

Ricardo Lukas Jung  
6227492  
Bachelor  
Empirische Sprachwissenschaft (Phonetik & Digital Humanities)  
15<sup>th</sup> Semester  
s2458588@stud.uni-frankfurt.de

**Thesis submitted in fulfilment of the requirements for the  
degree of Bachelor of Arts**

# **Lexicalizing a BERT Tokenizer**

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

Ricardo Lukas Jung

Date of Submission:  
December 21, 2022

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

Contents

1 Introduction 4

1.0.1 Background . . . . . 4

1.0.2 Motivation . . . . . 4

1.0.3 Explain the what and how . . . . . 5

1.0.4 Scope of this Thesis . . . . . 5

Bibliography 6

# 1 Introduction

This chapter covers the background, intentions and scope of the thesis.

## 1.0.1 Background

## 1.0.2 Motivation

There is an ongoing urge in the computational linguistics (CL) community to understand natural language. Research in the past decades shows use of frequentist and statistical methods (such as ZITATION) to their advantage, leading to the emergence of the first machine learning (ML) models. It became apparent that these ML models are the best currently available approach to an automated understanding of natural language. The structural parallels of machine learning to human learning have often been drawn (ZITATION)) to demonstrate how similar and more important: how different both can be. A powerful feature of ML (as opposed to human learning) is the possibility of actively controlling the the learning parameters in a supervised environment. To test the efficiency of ML parameters a variety of tasks (ZITATION) are designed and applied. A trained model will yield performance scores based on the quality of its training, much like humans on language tests. But the automated modeling of language is not the first instance language modelling in a broader sense. Traditional linguistics (DEFINITION has produced fundamental research the 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 somewhat satisfactory results, they are far from complete or perfect. The integration of aforementioned traditional linguistic knowledge into learning processes for machine learning is the underlying motivation of this thesis.

Language learners usually first learn a lexicon consisting of lexemes which they will have to analyze accurately in order to be productive in that target language. A ML model relies on a tokenizer to create such a vocabulary (ZITATION). It is programmed to segment tokens into subwords (if possible) and provide a vocabulary comprising all the components needed to analyze a given string. Ideally those subwords will be part of the functional vocabulary in the target language, so called morphemes ERKLÄRUNG. A morpheme is defined as the smallest unit carrying meaning in a language. The morphemes of a language and its generated tokenizer vocabulary rarely coincide. Typically, tokenzier vocabularies will contain a lot of noise and linguistically nonsensical segmentations or words. Following the guiding principle that **input quality is ouput quality** not only in language learning, the morpheme vocabulary is identified as the point of leverage in the upcoming section. Note: explain why i use tokens and words, they are interchaneable right? holistic, need less attention to produce satisfactory

### **1.0.3 Explain the what and how**

### **1.0.4 Scope of this Thesis**

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

cite (Glück and Rödel 2016)  
citeast (2016)

cite (Tenenbaum et al. 2011)  
citeast 2011)

The intent of this thesis is to inject linguistic bias into the machine learning framework of BERT to sharpen the analytical capacities of a masked language model. This is done by altering the

## Bibliography

- Glück, Helmut and Michael Rödel, eds. (2016). *Metzler Lexikon Sprache*. ger. 5th ed. Springer eBook Collection. Stuttgart: J.B. Metzler, Online-Ressource (XXVI, 814 S. 64 Abb., 12 Abb. in Farbe, online resource). 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>.
- Tenenbaum, Joshua, Charles Kemp, Thomas Griffiths, and Noah Goodman (2011). "How to Grow a Mind: Statistics, Structure, and Abstraction". In: *Science* 331.6022, pp. 1279–1285. DOI: 10.1126/science.1192788.