

Ricardo Lukas Jung  
6227492  
Empirische Sprachwissenschaft (B.A.)  
Phonetik & Digital Humanities  
15<sup>th</sup> 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 11, 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

## **Abstract**

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

# Contents

<b>List of Figures</b>	<b>II</b>
<b>List of Tables</b>	<b>III</b>
<b>List of Acronyms</b>	<b>IV</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Overview</b>	<b>2</b>
<b>3 Methodology</b>	<b>3</b>
1 Requirements . . . . .	3
1.1 Pipeline Structure & Transformers Library . . . . .	3
1.2 Data . . . . .	5
1.3 Benchmark . . . . .	5
2 Implementation . . . . .	5
2.1 Tokenizer . . . . .	5
2.2 Masked Language Model . . . . .	11
2.3 Benchmark . . . . .	11
<b>4 Results</b>	<b>13</b>
1 Benchmark . . . . .	13
2 Tokenization . . . . .	13
<b>5 Discussion</b>	<b>15</b>
<b>6 Conclusion</b>	<b>16</b>
<b>7 Testchapter</b>	<b>17</b>
1 Citing . . . . .	17
2 Quoting . . . . .	17
3 Referencing . . . . .	17



List of Figures

1	Schematic integration of Wordmap segmentation into the BERT architecture. Continuous lines show non-selective channels, while dotted lines only pass data selectively . . . . .	4
2	Permutations SEGMENTER generates from target i.e. verstehen. Dotted segments mark the segmentation chosen later by the maximization function during tokenization. . . . .	10

## List of Tables

1	Example wordmaps for <i>target</i> = verstehen . . . . .	8
2	ID references assigned to full model names . . . . .	11
3	List of all used models and their hyperparameters. LR = learning rate, WU = warmup steps. . . . .	12
4	Metrics for masked language model trained on the Oscar dataset with Wordmap infused tokenization. Evaluated on sequence classification task. . . . .	13
5	Metrics for masked language model trained on the GerParCor dataset with Wordmap infused tokenization. Evaluated on sequence classification task. . . . .	14
6	Metrics for masked language model trained on the GerParCor dataset with bert-base-german-cased (bbgc) tokenization. Evaluated on sequence classification task. . . . .	14
7	Metrics for masked language model trained on the GerParCor dataset with bbgc tokenization. Evaluated on sequence classification task. . . . .	14
14	Table 7.8: Test score summary for all evaluated models. . . . .	14

## List of Acronyms

**bbgc** bert-base-german-cased

**BERT** Bidirectional Encoders from Transformers

**BPE** Byte Pair Encoding

**CL** Computational Linguistics

**GerParCor** German Parliamentary Corpus

**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



# 1 Introduction

## 2 Overview

## 3 Methodology

in this section the whole methodology is covered. what do i use in this thesis, why do i use it and lastly, how? make sure the why covers methodological implications. (vergiss nicht alle pakete als quelle im Anhang)

### 1 Requirements

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

#### 1.1 Pipeline Structure & Transformers Library

Bidirectional Encoders from Transformers (BERT) is a language learning transformer model designed for Natural Language Processing (NLP) tasks (Vaswani et al. 2017). Upon release it achieved higher performance scores compared to previously used Long Short-Term Memory (LSTM) models (BERTHIGH1). Two main model characteristics can be observed for BERT. Firstly, 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 (at the time novel) Masked Language Model (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.

The NLP community has since developed BERT and adapted it to the needs of contemporary NLP problems (roberta, germanbert, mbert CITATION). Its wide 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 module (QUELLE?). 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 tokenizer is trained on a dataset it is specific to that dataset. This means the application of a tokenizer on another dataset may result in out-of-vocabulary issues and

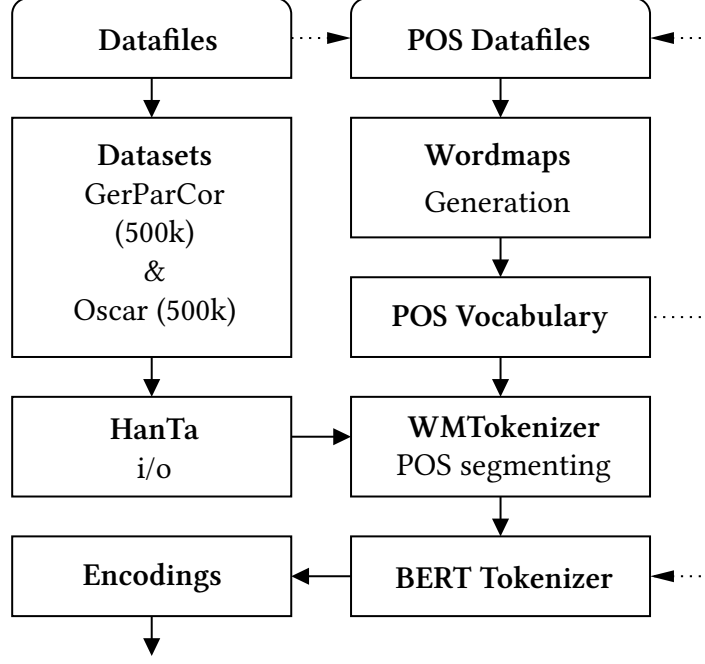


Figure 1: Schematic integration of Wordmap segmentation into the BERT architecture. Continuous lines show non-selective channels, while dotted lines only pass data selectively

different token/subtoken distributions.

Particularly relevant to this thesis is the option to train an own tokenizer from the base module. Usually, WordPiece generates its own set of subtokens called *vocabulary*, which is used by the standard tokenizer to generate unique input numerical 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 2.1. Once a string is tokenized it gets passed on to the transformer model for contextualized embedding. Schematic Figure 1 shows the processing of data up until token encoding.

In Figure 1, Datafiles contains the raw unprocessed data coming from corpora. The input is provided in txt-files holding sentences line by line. Two datasets are generated from this input, each amounting to 500000 sentences per corpus for fine-tuning. POS Datafiles ideally contains only words that are identified as the same Part of Speech (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. Otherwise a predefined list of words can be used as an external and generic source for Wordmap generation. In POS Vocabulary, naive 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 to the trade-off between vocabulary size and model performance, hence the addition of POS Vocabulary to the BERT Tokenizer vocabulary is used only when necessary. Datasets are then passed through Hanover Tagger (HanTa) to provide WMTTokenizer with flagged data to tokenize only selected tokens. Lastly BERT Tokenizer returns the canonical encodings known from language modeling.

## 1.2 Data

Two corpora where selected for fine-tuning: German Parliamentary Corpus (GerParCor) and oscar (**oscar**). FLOTA was trained on 12000 samples per category Hofmann, Schuetze, and Pierrehumbert 2022, which is why for this fine-tuning a sample size of 500000 seems sufficient. GerParCor is a “GerParCor 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 texts and transcripts of the german parliament. All 500k samples in the GerParCor dataset used for this thesis is 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 **oscarsub2**; (**oscarsub1**).

The POS vocabulary comes from **wiktionary** hanta-crawled verbs of gerparcor **hanta**

## 1.3 Benchmark

The benchmark used on the model is oLMpics (**olmpics**)

# 2 Implementation

Tatsächliche Anwendung der Methoden auf die Daten

## 2.1 Tokenizer

ESSENTIALLY DERIVING SENSIBLE SUBTOKENS TO REPRESENT LEXEMES The methods of tokenization is explained in the following. This is the main part of the thesis containing

two algorithms for vocabulary generation and token segmentation. Once the two algorithms have been explained, the last section of this chapter features the

Languages, as consistent form of communication, always 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 **hanta**. Initial experiments where done only with the wiktionary verb list since it provided sanitized input for the algorithm in development. As soon as the verbs where extracted from GerParCor via HanTa

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 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  $\delta$  difference in length between *s* and *l*.  $\delta$  functions as an offset for index-based comparisons WORDMAP.

WORDMAP is a naive 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>, <

---

**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
```

---

Table 1: 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	
(...)			

<s> will be more frequent than others). 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.

### BERT Tokenizer Modification

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. 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 a trees containing every permutation of the set



---

**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( $segmentations$ )
```

---

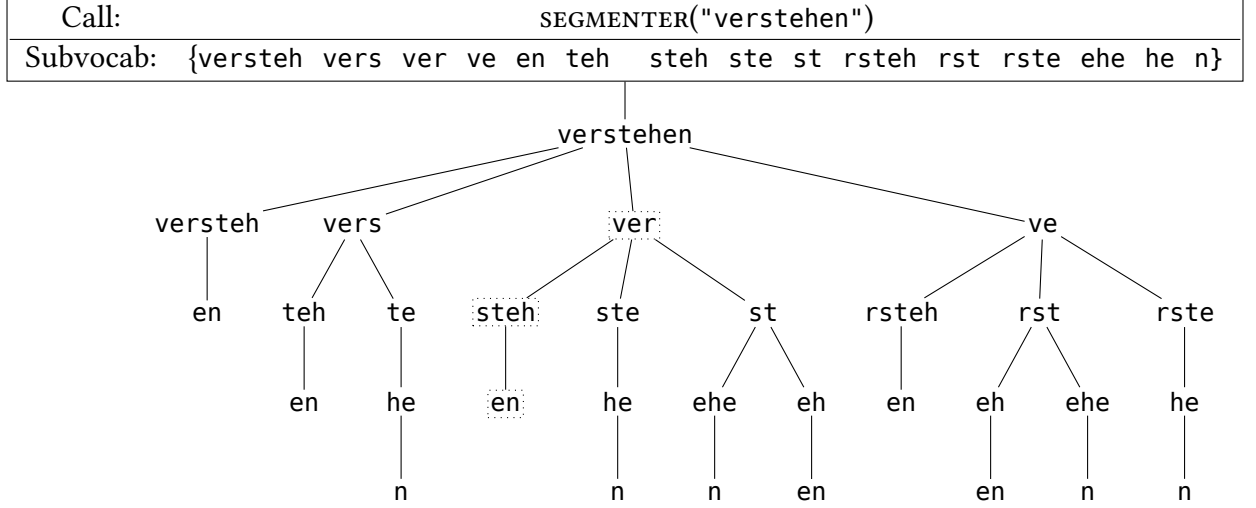


Figure 2: Permutations SEGMENTER generates from target i.e. verstehen. Dotted segments mark the segmentation chosen later by the maximization function during tokenization.

*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:

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

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

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

Table 2: ID references assigned to full model names

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  $s_i^{\text{th}}$  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 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

## 2.2 Masked Language Model

All fine-tuned models use the bbgc baseline by **bertbasegermancased**. For reference and readability they are given shortened IDs as follows:

The full model names in Table 2 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 500000 samples, while msg5 is the the **m**asked language model using **s**tandard BERT wordpiece tokenization fine-tuned on the **G**erParCor dataset of 500000 samples. All versions

## 2.3 Benchmark

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 <sup>1</sup>	BWP	0.0001	810000	1024	10000	-

Table 3: List of all used models and their hyperparameters. LR = learning rate, WU = warmup steps.

## 4 Results

### 1 Benchmark

Every benchmark done with the models are

### 2 Tokenization

Show specific examples of tokenization and analyze the qualitatively (maybe quantitatively)

mlm_wmt_oscar500k	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

Table 4: Metrics for masked language model trained on the Oscar dataset with Wordmap infused tokenization. Evaluated on sequence classification task.

<b>mlm_wmt_gpc500k</b>	<b>Epoch 1</b>	<b>Epoch 2</b>	<b>Epoch 3</b>	<b>Test score</b>
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

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

<b>mlm_std_oscar500k</b>	<b>Epoch 1</b>	<b>Epoch 2</b>	<b>Epoch 3</b>	<b>Test scores</b>
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

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

<b>mlm_std_gpc500k</b>	<b>Epoch 1</b>	<b>Epoch 2</b>	<b>Epoch 3</b>	<b>Test scores</b>
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

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

<b>bbgc</b>	<b>Epoch 1</b>	<b>Epoch 2</b>	<b>Epoch 3</b>	<b>Test scores</b>
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

Table 8: Metrics for masked language model baseline bert-base-german-cased<sup>1</sup>. Evaluated on sequence classification task.

<b>Summary</b>	<b>bbgc</b>	<b>std+oscar</b>	<b>std+gpc</b>	<b>wmt+oscar</b>	<b>wmt+gpc</b>
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 9: Test score summary for all evaluated models.

## **5 Discussion**

## **6 Conclusion**



## 7 Testchapter

### 1 Citing

Abrami et al. 2022

### 2 Quoting

“This is a quote by textquote” (DeepL 2021) “This is a quote by enquote”

### 3 Referencing

Short reference ??

Long reference ??

monofont for code or string monofont

## 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](http://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>.
- 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>.
- DeepL (Nov. 2021). *How does deepl work?* URL: <https://www.deepl.com/en/blog/how-does-deepl-work>. Last accessed: 28.12.2022.
- 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>.
- Lehman, Christian W. (Aug. 2022). *Indogermanisch*. URL: <https://www.christianlehmann.eu/ling/sprachen/indogermania/RomGesch/idg.php>.
- 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](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>.

- 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: Association for Computational Linguistics, pp. 1715–1725. DOI: 10.18653/v1/P16-1162. URL: <https://aclanthology.org/P16-1162>.
- Toraman, Cagri, Eyup Halit Yilmaz, Furkan Şahinuç, and Oguzhan Ozcelik (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>.
- 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>.