



**FACULTY  
OF MATHEMATICS  
AND PHYSICS**  
Charles University

## **MASTER THESIS**

Zuzana Šimečková

# **Entity Relationship Extraction**

Institute of Formal and Applied Linguistics

Supervisor of the master thesis: RNDr. Milan Straka, Ph.D.

Study programme: Computer Science

Study branch: IUI

Prague 2020

This is not a part of the electronic version of the thesis, do not scan!

I declare that I carried out this master thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In ..... date .....  
Author's signature

Dedication.

Title: Entity Relationship Extraction

Author: Zuzana Šimečková

Institute: Institute of Formal and Applied Linguistics

Supervisor: RNDr. Milan Straka, Ph.D., Institute of Formal and Applied Linguistics

Abstract: Abstract.

Keywords: key words

# Contents

<b>Introduction</b>	<b>2</b>
<b>1 Datasets</b>	<b>3</b>
1.1 TACRED dataset . . . . .	3
1.2 SEMEVAL 2010 task 8 dataset . . . . .	3
<b>2 Title of the second chapter</b>	<b>4</b>
2.1 Title of the first subchapter of the second chapter . . . . .	4
2.2 Title of the second subchapter of the second chapter . . . . .	4
<b>Conclusion</b>	<b>6</b>
<b>Bibliography</b>	<b>7</b>
<b>List of Figures</b>	<b>8</b>
<b>List of Tables</b>	<b>9</b>
<b>List of Abbreviations</b>	<b>10</b>
<b>A Attachments</b>	<b>11</b>
A.1 First Attachment . . . . .	11

# Introduction

There has been made noticeable progress in natural language processing since the first deep neural networks attempts. With multiple new approaches and inventions such as multitask learning, word embeddings, RNN, attention and the transformer architecture. Last year Devlin et al. [2018] created BERT and managed to achieve state-of-the-art performance in eleven natural language processing tasks, including GLUE (7.7% point absolute improvement), MultiNLI accuracy (4.6% absolute improvement) and SQuAD problems.

In this thesis, we will try to use those novel approaches to predict relation between two entities based on a Czech sentence. First part of this thesis will be focused on data. We will introduce some existing English datasets for Entity Relation Extraction. Than we will describe how we prepared data for Czech version of this task using distant supervision on Czech Wikipedia and Wikidata. Second part

previous work: Existing work on relation extraction (e.g., Zelenko et al., 2003; Mintz et al., 2009; Adel et al., 2016)

not a sentence

o čem bude druhá část

# 1. Datasets

---

## 1.1 TACRED dataset

The TAC Relation Extraction Dataset was introduced in Zhang et al. [2017].  
Authors claim

training data has often been too noisy for reliable training of relation extraction systems

... machine learning approaches have suffered from two key problems: (1) the models used have been insufficiently tailored to relation extraction, and (2) there has been insufficient annotated data available to satisfy the training of data-hungry models, such as deep learning models.

## 1.2 SEMEVAL 2010 task 8 dataset



## 2. Title of the second chapter

2.1 Title of the first subchapter of the second chapter

2.2 Title of the second subchapter of the second chapter

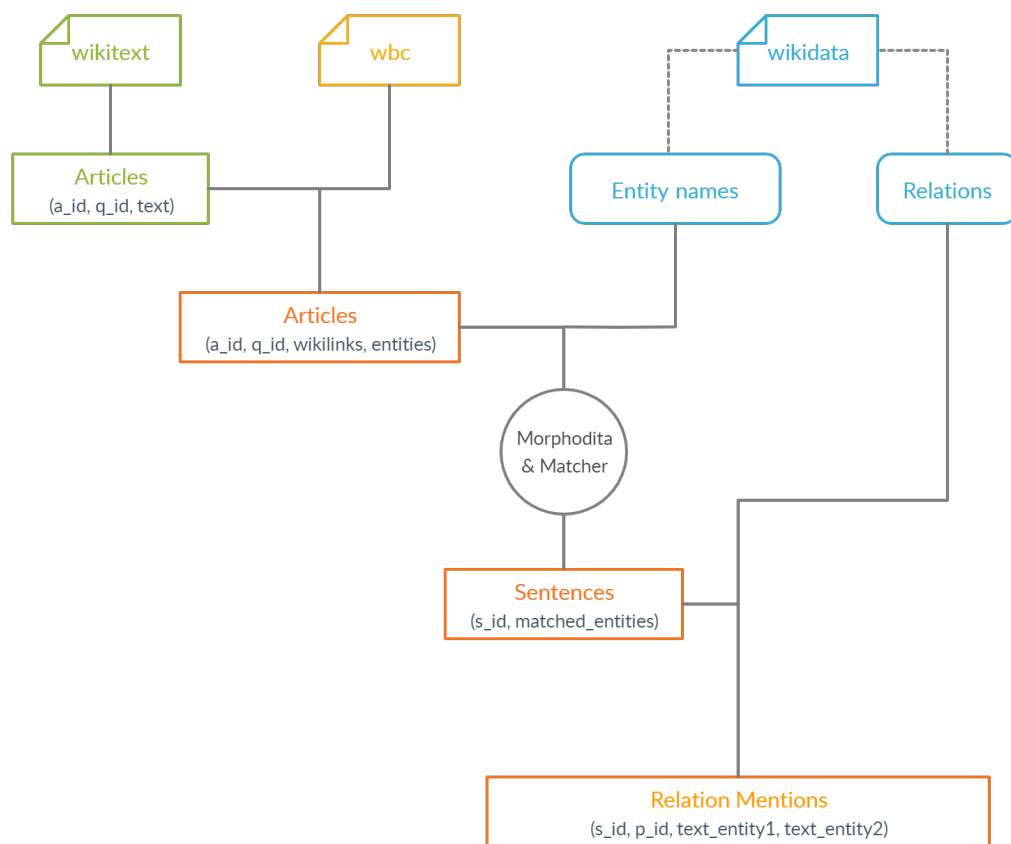


Figure 2.1: Zjednodušený diagram výroby korpusu

# Conclusion

# Bibliography

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. Position-aware attention and supervised data improve slot filling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)*, pages 35–45, 2017. URL <https://nlp.stanford.edu/pubs/zhang2017tacred.pdf>.

# List of Figures

2.1	Zjednodušený diagram výroby korpusu . . . . .	5
-----	---	---

# List of Tables

# List of Abbreviations

# A. Attachments

## A.1 First Attachment