, T. G. Dietterich, & V. Tresp (Eds.), *Nips* (pp. 932–938). MIT Press. http://dblp.uni-trier.de/db/conf/nips/nips2000.html#BengioDV00

Berasategi, A. (2020). Overview of tokenization algorithms in nlp (T. D. Science, Ed.) [Accessed: 2021-01-07]. https://towardsdatascience.com/overview-of-nlp-tokenization-algorithms-c41a7d5ec4f9

Bozinovski, Z., & Fulgosi, A. (1976). The influence of pattern similarity and transfer learning upon the training of a base perceptron b2. *Proceedings of Symposium Informatica*, *3-121-5*.

Chandar, A. P. S., Lauly, S., Larochelle, H., Khapra, M. M., Ravindran, B., Raykar, V. C., & Saha, A. (2014). An autoencoder approach to learning bilingual word representations. *CoRR*, *abs/1402.1454*. http://dblp.uni-trier.de/db/journals/corr/corr1402.html#PLLKRRS14

Christodoulopoulos, C., & Steedman, M. (2015). A massively parallel corpus: The bible in 100 languages. *Language Resources and Evaluation*, *49*, 375–395.

Chung, H. W., Garrette, D., Tan, K. C., & Riesa, J. (2020). Improving multilingual models with language-clustered vocabularies. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4536–4546. https://doi.org/10.18653/v1/2020.emnlp-main.367

Clark, J. H., Choi, E., Collins, M., Garrette, D., Kwiatkowski, T., Nikolaev, V., & Palomaki, J. (2020a). Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*.

Clark, J. H., Choi, E., Collins, M., Garrette, D., Kwiatkowski, T., Nikolaev, V., & Palomaki, J. (2020b). Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. *CoRR*, *abs/2003.05002*. http://dblp.uni-trier.de/db/journals/corr/corr2003.html#abs-2003-05002

Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., & Stoyanov, V. (2020). Unsupervised cross-lingual representation learning at scale. In D. Jurafsky, J. Chai, N. Schluter, & J. R. Tetreault (Eds.), *Acl* (pp. 8440–8451). Association for Computational Linguistics. http://dblp.uni-trier.de/db/conf/acl/acl2020.html#ConneauKGCWGGOZ20

Conneau, A., & Lample, G. (2019). Cross-lingual language model pretraining. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. B. Fox, & R. Garnett (Eds.), *Neurips* (pp. 7057–7067). http://dblp.uni-trier.de/db/conf/nips/nips2019.html#ConneauL19

Conneau, A., Lample, G., Ranzato, M., Denoyer, L., & Jégou, H. (2017). Word translation without parallel data. *CoRR*, *abs/1710.04087*. http://dblp.uni-trier.de/db/journals/corr/corr1710.html#abs-1710-04087

Dai, A. M., & Le, Q. V. [Quoc V]. (2015). Semi-supervised sequence learning. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information processing systems* (pp. 3079–3087). Curran Associates, Inc. https://proceedings.neurips.cc/paper/2015/file/7137debd45ae4d0ab9aa953017286b20-Paper.pdf

Deerwester, S. C., Dumais, S. T., Landauer, T. K., Furnas, G. W., & Harshman, R. A. (1990). Indexing by latent semantic analysis. *Journal of the American Society of Information Science*, *41*(6), 391–407.

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *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)*, 4171–4186. https://doi.org/10.18653/v1/N19-1423

Dufter, P., & Schütze, H. (2020). Identifying necessary elements for bert's multilinguality. http://nbn-resolving.de/urn/resolver.pl?urn=nbn:de:bvb:19-epub-72199-8

Faruqui, M., & Dyer, C. (2014). Improving vector space word representations using multilingual correlation. *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, 462–471. https://doi.org/10.3115/v1/E14-1049

Firth, J. R. (1957). A synopsis of linguistic theory 1930-55. *1952-59*, 1–32.

Gage, P. (1994). A new algorithm for data compression. *C Users Journal*.

Glavas, G., Litschko, R., Ruder, S., & Vulic, I. (2019). How to (properly) evaluate cross-lingual word embeddings: On strong baselines, comparative analyses, and some misconceptions. In A. Korhonen, D. R. Traum, & L. Màrquez (Eds.), *Acl (1)* (pp. 710–721). Association for Computational Linguistics. http://dblp.uni-trier.de/db/conf/acl/acl2019-1.html#GlavasLRV19

Graves, A. (2013). Generating sequences with recurrent neural networks. *CoRR*, *abs/1308.0850*. http://dblp.uni-trier.de/db/journals/corr/corr1308.html#Graves13

Hashimoto, H., Sonobe, Y., & Yagiura, M. (2010). A multilevel scheme with adaptive memory strategy for multiway graph partitioning. In C. Blum & R. Battiti (Eds.), *Learning and intelligent optimization* (pp. 188–191). Springer Berlin Heidelberg.

Heinzerling, B., & Strube, M. (2018). BPEmb: Tokenization-free Pre-trained Subword Embeddings in 275 Languages. In N. C. ( chair), K. Choukri, C. Cieri, T. Declerck, S. Goggi, K. Hasida, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, S. Piperidis, & T. Tokunaga (Eds.), *Proceedings of the eleventh international conference on language resources and evaluation (lrec 2018)*. European Language Resources Association (ELRA).

Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, *79*(8), 2554–2558. https://doi.org/10.1073/pnas.79.8.2554

Hoshen, Y., & Wolf, L. (2018). An iterative closest point method for unsupervised word translation. *CoRR*, *abs/1801.06126*. http://dblp.uni-trier.de/db/journals/corr/corr1801.html#abs-1801-06126

Hu, J., Ruder, S., Siddhant, A., Neubig, G., Firat, O., & Johnson, M. (2020). Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalization. *CoRR*, *abs/2003.11080*. http://dblp.uni-trier.de/db/journals/corr/corr2003.html#abs-2003-11080

K, K., Wang, Z., Mayhew, S., & Roth, D. (2020). Cross-lingual ability of multilingual bert: An empirical study. *ICLR*. http://dblp.uni-trier.de/db/conf/iclr/iclr2020.html#KWMR20

Kalchbrenner, N., Espeholt, L., Simonyan, K., van den Oord, A., Graves, A., & Kavukcuoglu, K. (2016). Neural machine translation in linear time. *CoRR*, *abs/1610.10099*. http://dblp.uni-trier.de/db/journals/corr/corr1610.html#KalchbrennerESO16

Karpathy, A., Johnson, J., & Li, F.-F. (2015). Visualizing and understanding recurrent networks. *CoRR*, *abs/1506.02078*. http://dblp.uni-trier.de/db/journals/corr/corr1506.html#KarpathyJL15

Karypis, G., & Kumar, V. (1995). Multilevel graph partitioning schemes. In K. A. Gallivan (Ed.), *Icpp (3)* (pp. 113–122). CRC Press. http://dblp.uni-trier.de/db/conf/icpp/icpp1995-3.html# KarypisK95

Klementiev, A., Titov, I., & Bhattarai, B. (2012). Inducing crosslingual distributed representations of words. *Proceedings of COLING 2012*, 1459–1474. https://www.aclweb.org/anthology/C12-1089

Kudo, T. (2018). Subword regularization: Improving neural network translation models with multiple subword candidates. In I. Gurevych & Y. Miyao (Eds.), *Acl (1)* (pp. 66–75). Association for Computational Linguistics. http://dblp.uni-trier.de/db/conf/acl/acl2018-1.html#Kudo18

Kudo, T., & Richardson, J. (2018). Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In E. Blanco & W. Lu (Eds.), *Emnlp (demonstration)* (pp. 66–71). Association for Computational Linguistics. http://dblp.uni-trier.de/db/conf/ emnlp/emnlp2018-d.html#KudoR18

Lazaridou, A., Dinu, G., & Baroni, M. (2015). Hubness and pollution: Delving into cross-space mapping for zero-shot learning. *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 270–280. https://doi.org/10.3115/v1/P15-1027

Lee, J., Cho, K., & Hofmann, T. (2017). Fully character-level neural machine translation without explicit segmentation. *Transactions of the Association for Computational Linguistics*, *5*, 365–378. https://doi.org/10.1162/tacl_a_00067

Levy, O., Søgaard, A., & Goldberg, Y. (2017). A strong baseline for learning cross-lingual word embeddings from sentence alignments. In M. Lapata, P. Blunsom, & A. Koller (Eds.), *Eacl (1)* (pp. 765–774). Association for Computational Linguistics. http://dblp.uni-trier.de/db/conf/ eacl/eacl2017-1.html#SogaardGL17

Li, Z., Wallace, E., Shen, S., Lin, K., Keutzer, K., Klein, D., & Gonzalez, J. E. (2020). Train large, then compress: Rethinking model size for efficient training and inference of transformers. *CoRR*, *abs/2002.11794*. http://dblp.uni-trier.de/db/journals/corr/corr2002.html#abs-2002-11794

Libovický, J., Rosa, R., & Fraser, A. (2019). How language-neutral is multilingual bert? *CoRR*, *abs/1911.03310*. http://dblp.uni-trier.de/db/journals/corr/corr1911.html#abs-1911-03310

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *CoRR*, *abs/1907.11692*. http://dblp.uni-trier.de/db/journals/corr/corr1907.html#abs-1907-11692

Macukow, B. (2016). Neural networks - state of art, brief history, basic models and architecture. In K. Saeed & W. Homenda (Eds.), *Cisim* (pp. 3–14). Springer. http://dblp.uni-trier.de/db/conf/ cisim/cisim2016.html#Macukow16

Mikolov, T., Karafiát, M., Burget, L., Cernock, J., & Khudanpur, S. (2010). Recurrent neural network based language model. *INTERSPEECH*, *2*, 3.

Mikolov, T., Le, Q. V. [Quoc V.], & Sutskever, I. (2013). Exploiting similarities among languages for machine translation. *CoRR*, *abs/1309.4168*. http://dblp.uni-trier.de/db/journals/corr/corr1309.html#MikolovLS13

Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global vectors for word representation. *EMNLP*, *14*, 1532–1543. https://nlp.stanford.edu/pubs/glove.pdf

Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). Deep contextualized word representations. In M. A. Walker, H. Ji, & A. Stent (Eds.), *Naacl-hlt* (pp. 2227–2237). Association for Computational Linguistics. http://dblp.uni-trier.de/db/conf/naacl/naacl2018-1.html#PetersNIGCLZ18

Pires, T., Schlinger, E., & Garrette, D. (2019). How multilingual is multilingual BERT? *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4996–5001. https://doi.org/10.18653/v1/P19-1493

Pyysalo, S., Kanerva, J., Virtanen, A., & Ginter, F. (2020). Wikibert models: Deep transfer learning for many languages. *CoRR*, *abs/2006.01538*. http://dblp.uni-trier.de/db/journals/corr/corr2006.html#abs-2006-01538

Radford, A., Józefowicz, R., & Sutskever, I. (2017). Learning to generate reviews and discovering sentiment. *CoRR*, *abs/1704.01444*. http://dblp.uni-trier.de/db/journals/corr/corr1704.html#RadfordJS17

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). *Language models are unsupervised multitask learners*. https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf

Rajpurkar, P., Zhang, J., Lopyrev, K., & Liang, P. (2016). Squad: 100, 000+ questions for machine comprehension of text. *CoRR*, *abs/1606.05250*. http://dblp.uni-trier.de/db/journals/corr/corr1606.html#RajpurkarZLL16

Rosenfeld, R. (2000). Two decades of statistical language modeling: Where do we go from here? *Proceedings of the IEEE*, *88*, 1270–1278.

Ruder, S., Vulic, I., & Søgaard, A. (2019). A survey of cross-lingual word embedding models. *J. Artif. Intell. Res.*, *65*, 569–631. http://dblp.uni-trier.de/db/journals/jair/jair65.html#RuderVS19

Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning internal representations by error propagation. *Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1: Foundations* (pp. 318–362). MIT Press.

Schönemann, P. (1966). A generalized solution of the orthogonal procrustes problem. *Psychometrika*, *31*, 1–10. https://doi.org/https://doi-org.proxy.library.uu.nl/10.1007/BF02289451

Schuster, M., & Nakajima, K. (2012). Japanese and korean voice search. *ICASSP*, 5149–5152. http://dblp.uni-trier.de/db/conf/icassp/icassp2012.html#SchusterN12

Sennrich, R., Haddow, B., & Birch, A. (2016). Neural machine translation of rare words with subword units. *ACL (1)*. http://dblp.uni-trier.de/db/conf/acl/acl2016-1.html#SennrichHB16a

Singh, J., McCann, B., Socher, R., & Xiong, C. (2019). BERT is not an interlingua and the bias of tokenization. *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, 47–55. https://doi.org/10.18653/v1/D19-6106

Smith, S. L., Turban, D. H. P., Hamblin, S., & Hammerla, N. Y. (2017). Offline bilingual word vectors, orthogonal transformations and the inverted softmax. *CoRR, abs/1702.03859*. http://dblp.uni-trier.de/db/journals/corr/corr1702.html#SmithTHH17

Søgaard, A., Agic, Z., Alonso, H. M., Plank, B., Bohnet, B., & Johannsen, A. (2015). Inverted indexing for cross-lingual nlp. *ACL (1)*, 1713–1722. http://dblp.uni-trier.de/db/conf/acl/acl2015-1.html#SogaardAAPBJ15

Søgaard, A., Ruder, S., & Vulić, I. (2018). On the limitations of unsupervised bilingual dictionary induction. *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 778–788. https://doi.org/10.18653/v1/P18-1072

Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in nlp [cite arxiv:1906.02243Comment: In the 57th Annual Meeting of the Association for Computational Linguistics (ACL). Florence, Italy. July 2019]. http://arxiv.org/abs/1906.02243

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, & K. Q. Weinberger (Eds.), *Nips* (pp. 3104–3112). http://dblp.uni-trier.de/db/conf/nips/nips2014.html#SutskeverVL14

Tela, A., Woubie, A., & Hautamäki, V. (2020). Transferring monolingual model to low-resource language: The case of tigrinya. *CoRR, abs/2006.07698*. http://dblp.uni-trier.de/db/journals/corr/corr2006.html#abs-2006-07698

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems 30* (pp. 5998–6008). Curran Associates, Inc. https://papers.nips.cc/paper/7181-attention-is-all-you-need

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T. L., Gugger, S., . . . Rush, A. M. (2020). Transformers: State-of-the-art natural language processing. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 38–45. https://www.aclweb.org/anthology/2020.emnlp-demos.6

Wu, S., & Dredze, M. (2019). Beto, bentz, becas: The surprising cross-lingual effectiveness of bert. *CoRR, abs/1904.09077*. http://dblp.uni-trier.de/db/journals/corr/corr1904.html#abs-1904-09077

Wu, Y., Schuster, M., Chen, Z., Le, Q. V. [Quoc V], Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K. et al. (2016). Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.

Xiao, M., & Guo, Y. (2014). Distributed word representation learning for cross-lingual dependency parsing. *Proceedings of the Eighteenth Conference on Computational Natural Language Learning*, 119–129. https://doi.org/10.3115/v1/W14-1613

Xing, C., Wang, D., Liu, C., & Lin, Y. (2015). Normalized word embedding and orthogonal transform for bilingual word translation. In R. Mihalcea, J. Y. Chai, & A. Sarkar (Eds.), *Hlt-naacl* (pp. 1006–1011). The Association for Computational Linguistics. http://dblp.uni-trier.de/db/conf/naacl/naacl2015.html#XingWLL15

Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou, R., Siddhant, A., Barua, A., & Raffel, C. (2021). MT5: A massively multilingual pre-trained text-to-text transformer. *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 483–498. https://doi.org/10.18653/v1/2021.naacl-main.41

# Appendix A

# Fine-Tuning Analysis

We report here training and valuation accuracy on the TyDiQA-GoldP dataset after one epochs and after two epochs of fine-tuning, for both the baseline and the ICEBERT model.
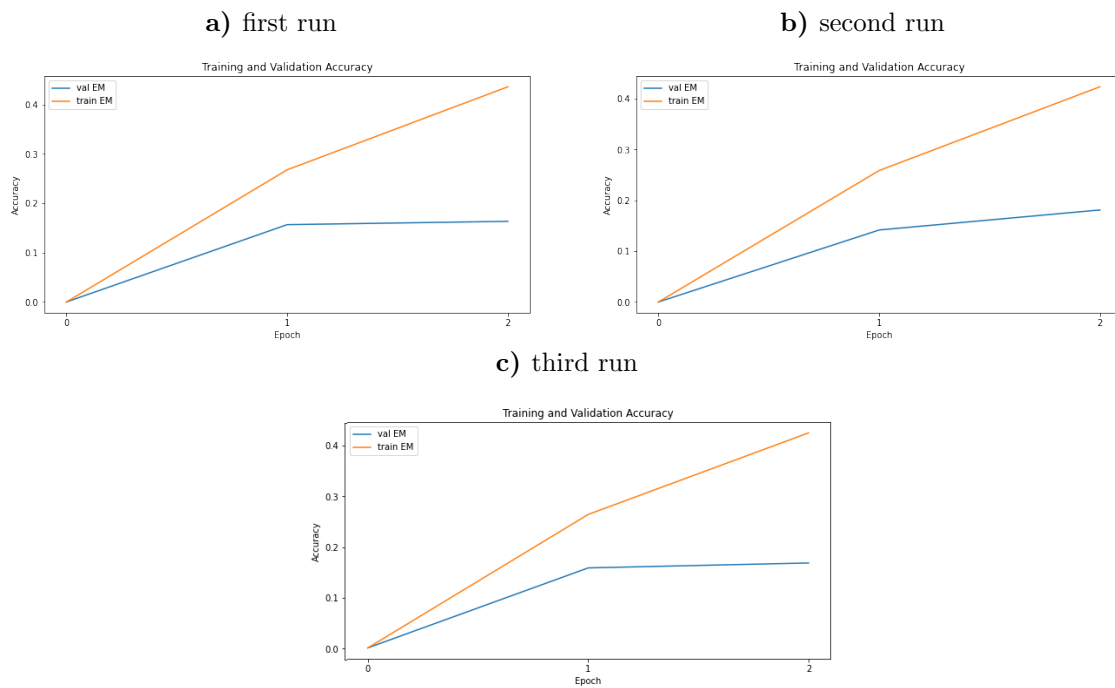
**a)** first run                                           **b)** second run



**c)** third run



Figure A.1: Training and validation accuracy during fine-tuning, baseline.

ICEBERT's accuracies are slightly higher than the baseline's ones, but no significant differences can be observed when comparing the curves.

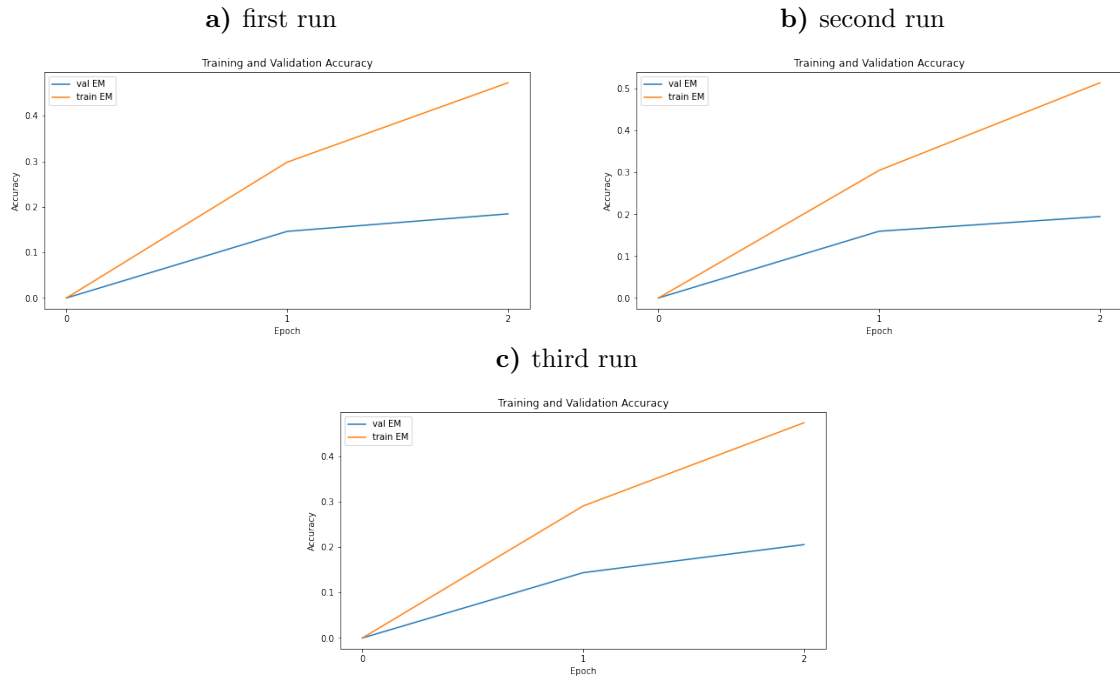**a)** first run

**b)** second run

**c)** third run



Figure A.2: Training and validation accuracy during fine-tuning, ICEBERT model.

It must be underlined, however, that while all ICEBERT's results were successful, two baseline fine-tuning yielded terrible results and had to be repeated. The plot of these fine-tunings, whose results were excluded from the computation of the average scores, are displayed in Figure A.3.

**a)** first failed run
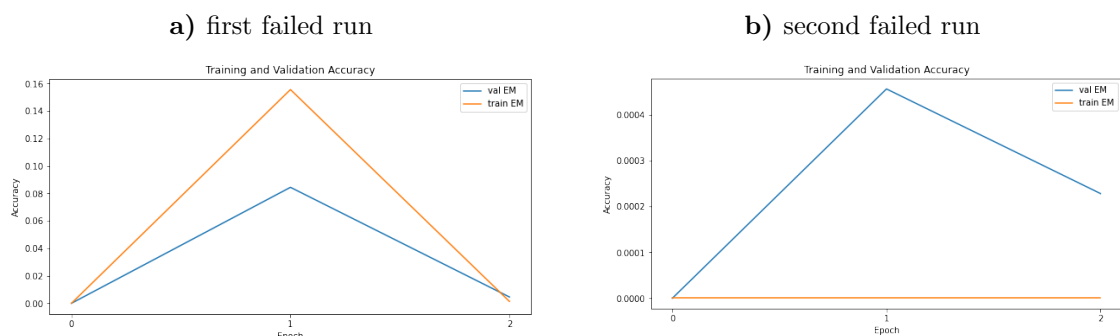
**b)** second failed run



Figure A.3: Unsuccessful baseline's fine-tunings.

The ICEBERT model, therefore, in addition to performing better than the baseline, also exhibits more stability in the fine-tuning.