

Ricardo Lukas Jung
6227492
Empirische Sprachwissenschaft (B.A.)
Phonetik & Digital Humanities
15th Semester
s2458588@stud.uni-frankfurt.de

Bachelor Thesis

Lexicalizing a BERT Tokenizer

**Building Open-End MLM for Morpho-Syntactically Similar
Languages**

Ricardo Lukas Jung

Date of Submission:
February 14, 2023

Text Technology Lab
Prof. Dr. Alexander Mehler
Dr. Zakharia Pourtskhvanidze

Erklärung

Hiermit bestätige ich, dass ich die vorliegende Arbeit selbstständig verfasst habe und keine anderen Quellen oder Hilfsmittel als die in dieser Arbeit angegebenen verwendet habe.

Ort, Datum

Unterschrift

I thank my supervisors Prof. Dr. Alexander Mehler and Dr. Zakharia Pourtskhvanidze for their competence and professionalism. Their support made this thesis across university departments possible. My explicit thanks go to Manuel Stoeckel and Giuseppe Abrami for their patience, advice and encouragement. Thanks to my close ones for their trust and support.

Abstract

This is the abstract: what is this about? what was done? what were the results?

Contents

List of Figures	I
List of Tables	II
List of Acronyms	III
1 Introduction	1
1.1 Motivation	1
1.2 Hypotheses	2
1.3 Scope and Structure	3
2 Overview	4
2.1 Related Works	4
2.2 Target Languages	5
3 Methodology	9
3.1 Requirements	9
3.1.1 Pipeline Structure & Transformers Library	9
3.1.2 Data	11
3.2 Implementation	11
3.2.1 Tokenizer	12
3.2.2 Model Training	18
3.2.3 Benchmark	19
4 Results	20
4.1 Benchmark	20
4.2 Wordmapping	23
5 Discussion	26
6 Conclusion	30
Bibliography	33

List of Figures

1	Schematic integration of Wordmap segmentation into the BERT architecture.	10
2	Segmenter output for 'verstehen'	17
3	Wordmap for 'viertelt'	23
4	Wordmap for 'anschauen'	24
5	Pareto distribution in POS vocabulary	25

List of Tables

1	List of similar languages	6
2	Example Wordmaps for <i>target</i> = verstehen	15
3	ID references assigned to full model names	18
4	List of all used models	19
5	Metrics for model mwo5	20
6	Metrics for model mwg5	21
7	Metrics for model mso5	21
8	Metrics for model msg5	22
9	Metrics for model bbgc	22
10	Test score summary for all evaluated models.	22

List of Acronyms

bbgc bert-base-german-cased

BERT Bidirectional Encoders from Transformers

BPE byte pair encoding

WP WordPiece

CL computational linguistics

GerParCor German Parliamentary Corpus

GPT Generative Pretrained Transformers

HanTa Hanover Tagger

LM language model

LSTM Long Short-Term Memory

ML machine learning

MLM masked language model

NLP Natural Language Processing

POS Part of Speech

TTLab Text Technology Lab

1 Introduction

This thesis investigates the use of specific intervention in tokenization subsystems of language modeling in machine learning. The intent of this thesis is to inject linguistic bias into the machine learning framework of Bidirectional Encoders from Transformers (BERT) to sharpen the analytical capacities of a masked language model (MLM). This introduction includes the research motivation and intentions, the resulting hypotheses and the scope of the thesis.

1.1 Motivation

There is a long-standing ambition in the computational linguistics (CL) community to understand natural language. Research in the past decades shows promising use of frequentist and statistical methods such as Generative Pretrained Transformers (GPT) (Radford et al. 2018) and BERT (Vaswani et al. 2017), leading to emergence of the powerful, problem solving machine learning (ML) models. It became apparent that these ML models are the state-of-the-art approach to an automated understanding of natural language. The structural parallels of machine learning to human learning have often been drawn (Jianlong Zhou (2018); Kühl et al. (2020)) to demonstrate their similarities and differences. An essential feature of ML (as opposed to human learning) is the possibility of actively controlling the learning parameters in a supervised environment. To test the efficiency of ML applications, a variety of tasks (e.g. classification, clustering, detection) are designed and applied. A trained model will yield performance scores based on the quality of its training, much like humans in language tests. Yet, statistical language modeling is not the first variant of language modeling. Traditional linguistics has accumulated extensive fundamental research prior to the discovery of ML architectures and their implementation. While generic ML frameworks seem appealing in the presumption that they require less work to reach fairly acceptable results, their performance is still unsatisfactory. The integration of aforementioned traditional linguistic knowledge into the learning processes for machine learning is the underlying motivation of this thesis.

While the research on neurophysical intake and representation of language is not entirely conclusive (Delogu, Brouwer, and Crocker (2019, pp. 1–3); Kimppa et al. (2019)), it is assumed that language learners typically build up a lexicon consisting of lexemes (Brennan

2022, pp. 82 – 100), which they have to analyze accurately in order to be proficient in that target language. A language model relies on a tokenizer to create such a vocabulary (attempts that bypass the vocabulary mechanism exist as well, e.g. Islam et al. (2022)). String-based tokenizers generally segment tokens into subwords (if possible) and provide a vocabulary comprising all the components needed to analyze a given string. Ideally, those subwords are part of the functional vocabulary in the target language, which would make them a morpheme. A morpheme is canonically defined as the smallest unit carrying meaning in a language (Colman 2009). Capturing lexical morphemes at the tokenization level is a problem of morphological analysis. A standard approach implemented in the BERT architecture is WordPiece tokenization, which takes pairs of n-grams and calculates a score based on the overall frequency in the text as well as cooccurrences. This idea was pioneered by Schuster and Nakajima (2012) in order to deal with the problem of near infinite vocabularies and was later adopted in a Google research project designing the final WordPiece algorithm (Wu et al. 2016). WordPiece is designed to find the most effective pairs of subwords in a corpus to optimize vocabulary size. Subwords, however, are not necessarily a productive component of inventory used in inflectional morphology. In practical use, the morphemes of a language and its generated tokenizer vocabulary rarely totally coincide. Following a revised principle of the interaction hypothesis that - put simply - *input quality is crucial for output quality* (Patsy M. Lightbown 2013, p. 114), the morpheme vocabulary is identified as the point of leverage in the upcoming chapter. Underlying is the notion not to just “push F scores”, but to explore a method of morphemic tokenization.

1.2 Hypotheses

The following research questions are formulated:

HYP1: Targeting a Part-of-Speech class during tokenization affects model performance.

This is tested through designing a tokenization process with POS overhead, fine-tuning a model on that tokenization and running a sequence classification task.

HYP2: Separating lexical from functional morphemes improves segmentation in lexemic tokens.

The design for hypothesis 1 includes creating a separate POS vocabulary. Model scores and qualitative evaluation will determine the plausibility of hypothesis 2.

1.3 Scope and Structure

The upcoming chapter 2 provides an overview of the current state of morphological tokenization to outline the research domain. German is declared as target language, verbs being the classified tokens of interest. The language is described paying special attention to its morphological complexity and peer languages. This serves as preface to the methodology, connecting characteristically matching languages to form a pool of possible target languages.

As the main part of this thesis, chapter 3 presents the methodology for a new tokenization method. It is sectioned into a theoretical part and the practical implementation of the process. The theoretical part focuses on the architecture of its components and the value they hold towards lexicalizing a tokenizer. In the implementation, each part of the tokenization process is explained, providing algorithms and their integration into the machine learning framework. Model training and testing are described. Results are presented in chapter 4, followed by a discussion of possible effects in chapter 5. The conclusion in chapter 6 puts the benchmark results and tokenization method into perspective and .

2 Overview

Morphological tokenization can be defined as the process of identifying segments in text that are productive in a given language, carrying meaning and hence fitting the definition of a morpheme. Recent or relevant findings on morphologically pretrained models in machine learning literature are discussed in this chapter.

2.1 Related Works

This thesis is mainly inspired through a paper by Hofmann, Schuetze, and Pierrehumbert, proposing more flexible segmentation with respect to preserving morphological structure (Hofmann, Schuetze, and Pierrehumbert 2022). In the past, many efforts towards morphological tokenization have been made. Earlier generalized attempts like morfessor (Creutz and Lagus 2002) have been outperformed by sequence based models that also use linguistic morphology (Peters and Martins 2022). Notably, the top-down generation of subword vocabularies has shown promising results for tokenization in fusional languages. This aligns with the notion that standard byte pair encoding (BPE) (Sennrich, Haddow, and Birch 2016) or WordPiece (Wu et al. 2016) tokenization effectivity suffers from complex morphology causing a big vocabulary.

The overall comparison (Peters and Martins 2022, p. 134) shows an increase in performance for languages of similar morphological complexity. Good benchmarks are reached for the agglutinative fusional peers, e.g. Italian, Latin, Spanish, Russian. However, the performance of this type of tokenization decreases with English. This effect could be attributed to the decline of English morphology in the last centuries. Toraman et al. find that the vocabulary size plays a special role in morphological tokenization and even define a vocabulary size ratio between 20% to 40% of the number of model parameters depending on the type of tokenizer (2022, pp. 11–12). Since tokenization and vocabulary are obviously interdependent, the vast amount of typological variety seen in languages raises the question: is there a right way of tokenizing? This issue is addressed by Rust et al., where a mid-scale investigation was done to see whether different languages actually need more specific tokenizers compared to generalized tokenizers. They report an improvement of model accuracy and F score across all tasks and languages (Rust et al. 2020). While this suggests to always consider choosing a

method tailored for single languages in Natural Language Processing (NLP), the answer to the problem of performance versus maintenance in models might be more straightforward than treating every language singularly. Every language is undoubtedly unique, however, that does not rule out simplification by means of further classifying and grouping target languages. In an effort to explain the provenience and relatedness of languages, many tools in the domain of typology, NLP and indo-european studies have been constructed.

Whether it be identification of morphological features (Comrie 1989, pp. 42–56), complexity measures (Çöltekin and Rama 2022) or connection through reconstruction (Bouckaert et al. 2012), the different linguistic disciplines suggest observable regularities by which to morphologically group languages. Leveraging the relatedness of languages in NLP is not a new idea in tokenization or Part of Speech (POS)-tagging, but is showing mixed results up to this day, even when using augmentation methods (Aepli and Sennrich 2021). The options seem to branch out quickly, but the mechanism of clearly separating lexemic and functional information seems inherent to most languages. They differ in how they combine grammatical functions in morphemes (fusion) and bind them to lexemic morphemes (synthesis). This seems to be the reason why stemming, lemmatization or other morphological analyses are very relevant to building good tokenizers for all languages on the isolating to synthetic spectrum (Schwartz et al. 2020, pp. 51–53).

2.2 Target Languages

This section identifies target languages that share common morphological features with German. It is assumed that, languages of the same morphological type will behave similarly when analyzed morphologically. German was selected to serve as the example language within the family of fusional languages. The aim is not to propose yet another case study of German, but to introduce German as a surrogate to further the scope of application on other languages of similar morphological complexity.

German (ISO639-3: deu) is a West Germanic language and the official language of Germany, Austria, Switzerland, Liechtenstein and Luxemburg (Glück and Rödel 2016). It is an inflectional synthetic language with approximately 130 million speakers¹. German is largely researched and is still paid much attention to in the domain of (computational) linguistics.

On the one side, linguistic typology has come up with many useful classifications for languages. On the other side, in the pursuit of reconstructing languages, Indo-European stud-

¹<https://de.statista.com/statistik/daten/studie/1119851/umfrage/deutschsprachige-menschen-weltweit/> Last accessed: 09.01.23

ies have established a widely accepted phylogenetic model of the diachronic dependency of Indo-European languages. Both disciplines contribute to language classifications that are used in this subsection.

Morphological complexity is a term to describe how languages use paradigms to connect grammatical information with lexemic information (Baerman, Brown, and Corbett 2017). Nevertheless, it needs to be taken into consideration that morphological complexity is a classifying category to describe gradients of function-to-morpheme correspondence and measure of morphematic agreement, not a qualitative assessment. The common denominators that make languages morphologically complex are their morphological features. The languages that use affixation, fusion, composition and derivation (among others) are all fit candidates compared to German. A summary of morphological typology is provided by Aikhenvald and Dixon (2017, pp. 78–93).

Within the scope of this thesis, German will be the exemplary target language for the experimental setup. Its morphological complexity can be compared to other related or non-related languages as shown in Table 1 (ISO identifiers provided at WALS.²). To this date, there is no universally accepted measure for the complexity of a language, but groupings exist on different parameters:

Table 1: Listing of languages similar to German given the type of similarity based off Lehman (2022) and Ehret et al. (2021).

Language	Similarity	ISO 639–3
Norwegian	Closely Related	isl
Danish	Closely Related	nor
Dutch	Closely Related	nld
English	Closely Related	eng
Icelandic	Closely Related	isl
Romanian	Morphology	ron
Spanish	Morphology	spa
Finnish	Morphology	fin
Italian	Morphology	ita
Hungarian	Morphology	hun

There have been interpolations between human judgement and statistical measures (on similarity of languages) which can be taken into consideration (Bentz et al. 2016). It should be emphasized, that while there is no definite proof of concept for tokenizations being effec-

²<https://wals.info/languoid>

tive when connecting target languages through morphological parameters, there is a strong suggestion in data and intuition of researchers that tokenization for morphologically similar languages should profit from these similarities.

With an arguable exception to English, the languages in Table 1 treat their lexemes with similar morphological processes. The upcoming interlinear glossings (as per Leipzig Glossing Rules³) provide examples for inflectional morphology in verbs within this group of languages. To outline word formation processes, a glossing from Nikanne (2017, p. 71) is considered:

- (1) Tytöt istu-i-vat tuolilla
 girls sit-PST-3PL chair.ADE
 ‘The girls sat on the chair.’

This Finnish sentence is a textbook example of agglutinative inflectional morphology. The verb {istu} is inflected by suffixing two morphemes marking the past tense {i} and the third person plural {vat} (curled brackets denote morpheme boundaries). In this case every morpheme expresses one grammatical function, apart from {istu} which contains the lexical information for the verb “to sit”. The functional morphemes in example (1) follow the word to be inflected. With many more functional morphemes present in Finnish, Nikanne reports that there is an order in which inflectional morphemes usually appear. In consequence, analyzing Finnish verbs results in different but reoccurring patterns depending on the degree of inflection. The lexemic morpheme {istu} can be modified by {i} alone to just express past tense and still be productive. Following up on the idea of agglutination, Hungarian applies a slightly different strategy to achieve inflection:

- (2) Tegnap meg-hallgattunk egy lányt.
 yesterday PRF-listen.PST.1PL a girl.ACC
 ‘Yesterday we interviewed a girl.’

Hungarian is also classified as an agglutinative language for its frequent use of affixes. It is additionally known to combine several grammatical functions into one morpheme as can be seen in example (2) as given by Kiss and Hegedus (2021, p. 262), making it a hybrid of agglutinative and fusional. The analysis of the verb in (2) does not allow canonical segmentation to the stem although there is an underlying form {hall;*} meaning “hear”. Instead, it exists in inflection paradigms like the given example {hallgatunk} combining the grammatical

³<https://www.eva.mpg.de/lingua/pdf/Glossing-Rules.pdf>

categories past tense and first person plural. The morphemes addressed so far where either suffixed or not segmentable. As lexical part {hallgatunk} receives a prefix {meg} expressing perfect tense. This use of morphemes can also be seen in Germanic or Romance languages, like Italian (adapted from Iacobini and Masini (2005, p. 163)):

- (3) far=se=la sotto
do-REFL.PRT-PRON.PRT under
‘To quake in one's boots’

Here the {se} and {la} carry two functions each and are suffixed to {far}, showing that there are morphological types in between agglutinating and fusioning. Arguably, {se} being a clitic pronoun that will appear in different positions acting as indirect object, but never independent of the verb.

After a partial look on the classified languages, two important empirical descriptive caveats remain: there are exceptions to almost every regularity in languages. No language is entirely consistent in following a morphosyntactical paradigm, meaning no language is entirely fusional or agglutinative (same applies to the synthetic and isolating spectrum). Judging from the word shapes in the data and literature, the way languages modify their stems or lexemic morphemes is largely based on affixation. In a tokenizer acknowledging lexemic parts of words, the knowledge of word formation in the target language should be conveyed.

3 Methodology

This chapter first outlines the theoretical framework as well as required input data. Subsequently, the implementation is explained with references to methodological implications of the choices made.

3.1 Requirements

A series of tools will help to achieve lexicalized tokenization. They will be described in this section along with their methodological edge.

3.1.1 Pipeline Structure & Transformers Library

BERT is a language learning transformer model designed for NLP tasks (Vaswani et al. 2017). Upon release, it achieved higher performance scores compared to previously used Long Short-Term Memory (LSTM) models (Devlin et al. 2018). Two main model characteristics can be observed for BERT: it is the first language model (LM) to implement simultaneous attention heads, allowing for bidirectional reading. The methodological implication of reading to the left and right of a token is to include more information about the language in single embeddings. Secondly, BERT introduced the still relevant MLM method for training. The method involves masking a specified amount (default 15%) of random tokens in the input sequence. Masked tokens are guessed by the model, which can then update its weights according to success or failure. Bert-base-german-case is a BERT-base model with 12 encoding layers, 768 hidden layers, 12 heads and 110M parameters¹.

The NLP community has since developed BERT and adapted it to the needs of contemporary NLP problems (e.g. roberta, germanbert, mbert). Its extensive support, comparability and versatility make BERT the model of choice for this thesis. Another notable feature in BERT is the implementation of the WordPiece tokenizer module². Default BERT WordPiece tokenization is predominantly heuristic by combining strings based on a precalculated score. A variety of pre-trained tokenizers are available, although they come with a caveat: once a

¹https://huggingface.co/transformers/v2.9.1/pretrained_models.html

²The source code was not published but is inferred on huggingface

tokenizer is trained on a dataset, it is specific to that dataset. This means that the application of a tokenizer on another dataset may result in out-of-vocabulary issues and different token/subtoken distributions.

Particularly relevant to this thesis is the option to train an individual tokenizer from the base module, which the bert-base-german-cased (bbgc) model implemented with German data. Usually, WordPiece generates its own set of subtokens called *vocabulary*, which is used by the standard tokenizer to generate unique numerical input IDs. These identifiers correspond to tokens or parts of tokens found in the original dataset that the model was trained on. Wordmap partly takes over the tokenization depending on the the type of input it receives, the exact process is given in subsection 3.2.1. Once a string is tokenized, it is passed on to the transformer model for contextualized embedding. Figure 1 shows the processing of data up until token encoding.

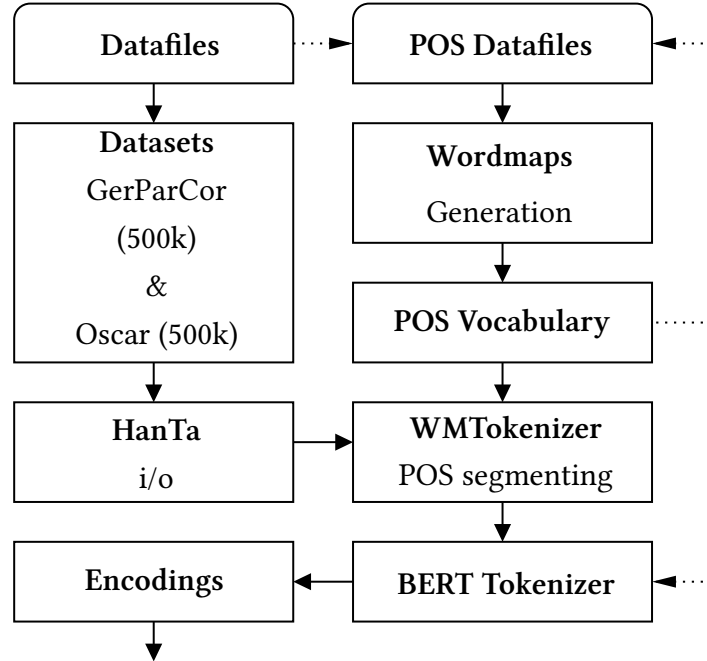


Figure 1: Schematic integration of Wordmap segmentation into the BERT architecture. Continuous lines show non-selective channels, while dotted lines only pass data selectively

In Figure 1, **Datafiles** contains the raw unprocessed data coming from corpora. The input is provided in txt-files containing sentences line by line. Two datasets are generated from this input, each amounting to 500k sentences per corpus for fine-tuning. **POS Datafiles** ideally contains only words that are identified as the same POS category. To increase the resulting

POS vocabulary, POS Datafiles can be augmented by two channels: input can come from the corpus itself, but would need to be tagged first. Alternatively, a predefined list of words can be used as an external and generic source for Wordmap generation. In POS Vocabulary, adjustments like removing outliers and interferences are performed on the set of subwords yielded by Wordmap generation. It is important to note that the POS vocabulary is separate from the the tokenizer vocabulary. This serves the specific purpose of keeping the vocabulary on which the model is trained clean from unwanted or unused subwords. Machine learning practice generally points out the trade-off between vocabulary size and model performance, hence the addition of POS Vocabulary to the BERT Tokenizer vocabulary is done only when necessary. Datasets are then passed through Hanover Tagger (HanTa) to provide WMTOKENIZER with flagged data to tokenize only selected tokens. Lastly, BERT Tokenizer returns the current standard encodings used in BERT language modeling.

3.1.2 Data

Two corpora were selected for fine-tuning: German Parliamentary Corpus (GerParCor) and Oscar (OSCAR 2022). FLOTA was trained on 12k samples per category (Hofmann, Schuetze, and Pierrehumbert 2022), which is why a sample size of 500k is deemed sufficient for this for this fine-tuning. GerParCor is a “genre-specific corpus of (predominantly historical) German-language parliamentary protocols from three centuries and four countries, including state and federal level data.” (Abrami et al. 2022, p. 1). Of all subcorpora included in GerParCor, one specifically includes texts and transcripts of the German Parliament. All 500k samples in the GerParCor dataset used for this thesis are found in the Bundestag subcorpus of GerParCor. Oscar offers other subcorpora partly hosted on huggingface. In this case, a web-crawled German corpus called unshuffled_deduplicated_de was chosen, as introduced by its curators (Ortiz Su’arez, Sagot, and Romary (2019); (Ortiz Su’arez, Romary, and Sagot (2020)). The POS vocabulary put through the Wordmapping pipelines comes from Wiktionary (2022) and HanTa-crawled verbs of GerParCor.

3.2 Implementation

The implementation of aforementioned methodology is presented in this section. The core part is the conceptualization and implementation of the tokenizer, which is explained in three steps: the Wordmap algorithm generating a vocabulary, the segmentation algorithm for token segmentation and maximization, and the integration into BERT. After that, a short description of model hyperparameterization and training is given. The same is done for the

model evaluation task in the last subsection of this chapter.

3.2.1 Tokenizer

Being the main part of the thesis, the methods of tokenization are explained more in detail in the upcoming subsection. It contains two algorithms - one for vocabulary generation and one for token segmentation. Throughout this section, the words *target* and *token* are used to describe a word that is analyzed. One denotes the argument of the wordmapping and segmenter function (*target*), the other describes a word occurring in the corpus (*token*).

Generating a custom pre-training vocabulary

Embedding those subwords which take part in inflectional processes essentially means deriving sensible subwords to represent actual morphemes. The vocabulary generated for this chapter is a list of inflected and non-inflected verbs (as to state the example) provided by two sources: (1) the german verb wiktionary and (2) a crawl of the GerParCor corpus with the HanTa tagger Wartena 2019. Initial experiments were done only with the wiktionary verb list since it provided sanitized input for the algorithm in development. As soon as the verbs were extracted from GerParCor via HanTa they were added to increase the sample size at the cost of adding noise. A naive regular expression detecting syllable nuclei was used to sort the verbs into their respective syllable count. It provides an overview over the POS dataset and can potentially decrease iterations. The whole POS dataset was iterated over nevertheless, words with syllable count 2 or less are only analyzed to the right to prevent the algorithm from analyzing words that are likely to not have a prefix. Through the positional distribution and morphological variety mentioned in section 2.2 this did not pose a problem.

The Wordmap algorithm as shown in Algorithm 1 is the first step to extracting morphemes from a token. Its purpose is to compare two strings and store their intersections in a map of boolean values.

Wordmap requires two **inputs** *verbs* and *target*. The resulting wordmap will be generated for *target*, while *verbs* serves as comparison. Any map generated from this also has the same length as *target*. *verbs* is a set of tokens pertaining to the same POS category. Note that *verbs* should only contain those POS-tagged tokens that are expected to carry lexical information (e.g. verbs, adjectives, etc.). The set is previously extracted from the corpus by POS-tagging. Optionally, the set can be augmented by manually adding POS matching tokens from external sources. The 2-tuple *pair* are the strings to be compared. It is passed on to (1) `SHORTER`, a function returning the shorter of both strings (2) `LONGER`, expectedly returning the longer of both strings (3) `MATCH_CASE`, a function to determine the behavior

Algorithm 1 Wordmap generation

Input: $verbs = \{v : POS\}, target$ $\triangleright verbs$: set of single-POS lexemic tokens**Output:** $maps = (map_1, \dots, map_{|verbs|})$

```
1: function WORDMAP( $w1, w2, d=0$ )
2:    $f_1 : c1, c2 \mapsto c1 == c2$ 
3:    $f_2 : \left( \sum_{i=0+d}^{|w1|} w1[i], w2[i] \mapsto f_1(w1[i], w2[i]) \right)$ 
4:   return  $f_2$ 
5: end function
6:
7: for  $v \in verbs$  do
8:    $pair = (target, v)$ 
9:    $s = \text{SHORTER}(pair)$ 
10:   $l = \text{LONGER}(pair)$ 
11:   $case = \text{MATCH\_ENDS}(pair)$   $\triangleright$  Returns if strings match in the last or first position
12:   $len = \text{LEN}(l)$ 
13:   $\delta = \Delta(len, \text{LEN}(s))$ 
14:
15:  if  $case$ : any match then
16:    if  $\delta$  then
17:      if  $case$ : left match then
18:        WORDMAP( $l, s$ )
19:        Pad map from right side with 0s to match  $len$ 
20:      end if
21:      if  $case$ : right match then
22:        WORDMAP( $l, s, \delta$ )
23:        Pad map from left side with 0s to match  $len$ 
24:      end if
25:    else
26:      WORDMAP( $l, s$ )
27:    end if
28:  end if
29: end for
```

of the algorithm later on. As two strings are compared `MATCH_CASE` captures three cases: *pair* matches in the first, last or both positions. Finally *len* denotes the length of the longest string and δ difference in length between *s* and *l*. δ functions as an offset for index-based comparisons `WORDMAP`, the mapping function generating the Wordmaps.

Once every *v* has been compared to *target*, *maps* stores boolean counts of characters occurring in their respective positions in *target*. Every map is cleaned with a regular expression to reduce noise caused by natural character occurrence (some characters like `<n>`, `<s>` will be more frequent than others). Omitting this step will result cause a ceiling effect rendering the maps useless. Continuous concatenations of leading or trailing matches stay, while every match enclosed by 0 will be replaced the 0. As an example, *wordmap* = 11101100101 contains three matches on the inside which will result in 11100000001 as final output. In the penultimate stage of mapping *target*, all maps are summed up to receive the number of absolute positional occurrences of every character in *target*. This positional mapping allows for detecting relevant segments in a token based on a threshold. Characters in range of the predefined threshold are selected for the mapping of a target token. Functional morphemes (morphemes that are carriers of grammatical features) are typically much more frequent than their lexical counterparts. Consequently, lexical morphemes in the family of inflectional languages are - by definition - modified by functional morphemes, they occur much less frequently. In this case, the activation function for a concatenation of same boolean values to be selected as segment is the normalizing z-score function defined as: $z_i = \frac{x_i - \bar{x}}{S}$, where *z* is the z-score, \bar{x} and *S* are the sample mean and the sample standard deviation. Any segment below a z-score of zero is detected. For every verb that is mapped, the detected affixes are added to a vocabulary of functional morphemes. The rest of the string is saved into the vocabulary of lexemic segments, assuming that it still contains the lexical information in the verb.

Table 2: Example Wordmaps for *target* = verstehen

String	Wordmap	Case	Padding
verarbeiten	111000000	left match	Yes
	001000011	right match	Yes
variiert	101001000	left match	Yes
vormachen	101000111	left match	No
	101000111	right match	No
anstrebttest		no match	
(...)			

Both vocabularies are then cleared of outliers in different ways. The functional vocabulary drops every string longer than 2.5 times the standard deviation of its own population. Since more variance in length is expected in the lexemic vocabulary, the outlier function is limited to drop everything above the length of 2 times the mean absolute deviation as a staple method for more volatile datasets. The upper segment of n-length strings in the lexemic vocabulary has thus been covered, but the low segment (short strings) remains untouched. To solve this problem, probably one of the more reaching but important measures is taken: removing any string of the lexemic vocabulary that is smaller than the mean length of strings in the functional vocabulary. Removing this interference in vocabularies is subject to the notion that functional morphemes will be shorter than lexemic morphemes. After minimizing outliers and interference, both vocabularies are joined again to form the segmenters vocabulary.

BERT Tokenizer Modification

The BERT tokenizer module tokenizes using subwords from its own vocabulary during pre-training. By loading the bbgc pre-trained tokenizer with and an instance of HanTa as tag-gint overhead a sample from the dataset batching function can be fully encoded for training. While iterating over the sample HanTa passes on detected verbs to Algorithm 2. All tokenized subsentences are then joined and added to the encoded batch.

Algorithm 2 recursively computes every possible segmentation for a string *target* from a given vocabulary *pos_vocab* from left to right. The vocabulary contains all the segments that were identified in the previous step by Algorithm 1. Every segmentation has to be complete so that its segments corresponds to non overlapping substrings of *target*. For every *target* a subvocabulary *morpheme* is defined, containing all strings that are in the vocabulary of POS! (POS!) members. This task is a weighted coverage problem in the NP-hard domain.

Algorithm 2 Target Segmentation

Input: $target, pos_vocab$ \triangleright Vocabulary consisting of verbs only**Output:** $\{(tuples\ of\ subwords)\} \approx \{t \in \mathcal{P}(s \in pos_vocab : s \in target) : t \equiv target\}$

```
1:
2:  $segmentations = ()$ 
3:
4: function SEGMENTER(token, stop, start=0, segments)
5:   if  $start = stop$  then
6:     Add  $segments$  to  $segmentations$ 
7:   else
8:      $morphemes = (m \in pos\_vocab : target.STARTSWITH(m) \wedge |m| > 1)$ 
9:     for  $m \in morphemes$  do
10:       $start += LEN(m)$ 
11:      Add  $m$  to  $segments$ 
12:       $rest = target[LEN(m) :]$ 
13:      if  $LEN(rest) == 1$  then
14:         $start = stop$ 
15:        Add  $rest$  to  $segments$ 
16:        Add  $segments$  to  $segmentations$ 
17:        Decrement  $start$ , crop  $segments$ 
18:      else
19:        SEGMENTER( $target = rest, stop, start = start, segments$ )
20:        Decrement  $start$ , crop  $segments$ 
21:      end if
22:    end for
23:  end if
24: end function
25: MAXIMIZE_SEGMENTS( $segmentations$ )
```

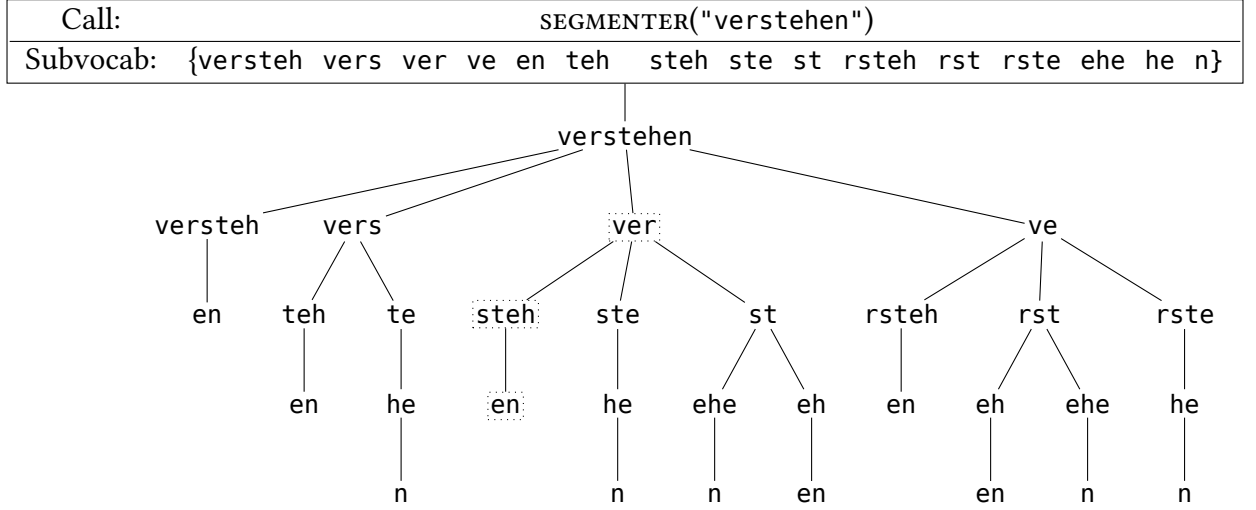


Figure 2: Permutations SEGMENTER generates from target i.e. verstehen. Dotted segments mark the segmentation chosen later by the maximization function during tokenization.

Unary morphemes to the left are excluded from the pre-selected subvocabulary to drastically reduce the number of possible permutations, as they can be embedded in n-ary tokens as well. The recursion can be seen as n-ary trees containing every permutation of the set *morphemes* where the sum of branches all satisfy *target*.

As shown in Figure 2, all morphemes that *target* starts with are stored to form the first nodes of the permutation tree. Each time a morpheme is selected the index *start* is incremented by the length of the morpheme to indicate when the string has been completely segmented. The new recursion is called with the updated index *start* and *target* sliced by the length of the morpheme contained in the parent node. Incomplete segmentations that miss exactly one character to the right are accepted with the added missing string. If the vocabulary cannot satisfy a segmentation because it is missing the necessary strings, segmentation is omitted and the original input token is returned as such.

Then, of all *segmentations* a single segmentation is selected by a maximization function MAXIMIZE_SEGMENTS calculating weights for every segment. The maximization function for one segment is the defined as:

$$\arg \max_{s \in S} f(s) := \left\{ s \in S : \sum_{i=1}^{|s|} \sqrt[|s|]{\frac{s_i}{t} \div |s|} \right\} \quad (3.1)$$

Where S is a set of segmentation tuples s , $|s|$ is the length of the tuple (read: number of

segments) and s_i indicates the length of the segment at position i in tuple s . For every segment (vertex) in a segmentation tuple (branch) the segment’s length is divided by the token length t to get the coverage the morpheme provides towards the target string. The number of segments $|s|$ is a divisor meant to cap the number of segments in a segmentation. It prevents choosing a segmentation overflowing with too many short segments. Short segments are convenient for completing a segmentation, but will increase the chance of slicing a token where it linguistically lacks sense to do so. Lastly the s_i^{th} root implements a bias against segments that are too long. While the accuracy decreases for longer words, this method of maximization performs reasonably well in the range of 1–4 syllable tokens, which make up around 87% of the verb vocabulary.

Ultimately, segmentation is a matter of interpretation. If there are several possible interpretations by which to segment a word this tokenization method relies on the assumption that the comparison with POS- members displays the probable shape of segmentation for a token. As mentioned in 3.2.2, the default WordPiece Tokenizer lacks a linguistically informed

The field of NLP (Glück and Rödel 2016) has been expanded ever since the emergence of the language models. Natural language processing is understood as the

3.2.2 Model Training

All fine-tuned models use the bbgc baseline by Chan, Schweter, and Möller (2019). For reference and readability they are given shortened IDs as follows:

Table 3: ID references assigned to full model names

ID	Full name
mwg5	mlm_wmt_gpc500k
mwo5	mlm_wmt_oscar500k
msg5	mlm_std_gpc500k
mso5	mlm_std_oscar500k
bbgc	bert-base-german-cased

The full model names in Table 3 captures broad characteristics of the model. For example, `mwo5` is the `m`asked language model using `W`ordmap tokenization fine-tuned on the `O`scar dataset of `5`00k samples, while `msg5` is the the masked language model using standard BERT wordpiece tokenization fine-tuned on the GerParCor dataset of 500k samples. All models except for `bbgc` are trained on the same parameters as seen in table Table 4.

Table 4: List of all used models and their hyperparameters. LR = learning rate, WU = warmup steps.

ID	Corpus	Tokenization	LR	Steps	Batchsize	WU	Base
mwg5	GerParCor	BWP + WM	0.0003	31250	16	500	bbgc
mwo5	Oscar	BWP + WM	0.0003	31250	16	500	bbgc
msg5	Gerparcor	BWP	0.0003	31250	16	500	bbgc
mso5	Oscar	BWP	0.0003	31250	16	500	bbgc
bbgc	Mixed	BWP	0.0001	810k/30k	1024	10000	-

The first four models use bbgc as baseline, meaning that they inherit its structural properties. DeepsetAI released bbgc training it in two phases (see column Steps in Table 4) with differing sequence lengths 128/512 respectively. The maximum sequence length for each model is 512, the base vocabulary size is 30k. Models mwg5, mwo5, msg5, and mso5 were trained with less steps and thus higher learning rate to converge quicker with the given datasets. The models were implemented in PyTorch (Paszke et al. (2019)) and trained on the GPUs Quadro RTX 8000 (48Gb) and NVIDIA GeForce GTX 1080 Ti (11Gb) provided at Text Technology Lab (TTLab).

3.2.3 Benchmark

The original draft for this thesis featured a benchmark with a translated multiple choice question answering task from the oLMpics benchmark (Talmor et al. 2020). Due to incompatible versioning dependencies this was not feasible and a substitute had to be found. To test the performance of trained models a one-shot sequence classification task is set up instead. The task consists of ~11k samples containing a sentence and a single label. Each sentence is the title of a news article found on the german Wikinews site, belonging to one of the categories listed on their topic section³. Samples are modified so that one sentence is assigned a single category label. Only those categories which appeared at least 200 times were selected for this task, leaving a total of 24 labels to choose from. Every model is then trained and validated on a split (train: 70%, test: 23%, validation: 7%) with a training and evaluation batch size of 16, a less sharp learning rate at 0.00002, and three epochs.

³<https://de.wikinews.org/wiki/Kategorie:Themenportal>

4 Results

4.1 Benchmark

This section is a comparative description of the results from the sequence classification task described in subsection 3.2.3. For every benchmark of three epochs precision, recall and F1 are given. Precision reflects the amount of correct predictions a model has made for that particular class. Recall shows how often the model correctly predicts a class in relation to all positive predictions. F1 is a score compromising precision and recall through the *harmonic mean* of precision and recall, measuring the models accuracy.

There is a common trend in all runs displaying a growth in precision, recall and F1. This means that all models have continuously improved predicting the class of a news title. The highest F1 in the last epoch achieved for any of the smaller models (mso5, msg5, mwo5, mwg5) is mwo5 with a score of 0.69194 (Table 5). Meanwhile, the lowest F1 is found in Table 6 at 0.590542 mwg5. A similar arrangement is found when comparing final test scores in the set of smaller models: mwo5 > mso5 > msg5 > mwg5 ($0.442827 > 0.405168 > 0.392490 > 0.389116$).

Table 5: Metrics for masked language model trained on the Oscar dataset with Wordmap infused tokenization. Evaluated on sequence classification task.

mwo5	Epoch 1	Epoch 2	Epoch 3	Test score
Precision	0.292614	0.446338	0.71387	0.449735
Recall	0.329531	0.552598	0.73384	0.474525
F1	0.242739	0.473851	0.69194	0.442827

Of all small models, mwo5 has the best results for all three metrics precision, recall and F1. It seemingly learns the fastest out of the four, starting with a recall of 0.329521 and ending on 0.73384 on the last epoch, overtaking its standard counterpart mso5 at the second epoch.

Table 6: Metrics for masked language model trained on the GerParCor dataset with Wordmap infused tokenization. Evaluated on sequence classification task.

msg5	Epoch 1	Epoch 2	Epoch 3	Test score
Precision	0.237664	0.399781	0.603534	0.441891
Recall	0.244613	0.463878	0.637516	0.440304
F1	0.163024	0.389905	0.590542	0.389116

Contrary to the pattern of mwo5, the wordmap model trained on GerParCor data predicts less accurately than the standard msg5 model. The initial epoch shows the lowest scores of all models for precision, recall and logically F1. None of the predictions keep up with the accuracy of the other small model at any point in training.

Table 7: Metrics for masked language model trained on the GerParCor dataset with bbgc tokenization. Evaluated on sequence classification task.

mso5	Epoch 1	Epoch 2	Epoch 3	Test scores
Precision	0.269615	0.422096	0.596987	0.395879
Recall	0.351077	0.501901	0.657795	0.446388
F1	0.266260	0.412598	0.604824	0.405168

With an F1 of 0.604824 at the final epoch mso5 is the lowest scoring standard model, despite starting with similar scores in epoch 1 (Table 7 vs. Table 8). The standard F1 test scores are close to each other, deviating only by one second decimal point indicating comparable benchmark performance. A comparison with the scores of the other Oscar-trained model mwo5 shows a bigger difference in accuracy, favoring the Wordmap version. The msg5 iteration is has a slightly steeper adaptation to the task, but fails to carry its performance over to the test scores:

Table 8: Metrics for masked language model trained on the GerParCor dataset with bbgc tokenization. Evaluated on sequence classification task.

msg5	Epoch 1	Epoch 2	Epoch 3	Test scores
Precision	0.297466	0.517110	0.656808	0.439873
Recall	0.359949	0.544994	0.676806	0.439924
F1	0.267111	0.480420	0.626593	0.392490

As can be seen, msg5 converges more quickly than its oscar standard counterpart (Table 8). Although it has better precision, recall and F1 throughout all epochs final its test F1 is minimally inferior to mso5.

Table 9: Metrics for masked language model baseline bert-base-german-cased. Evaluated on sequence classification task.

bbgc	Epoch 1	Epoch 2	Epoch 3	Test scores
Precision	0.646150	0.768675	0.860180	0.622436
Recall	0.709759	0.804816	0.883397	0.637262
F1	0.660588	0.778371	0.868166	0.624789

Of all models tested, bert-base-german-cased achieved the highest scores in all categories (Table 9). Its first epoch scores are as high as most of the smaller models final epoch. It is the only model that manages to achieve test scores above 0.5, but unlike the smaller models also performs below its first epoch scores.

Table 10: Test score summary for all evaluated models.

Summary	bbgc	mso5	msg5	wmo5	wmg5
Precision	0.622436	0.395879	0.439873	0.441891	0.449735
Recall	0.637262	0.446388	0.439924	0.440304	0.474525
F1	0.624789	0.405168	0.392490	0.389116	0.442827

Table 10 summarizes the final test scores for all evaluated models. All small model scores are relatively close to each other, but ultimately outperformed by their baseline model for around 0.2 points. Wmg5 comes closest to the baseline model performance at 0.442827 with a difference of 0.181962. It beats the smaller model test scores by a small margin of 0.4.

4.2 Wordmapping

The Wordmap algorithm calculated segment thresholds for each given verb and split the segments into two vocabularies, functional and lexemic. Sample results for mappings of tokens 'viertelt' and 'anschauen' are found in Figure 3 and Figure 4.:

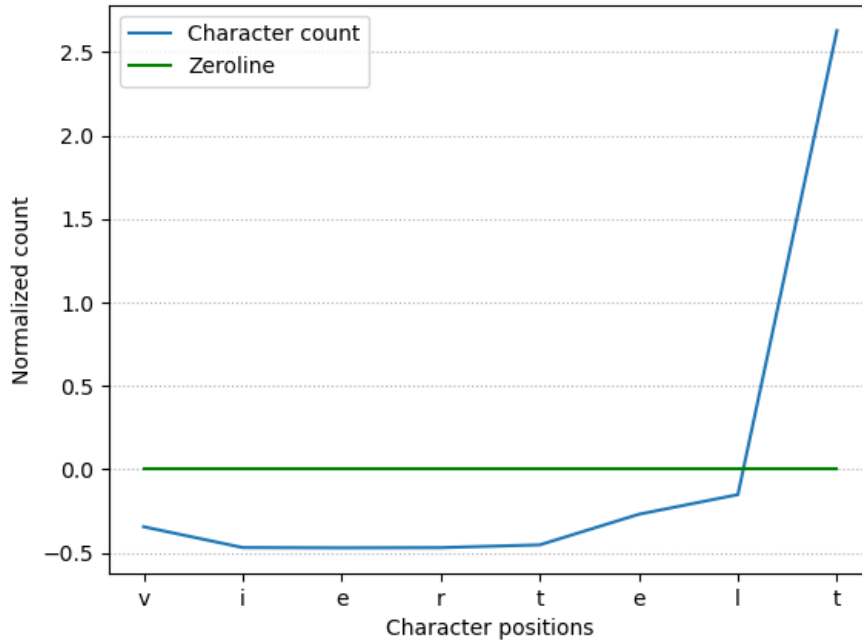


Figure 3: Wordmap thresholds for token 'viertelt'

Figure 3 shows the normalized count of each character relative to the POS corpus. The Wordmap algorithm interprets 'viertelt' as two separate subwords 'viertel' and 't', which correspond to the canonical segmentation of {viertel} and {t}. There is a slight fluctuation in lexemic segment ranging between -0.5 and -0.1 and a strong increase for the functional segment above 2.5 times the standard deviation of counts.

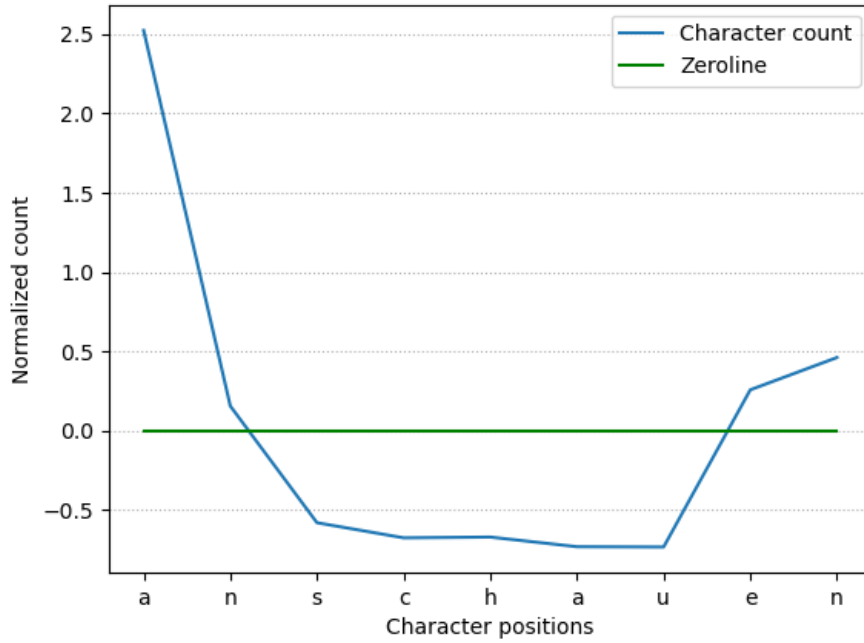


Figure 4: Wordmap thresholds for token 'anschauen'

The result for the 3-syllabic token 'anschauen' is similar, again showing the lexemic part way under the mean count (more than -0.5 standard deviations). Its prefix 'an' and suffix 'en' are above zeroline and are thus added to the POS vocabulary. A difference in affixes can be seen in this example, as the {an} combines counts for the morpheme and the single character <a> at the same time. This part of character bias mentioned in subsection 3.2.1 is not captured by the regular expression for character noise reduction. For every Wordmap, initial and final characters in tokens receive an overall count bias because only those matches are computed in Algorithm 1. Finally, functional vocabulary holds 138 types of affixes, including subwords, adverbials, prepositions and morphemes found in the German inflection paradigm. As expected, the lexemic vocabulary contains many more types (over 40k) than the functional one. Many stems are found multiple times with either subwords or morphemes attached to them. After combining both vocabularies, a quick search yielded most of the common atomic (non-segmentable) stems like e.g. {schlaf} and even their changed equivalent {schliefe}. The vocabulary was iterated over n^2 times with different versions of the same verb occurring many times, leading to duplicates. Ignoring duplicates and duplicate frequency grants a look at the distribution of frequencies in the POS dataset:

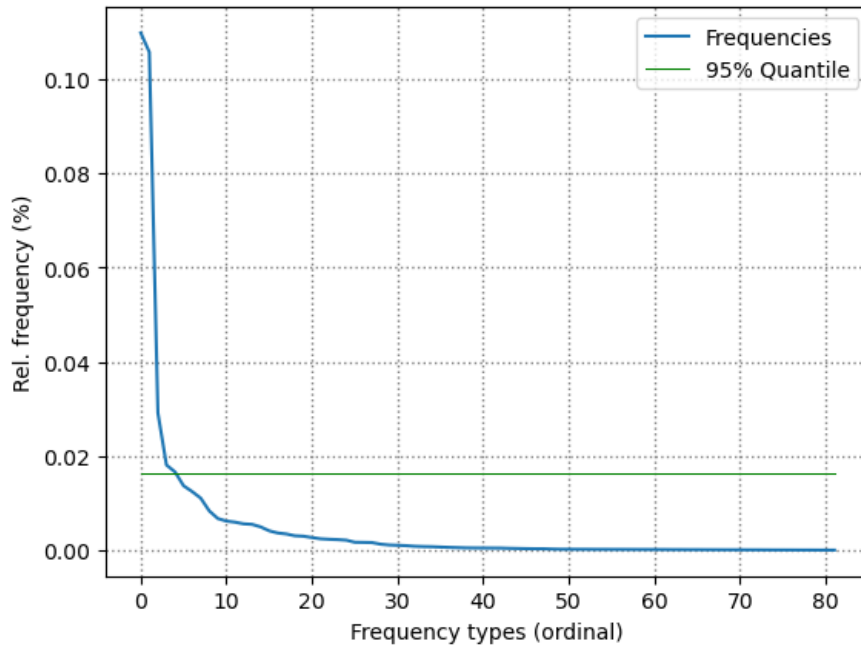


Figure 5: Frequency values mapped to unique frequencies

Plotting the relative frequency value to the set of unique frequencies in the POS vocabulary shows a pareto-like distribution (Figure 5), hinting at Zipf's law applied to the behavior of frequency in natural language. The 95% quantile visualizes the amount of non-frequent sub-words contained in the vocabulary - this becomes extremely distorted when plotting tokens to their frequency.

5 Discussion

Results of the pre-trained model benchmarks and tokenization are analyzed in this chapter. A closer look is taken at the training procedure and the benchmark results to identify possible effects in data or model architecture. For purposes of readability the models are referred to by their IDs introduced in subsection 3.2.2.

The first apparent difference when comparing the models presented in this study are the test scores. An overall sub-par performance can be observed for the non-ggbc models, while their baseline model clearly performed better at the task. There are many causes for low model performance. A model runs through many direct (hyperparameters, implementation) or indirect (data quality, task choice) modeling steps. For the course of this thesis a smaller set of experiment parameters were chosen to establish a controlled environment for hypothesis testing. Relating to the models this includes the following design measures: (1) smaller model and sample size, (2) constructing Wordmap tokenization and using it on singular tokens during batching, and (3) a single evaluation task. Since the larger model performed better than the smaller ones fine-tuned on Oscar and GerParCor one could assume that the performance was influenced by the sheer size of input data and model parameters. However, there are enough accounts for outstanding performances in current machine learning applications for smaller models (Schick and Schütze (2020); Izacard et al. (2022)). Kaplan et al. provide an extensive study of model scaling stating that while larger models tend to be more sample-efficient, any of the factors computational capacity, model size and amount of data can compensate for the other 2020, p. 3. As a consequence the smaller models trained on Oscar and GerParCor should generally be able to achieve better scores than their evaluation indicates.

The impact of modifying the tokenization on a single POS category in masked language learning must be close to none, given that the tasks do not explicitly focus on that category. In other endeavors Tan, Joty, Varshney, et al. used the Penn Treebank to model each inflection with designated tokens instead of analyzing the textual surface form, showing that an improvement of WordPiece tokenization is morphologically augmentable with POS! 2020. Their preceding use of morphological concepts in language learning is found in their earlier research through adversarial shifting of morphological paradigms to improve robustness of

existing language models (Tan, Joty, Kan, et al. 2020, p. 2927). Implementing supporting processes for morphological modeling on an integral scale seems to be a condition for training impactful models. This could be one explanation for the performance of the Wordmap based models but does not cover the standard models.

The classification task used to see if the model predictions are noteworthy or negligible was partly inconclusive. As stated before the smaller models do not compare to the bbgc model. Since this one shot evaluation made even bbgc perform worse than its other classification tasks (Wiegand (2019); GNAD¹) there may be an interfering factor in the task design. The high number false positives may be attributed to imbalanced sampling (Brownlee 2020, pp. 5–6). This could be an intrinsic property of the dataset or happen during the modification of the sample labels. On a side note the Wordmap models were designed to explicitly tokenize verbs with the aim of encoding an important core element of language, but the items in the Wikinews dataset featured news titles sometimes using elliptic phrasing and entirely dropping verbs. To explain the uniform performance drop in all models additional constraints have to be taken into consideration. The models are all masked language models, being a fine-tuned version of bbgc. Instead of unfreezing all layers and retraining a more suitable approach would be to try transfer learning instead, preserving useful ggbc encoding layers. To find out which of the structural modeling factors prevented the model from performing satisfyingly more training would have to be done, paying special attention to backpropagation, controlling the loss function in both fine-tuning and benchmark testing (ideally on runtime).

Considering the overall low performance it is surprising to see mwg5 performing at the given rate. Whether this can be ascribed to the aforementioned sample imbalance or learning rate jitter is impossible to tell in retrospect, although the latter is rather unlikely in view of the other small model performances using the same small learning rate. Positively the transfer of new subwords from the POS vocabulary to the BERT tokenizer vocabulary does not visibly impair performance. All the input verbs come from sanitized data and GerParCor, and since mwo5 (fine-tuned on Oscar) was the best performing small model there is even an argument for saying it did not carry over a vocabulary bias. Ultimately it is hard to classify the performance in this iteration of fine-tuning and testing. Judging solely by the scores it should not be compared to recent methods like the original inspiration FLOTA (Hofmann, Schuetze, and Pierrehumbert 2022).

These observations allow answering the first hypothesis: the classification task demonstrates irregular scoring behavior for Wordmap models surpassing (mwo5) and also failing to reach

¹<https://tblock.github.io/10kGNAD/>

(mwig5) comparable standard models msg5 and mso5. No effect could be connected to the tokenization with Wordmap, the null hypothesis remains.

Concerning the Wordmap vocabulary generation it can be said that Wordmap shows effective extraction and analysis of segments in german verbs. Its segmenting capabilities hold up to the FLOTA algorithm under the assumption that other lexical POS categories yield similar results. This assumption is favored by the idea of morphological complexity given in section 2.2. From a linguistic point of view there are a few reservations that should be addressed. The Wordmap algorithm does not work on any grammatical framework and is unable to connect morphemes that would otherwise be analyzed as one, for example circumfixes. This makes for lighter and easier implementation but outsources the connection of spatially separated morphemes and their functional exponence to the layers of the neural network (unlike e.g. Tan, Joty, Varshney, et al. (2020)). Recent language models deliver for languages with allomorphisms or multiple exponence² despite not analyzing segmental dependencies at all and perform very well (this does not mean they cannot be improved by doing so).

Naturally there is still room for improvement regarding the maximization function, scope and especially testing. As of now the segmentation algorithm behaves similar to FLOTA in that it discriminates subwords by length, among other features. The lexical bias is mainly provided by the stand-alone vocabulary generated beforehand. This can most likely be improved by incorporating relative subword frequency into the maximization weight, following the example of successful tokenizers like WordPiece (WP) and BPE do nowadays. More ample implications pose questions as to how Wordmap would fare under circumstances of reanalysis, loaning, grammaticalization or lexicalization. These processes usually take a long time and require many (linguistically) productive cycles, so that affected words will probably appear in corpora for a Wordmap generation algorithm to catch. Another remark to be made in respect to character frequency. The problem of identifying segments and discarding characters that do not belong to an adjacent morpheme is not trivial. For Wordmap this poses the greatest challenge, leading to most of the unwanted fragments in both functional and lexemic vocabulary. A suggestion is to do character sampling for retrieving (maybe even POS specific) character frequencies and using those as dampening for the Wordmap segment detection. The detection itself is still binary and should be revised to detect segment counts that are in a range instead. During development, a good indicator of segment boundaries where values around the mean, which almost are consistently in proximity of each other. Other interesting operations with Wordmaps worth mentioning are derivatives. In general,

²Brief summary on morphological exponence: https://home.uni-leipzig.de/jtrommer/netzwerk_a_ntrag.pdf

the first (mathematic) derivative of a word map seems to coincide with a segment boundary. The second derivative can also mark morpheme boundaries, though this is just a presumption based on observations during plotting. In most cases the resolution of the Wordmap is too volatile. The application scope and testing cycle determine each other. Examining Wordmap performance for more POS categories and languages would broaden its scope and probably reach their limit at the domain of incorporating languages, at which point the method needs fundamental expansion. At the given the test scores are the only statistical measure provided that may be used to verify hypothesis 2, except for the test scores. The previous discussion of these scores revealed that they cannot be used to draw reliable conclusions about the Wordmap tokenizer performance. A qualitative assessment suggests that segmentations are less noisy for the intended POS class. It would be in good faith to say that the tokenization is improved with Wordmap. H0 stands unfalsified.

6 Conclusion

In this thesis a data-driven tokenization method was explored to see whether the distinction of lexical and functional morphemes is beneficial towards tokenization. The scope of the thesis was set to a single language and a single Part-of-Speech category, comparing word shapes and morphological complexity in similar languages. Crawled and sanitized verbs were used to heuristically generate separate vocabularies that were then modified to remove competing subwords. A segmentation algorithm was implemented to judge varying segmentations by a maximization function based around the token coverage, subword length and number of segments. These procedures were integrated into a pre-trained BERT WordPiece tokenizer. Different versions (standard and Wordmap) of the same BERT base model were fine-tuned as masked language model on two datasets and tested on a sequence classification task. The model benchmark results were inconclusive and the below-average performance was partly accounted for. The tokenization process was found to be functional and seemed to produce results comparable to standard tokenization methods. Final measurements of effectiveness have yet to be determined to confirm this. Neither hypothesis was verified in the process.

When conducting further research on morphological tokenization the linguistic split of data is encouraged. Optimizing the maximization function was a small part playing a core role: its weights are numerical characteristic of morphemes and finding appropriate characteristics was a welcomed challenge that offers deep research potential for every language. It is advised to apply loss control while fine-tuning and testing to ensure the model parameters update in favor of the desired learnable feature. This way the results will be much better interpretable.

In all statistical modeling the goal is to model reality as precisely as feasible. Language models are no exception, and compromises are always made. Human language understanding builds on maximizing relevant input, without it there is no way to generate meaningful output. Machine learning and human learning offer different levers to achieve this goal. While computation cost is an upcoming problem in natural language processing and machine learning, it does not take much space to split language model input into more classified sets of data. The domain of language modeling will keep evolving as will language. Though maturity is in arm's reach at the presence of highly parameterized or specialized models, the data they are performing on ultimately traces back to humans. Our tendencies to organize

language data are a basic property of our cognition, just like the distinction between what expresses a grammatical function and what is a lexical reference. The discriminatory functions we have learned to speak, write and read by do not have to be learned from scratch. They can be an integral part of how to accelerate neural networks and an opportunity to keep a close eye on the data we select for machine learning.

The methods presented in this thesis are minute in the face of this vision, but hopefully every morphological tokenizer is a step towards it.

Code repository: <https://github.com/s2458588/wsm-tokenizer.git>

Thesis repository: **Used Python Packages:**

```

aiohttp==3.8.1 Jinja2==3.1.2 pyOpenSSL==22.0.0
aiosignal==1.2.0 joblib==1.1.0 pyparsing==3.0.9
argon2-cffi==21.3.0 jsonschema==4.16.0 PyQt5==5.15.7
argon2-cffi-bindings==21.2.0 jupyter==1.0.0 PyQt5-sip==12.11.0
async-timeout==4.0.2 jupyter_client==7.3.5 pyrsistent==0.18.1
asynctest==0.13.0 jupyter-console==6.4.4 PySocks==1.7.1
attrs==22.1.0 jupyter_core==4.11.1 python-dateutil==2.8.2
backcall==0.2.0 jupyterlab-pygments==0.2.2 pytz==2022.2.1
backports.functools-lru-cache==1.6.4 jupyterlab-widgets==3.0.3 PyYAML==6.0
beautifulsoup4==4.11.1 kiwisolver==1.4.4 pyzmq==23.2.1
bleach==5.0.1 lxml==4.9.1 qtconsole==5.3.2
Bottleneck==1.3.5 MarkupSafe==2.1.1 QtPy==2.2.0
brotlipy==0.7.0 matplotlib==3.5.3 regex==2022.9.13
certifi==2022.12.7 matplotlib-inline==0.1.6 requests==2.28.1
cffi==1.15.1 mistune==2.0.4 responses==0.18.0
charset-normalizer==2.1.1 more-itertools==9.0.0 sacremoses==0.0.53
click==8.1.3 mpmath==1.2.1 scikit-learn==1.0.2
colorama==0.4.5 multidict==6.0.2 scipy==1.7.3
cryptography==37.0.4 multiprocessing==0.70.13 Send2Trash==1.8.0
cyclor==0.11.0 munkres==1.1.4 setuptools==65.3.0
dataclasses==0.8 nbclient==0.6.8 sip==6.6.2
datasets==2.4.0 nbconvert==7.0.0 six==1.16.0
debugpy==1.6.3 nbformat==5.4.0 soupsieve==2.3.2.post1
decorator==5.1.1 nest-asyncio==1.5.5 sympy==1.10.1
defusedxml==0.7.1 nltk==3.8.1 terminado==0.15.0
dill==0.3.5.1 notebook==6.4.12 threadpoolctl==3.1.0
entrypoints==0.4 numexpr==2.8.3 tinycss2==1.1.1
fastjsonschema==2.16.1 numpy==1.21.6 tokenizers==0.12.1
filelock==3.8.0 packaging==21.3 toml==0.10.2
findiff==0.8.9 pandas==1.3.5 torch==1.12.1.post200
flit_core==3.7.1 pandocfilters==1.5.0 tornado==6.2
fonttools==4.25.0 parso==0.8.3 tqdm==4.64.1
frozenlist==1.3.1 pexpect==4.8.0 traitlets==5.3.0
fsspec==2022.8.2 pickleshare==0.7.5 transformers==4.21.3
gmpy2==2.1.2 Pillow==9.2.0 typing_extensions==4.3.0
HanTa==1.0.0 pip==22.2.2 urllib3==1.26.11
huggingface-hub==0.9.1 pkgutil_resolve_name==1.3.10 wcwidth==0.2.5
idna==3.3 ply==3.11 webencodings==0.5.1
importlib-metadata==4.11.4 prometheus-client==0.14.1 wheel==0.37.1
importlib-resources==5.9.0 prompt-toolkit==3.0.31 widgetsnbextension==4.0.3
ipykernel==6.15.2 psutil==5.9.2 xxhash==0.0.0
ipython==7.33.0 ptyprocess==0.7.0 yarl==1.7.2
ipython-genutils==0.2.0 pyarrow==8.0.0 zipp==3.8.1
ipywidgets==8.0.2 pycparser==2.21
jedi==0.18.1 Pygments==2.13.0

```

Bibliography

- Abrami, Giuseppe, Mevlüt Bağcı, Leon Hammerla, and Alexander Mehler (June 2022). “German Parliamentary Corpus (GerParCor)”. In: *Proceedings of the Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, pp. 1900–1906. URL: <https://aclanthology.org/2022.lrec-1.202>.
- Aeppli, Noëmi and Rico Sennrich (2021). *Improving Zero-shot Cross-lingual Transfer between Closely Related Languages by injecting Character-level Noise*. DOI: 10.48550/ARXIV.2109.06772. URL: <https://arxiv.org/abs/2109.06772>.
- Aikhenvald, Alexandra and R Dixon (2017). *The Cambridge Handbook of Linguistic Typology*. Cambridge Handbooks in Language and Linguistics. Cambridge: Cambridge University Press. DOI: 10.1017/9781316135716.
- Baerman, Matthew, Dunstan Brown, and Greville G. Corbett, eds. (2017). *Morphological Complexity*. Cambridge studies in Linguistics 153. Cambridge University Press. ISBN: 1107120640. DOI: 10.1017/9781316343074. URL: www.cambridge.org/9781107120648.
- Bentz, Christian, Tatyana Ruzsics, Alexander Koplenig, and Tanja Samardžić (Dec. 2016). “A Comparison Between Morphological Complexity Measures: Typological Data vs. Language Corpora”. In: *Proceedings of the Workshop on Computational Linguistics for Linguistic Complexity (CL4LC)*. Osaka, Japan: The COLING 2016 Organizing Committee, pp. 142–153. URL: <https://aclanthology.org/W16-4117>.
- Bouckaert, Remco et al. (2012). “Mapping the Origins and Expansion of the Indo-European Language Family”. In: *Science* 337.6097, pp. 957–960. DOI: 10.1126/science.1219669. eprint: <https://www.science.org/doi/pdf/10.1126/science.1219669>. URL: <https://www.science.org/doi/abs/10.1126/science.1219669>.
- Brennan, J.R. (2022). *Language and the Brain: A Slim Guide to Neurolinguistics*. Oxford linguistics. Oxford University Press. ISBN: 9780198814757. URL: <https://books.google.de/books?id=KLe5zgEACAAJ>.
- Brownlee, J. (2020). *Imbalanced Classification with Python: Better Metrics, Balance Skewed Classes, Cost-Sensitive Learning*. Machine Learning Mastery. URL: <https://books.google.de/books?id=jaXJDwAAQBAJ>.
- Chan, Branden, Stefan Schweter, and Timo Möller (June 2019). *GermanBERT*. URL: <https://huggingface.co/bert-base-german-cased>.

- Colman, Andrew M. (2009). *morpheme*. DOI: 10.1093/acref/9780199534067.013.5219. URL: <https://www.oxfordreference.com/view/10.1093/acref/9780199534067.001.0001/acref-9780199534067-e-5219>.
- Çöltekin, Çağrı and Taraka Rama (2022). “What do complexity measures measure? Correlating and validating corpus-based measures of morphological complexity”. In: *Linguistics Vanguard*. DOI: doi:10.1515/lingvan-2021-0007. URL: <https://doi.org/10.1515/lingvan-2021-0007>.
- Comrie, Bernard (1989). “Morphological Typology”. In: *Language universals and linguistic typology*. 2nd ed. University of Chicago Press, pp. 42–56.
- Creutz, Mathias and Krista Lagus (2002). “Unsupervised Discovery of Morphemes”. In: *CoRR* cs.CL/0205057. URL: <https://arxiv.org/abs/cs/0205057>.
- Delogu, Francesca, Harm Brouwer, and Matthew W. Crocker (2019). “Event-related potentials index lexical retrieval (N400) and integration (P600) during language comprehension”. In: *Brain and Cognition* 135, pp. 1–3. ISSN: 0278-2626. DOI: <https://doi.org/10.1016/j.bandc.2019.05.007>. URL: <https://www.sciencedirect.com/science/article/pii/S0278262618304299>.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. DOI: 10.48550/ARXIV.1810.04805. URL: <https://arxiv.org/abs/1810.04805>.
- Ehret, Katharina, Alice Blumenthal-Dramé, Christian Bentz, and Aleksandrs Berdicevskis (2021). “Meaning and Measures: Interpreting and Evaluating Complexity Metrics”. In: *Frontiers in Communication* 6. ISSN: 2297-900X. DOI: 10.3389/fcomm.2021.640510. URL: <https://www.frontiersin.org/articles/10.3389/fcomm.2021.640510>.
- Glück, Helmut and Michael Rödel, eds. (2016). *Metzler Lexikon Sprache*. ger. 5th ed. Springer eBook Collection. Stuttgart: J.B. Metzler, pp. 141–142. ISBN: 978-3-476-05486-9. DOI: 10.1007/978-3-476-05486-9. URL: <http://dx.doi.org/10.1007/978-3-476-05486-9>.
- Hofmann, Valentin, Hinrich Schuetze, and Janet Pierrehumbert (May 2022). “An Embarrassingly Simple Method to Mitigate Undesirable Properties of Pretrained Language Model Tokenizers”. In: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Dublin, Ireland: Association for Computational Linguistics, pp. 385–393. DOI: 10.18653/v1/2022.acl-short.43. URL: <https://aclanthology.org/2022.acl-short.43>.
- Iacobini, Claudio and Francesca Masini (Sept. 2005). “Verb-particle Constructions and Prefixed Verbs in Italian: Typology, Diachrony and Semantics”. In: *Proceedings of the Fifth Mediterranean Morphology Meeting (MMM5)*. Università degli Studi di Bologna. URL: <https://>

- [//www.academia.edu/1183339/Verb_particle_constructions_and_prefixed_verbs_in_Italian_typology_diachrony_and_semantics](http://www.academia.edu/1183339/Verb_particle_constructions_and_prefixed_verbs_in_Italian_typology_diachrony_and_semantics).
- Islam, Md Mofijul et al. (2022). “A Vocabulary-Free Multilingual Neural Tokenizer for End-to-End Task Learning”. In: DOI: 10.48550/ARXIV.2204.10815. URL: <https://arxiv.org/abs/2204.10815>.
- Izacard, Gautier et al. (2022). *Atlas: Few-shot Learning with Retrieval Augmented Language Models*. DOI: 10.48550/ARXIV.2208.03299. URL: <https://arxiv.org/abs/2208.03299>.
- Jianlong Zhou, Fang Chen (2018). *Human and Machine Learning: Visible, Explainable, Trustworthy and Transparent*. 1st ed. 2018. Human–Computer Interaction Series. Springer. ISBN: 3319904027.
- Kaplan, Jared et al. (2020). “Scaling Laws for Neural Language Models”. In: *CoRR* abs/2001.08361. arXiv: 2001.08361. URL: <https://arxiv.org/abs/2001.08361>.
- Kimppa, Lilli et al. (2019). “Acquisition of L2 morphology by adult language learners”. In: *Cortex* 116. Structure in words: the present and future of morphological processing in a multidisciplinary perspective, pp. 74–90. ISSN: 0010-9452. DOI: <https://doi.org/10.1016/j.cortex.2019.01.012>. URL: <https://www.sciencedirect.com/science/article/pii/S0010945219300310>.
- Kiss, Katalin É. and Veronika Hegedus, eds. (2021). *Postpositions and Postpositional Phrases*. Amsterdam: Amsterdam University Press, p. 262. ISBN: 9789048544608. DOI: [doi:10.1515/9789048544608](https://doi.org/10.1515/9789048544608). URL: <https://doi.org/10.1515/9789048544608>.
- Kühl, Niklas, Marc Goutier, Lucas Baier, Clemens Wolff, and Dominik Martin (2020). “Human vs. supervised machine learning: Who learns patterns faster?” In: *CoRR* abs/2012.03661. arXiv: 2012.03661. URL: <https://arxiv.org/abs/2012.03661>.
- Lehman, Christian W. (Aug. 2022). *Indogermanisch*. URL: <https://www.christianlehmann.eu/ling/sprachen/indogermania/RomGesch/idg.php>.
- Nikanne, Urpo (Dec. 2017). *Finite sentences in Finnish: Word order, morphology, and information structure*. DOI: 10.5281/zenodo.1117710. URL: <https://doi.org/10.5281/zenodo.1117710>.
- Ortiz Su’arez, Pedro Javier, Laurent Romary, and Benoit Sagot (June 2020). “A Monolingual Approach to Contextualized Word Embeddings for Mid-Resource Languages”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, pp. 1703–1714. URL: <https://www.aclweb.org/anthology/2020.acl-main.156>.
- Ortiz Su’arez, Pedro Javier, Benoit Sagot, and Laurent Romary (2019). “Asynchronous pipelines for processing huge corpora on medium to low resource infrastructures”. en. In: ed. by Pi-

- otr Bański et al. Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-7) 2019. Cardiff, 22nd July 2019. Mannheim: Leibniz-Institut für Deutsche Sprache, pp. 9–16. DOI: 10.14618/ids-pub-9021. URL: <http://nbn-resolving.de/urn:nbn:de:bsz:mh39-90215>.
- OSCAR (Jan. 2022). *The OSCAR project (Open Super-large Crawled Aggregated coRpus)*. URL: <https://oscar-project.org/>.
- Paszke, Adam et al. (2019). *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. DOI: 10.48550/ARXIV.1912.01703. URL: <https://arxiv.org/abs/1912.01703>.
- Patsy M. Lightbown, Nina Spada (2013). *How Languages Are Learned*. 4th. Oxford University Press, p. 114. ISBN: 9780194541268. URL: <https://global.oup.com/academic/product/how-languages-are-learned-9780194406291?lang=en&cc=gb>.
- Peters, Ben and Andre F. T. Martins (July 2022). “Beyond Characters: Subword-level Morpheme Segmentation”. In: *Proceedings of the 19th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*. Seattle, Washington: Association for Computational Linguistics, pp. 131–138. DOI: 10.18653/v1/2022.sigmorphon-1.14. URL: <https://aclanthology.org/2022.sigmorphon-1.14>.
- Radford, Alec, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever (2018). “Improving language understanding by generative pre-training”. In: URL: https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf.
- Rust, Phillip, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych (2020). *How Good is Your Tokenizer? On the Monolingual Performance of Multilingual Language Models*. DOI: 10.48550/ARXIV.2012.15613. URL: <https://arxiv.org/abs/2012.15613>.
- Schick, Timo and Hinrich Schütze (2020). *It’s Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners*. DOI: 10.48550/ARXIV.2009.07118. URL: <https://arxiv.org/abs/2009.07118>.
- Schuster, Mike and Kaisuke Nakajima (2012). “Japanese and Korean voice search”. In: *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5149–5152. DOI: 10.1109/ICASSP.2012.6289079.
- Schwartz, Lane et al. (2020). *Neural Polysynthetic Language Modelling*. DOI: 10.48550/ARXIV.2005.05477. URL: <https://arxiv.org/abs/2005.05477>.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch (Aug. 2016). “Neural Machine Translation of Rare Words with Subword Units”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Asso-

- ciation for Computational Linguistics, pp. 1715–1725. DOI: 10.18653/v1/P16-1162. URL: <https://aclanthology.org/P16-1162>.
- Talmor, Alon, Yanai Elazar, Yoav Goldberg, and Jonathan Berant (Dec. 2020). “oLMpics-On What Language Model Pre-training Captures”. In: *Transactions of the Association for Computational Linguistics* 8, pp. 743–758. ISSN: 2307-387X. DOI: 10.1162/tacl_a_00342. eprint: https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00342/1923716/tacl_a_00342.pdf. URL: https://doi.org/10.1162/tacl%5C_a%5C_00342.
- Tan, Samson, Shafiq Joty, Min-Yen Kan, and Richard Socher (July 2020). “It’s Morphin’ Time! Combating Linguistic Discrimination with Inflectional Perturbations”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, pp. 2920–2935. DOI: 10.18653/v1/2020.acl-main.263. URL: <https://aclanthology.org/2020.acl-main.263>.
- Tan, Samson, Shafiq Joty, Lav Varshney, and Min-Yen Kan (Nov. 2020). “Mind Your Inflections! Improving NLP for Non-Standard Englishes with Base-Inflection Encoding”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online: Association for Computational Linguistics, pp. 5647–5663. DOI: 10.18653/v1/2020.emnlp-main.455. URL: <https://aclanthology.org/2020.emnlp-main.455>.
- Toraman, Cagri, Eyup Halit Yilmaz, Furkan Şahinuç, and Oguzhan Ozelik (2022). *Impact of Tokenization on Language Models: An Analysis for Turkish*. DOI: 10.48550/ARXIV.2204.08832. URL: <https://arxiv.org/abs/2204.08832>.
- Vaswani, Ashish et al. (2017). *Attention Is All You Need*. DOI: 10.48550/ARXIV.1706.03762. URL: <https://arxiv.org/abs/1706.03762>.
- Wartena, Christian (2019). “A Probabilistic Morphology Model for German Lemmatization”. In: *Proceedings of the 15th Conference on Natural Language Processing (KONVENS 2019)*, pp. 40–49. DOI: 10.25968/opus-1527. URL: <http://nbn-resolving.de/urn:nbn:de:bsz:960-opus4-15271>.
- Wiegand, Michael (2019). *GermEval-2018 Corpus (DE)*. Version V1. DOI: 10.11588/data/0B5VML. URL: <https://doi.org/10.11588/data/0B5VML>.
- Wiktionary (Sept. 2022). *Verzeichnis: Deutsch*. URL: <https://de.wiktionary.org/wiki/Verzeichnis:Deutsch>.
- Wu, Yonghui et al. (2016). *Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation*. DOI: 10.48550/ARXIV.1609.08144. URL: <https://arxiv.org/abs/1609.08144>.