

Exploring Transformer-Based Models for Named Entity Recognition in Ukrainian

Master's Thesis

Timo Junolainen

Februrary 2023

Department of Computer Science
University of Turku

Abstract

Named entity recognition (NER)[2] is an important task in natural language processing, consisting of identifying and classifying named entities in text. With the recent success of transformer-based models for NER in high-resource languages, this thesis explores the effectiveness of these models for NER in a low-resource[1] language: Ukrainian[3]. The research questions addressed in this thesis include the effectiveness of transformer-based models for NER in Ukrainian, the impact of different pre-training and fine-tuning strategies on model performance, and a comparison of the results to existing state-of-the-art methods for NER in Ukrainian. The objectives of this thesis are to develop transformer-based models for NER in Ukrainian, evaluate their performance on standard benchmarks, analyze the impact of different pre-training and fine-tuning strategies, and compare the results to existing state-of-the-art methods. The experimental results show that transformer-based models can achieve competitive results on NER in Ukrainian with appropriate pre-training and fine-tuning strategies. This thesis contributes to the development of NER models for low-resource languages and provides insights into the effectiveness of transformer-based models for NER in such languages.

Contents

Abstract	1
1 Introduction	3
1.1 Motivation behind the work	3
2 Literature Review	5
3 Methodology	6
4 Results	7
5 Conclusions	8

Chapter 1

Introduction

1.1 Motivation behind the work

Named entity recognition (NER)[2] is one of the primary tasks in natural language processing, and consists of classifying named entities in text, such as people, organizations, and locations.

Last years, with use of deep learning methods, with transformers, state-of-the-arts results have been achieved with English, Spanish and Chinese languages (ref ?)

It must be observed, that English, Spanish and Chinese, are high-resource languages, meaning essentially that a lot of computing power went into finetuning models for English, Spanish and Chinese.

However, there are much so called “low-resource” languages, as for example:

1. Basque: spoken in the Basque Country, an autonomous region of northern Spain and southwestern France.
2. Belarusian: spoken in Belarus, a country in Eastern Europe.
3. Fijian: spoken in Fiji, an island country in the South Pacific.
4. Irish: spoken primarily in Ireland and Northern Ireland.
5. Kyrgyz: spoken in Kyrgyzstan, a country in Central Asia.
6. Luxembourgish: spoken in Luxembourg, a small country in Europe.
7. Maltese: spoken in Malta, a small island nation in the Mediterranean.
8. Samoan: spoken in Samoa, an island nation in the South Pacific.

9. Scottish Gaelic: spoken primarily in Scotland.
10. Yoruba: spoken in Nigeria and other West African countries.

Chapter 2

Literature Review

This chapter reviews the literature on NER, with a focus on transformer-based models for NER and existing work on

Chapter 3

Methodology

Chapter 4

Results

Chapter 5

Conclusions

Bibliography

- [1] Felix Laumann. Low resource languages, Jun 2022.
- [2] Wikipedia. Named-entity recognition — Wikipedia, the free encyclopedia. <http://en.wikipedia.org/w/index.php?title=Named-entity%20recognition&oldid=1136823104>, 2023. [Online; accessed 01-March-2023].
- [3] Wikipedia. Ukrainian language — Wikipedia, the free encyclopedia. <http://en.wikipedia.org/w/index.php?title=Ukrainian%20language&oldid=1141828008>, 2023. [Online; accessed 01-March-2023].