



FACULTY
OF ARTS

Masaryk University

Recent Advancements on the New Online Proofreader of Czech

RASLAN 2019

Vojtěch Mrkývka

mrkyvka@phil.muni.cz

December 6, 2019





Let's state the facts

1. The best proofreader of Czech is part of proprietary system
2. There is no alternative (free of proprietary) providing similarly good results



The project

In 2018 a new project was submitted to and accepted by the Technology Agency of the Czech Republic (TA ČR).

- Goal: New web-based proofreading software
- Started in February 2019, ending in 2022

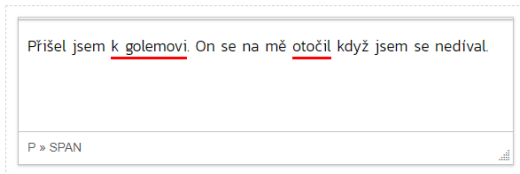
Used tools

- Tools developed in the NLP Center, FI MU
- Rulesets made by students
(bachelor/master theses, doctoral studies)
- Additional external datasets or tools

The current interface



Rozhraní korektoru

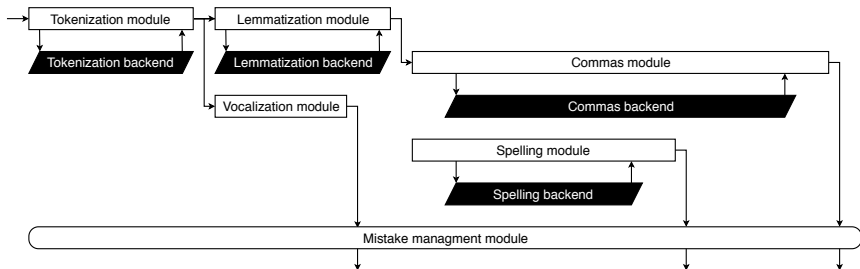




The first version

- Based on on-line text processor TinyMCE
- Series of separate modules
- Javascript & AJAX connection to backend (usually in Python) if necessary

The first version



The pros & cons of the first version

Pros

- Easy development of separate modules
- Little to no bottlenecking due to asynchronous processing

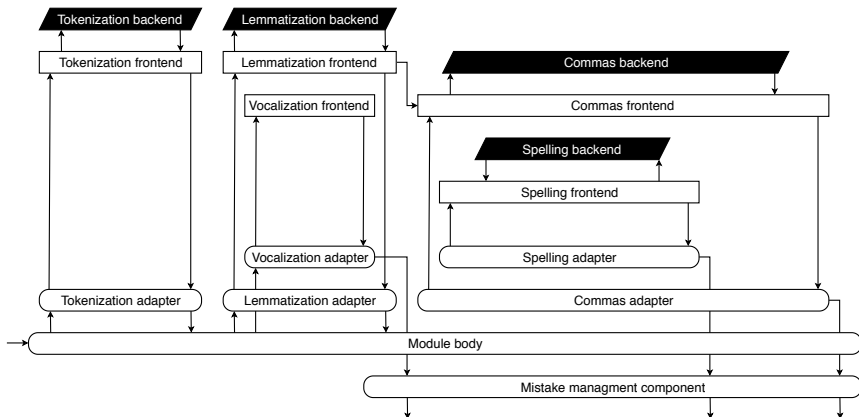
Cons

- Too tightly bound to the single piece of software (TinyMCE)
- No control over potential extensions
- User-unfriendly installation process

The second version

- Proofreader code independent on TinyMCE
- TinyMCE edition contained single module
- Javascript & AJAX connection to backend (usually in Python) if necessary
- Unfinished due to change of the assignment

The second version



The pros & cons of the second version

Pros

- Keeps most of the advantages of the first version
- Strict separation between proofreader and text processor code
⇒ portability
- Improvement in user-friendliness as the proofreader (TinyMCE version) could have been distributed as a single package

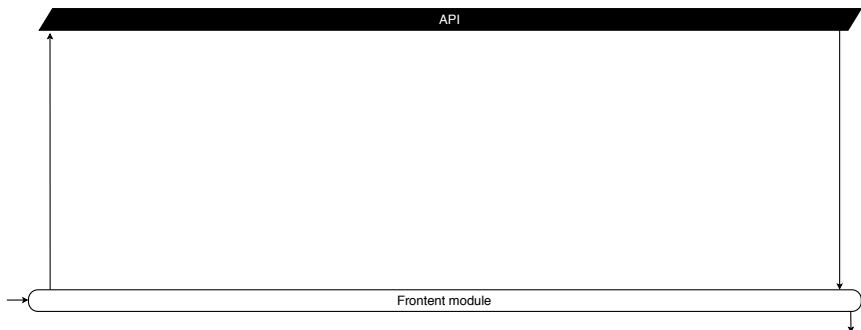
Cons

- Not very suitable for non-web based applications

The third and current version

- Brand new \Rightarrow first preview available from November
- We were addressed by potential customer
 \Rightarrow slight alteration of the assignment
- More suitable for non-web based applications
- Single API written in Python

The third and current version



The pros & cons of the third version

Pros

- Easy to use for web based apps and non-web based apps alike
- Portability improvements due to single input
- Better for collecting statistical data for future improvement
- Easy to be used for automatic testing

Cons

- Single output \Rightarrow possibly lower performance



Currently included and worked on modules

Current modules

- *commas, non-grammatical constructions, "other"*
usage of SET parser for proofreading

Worked on modules

- *spellchecker*
data from Internet Language Manual
(Internetová jazyková příručka, IJP)
- *interpunction*
regular expressions designed by Zbyněk Michálek

Processing speed

- Optimization of used data structures (named tuples instead of dictionaries)
- Asynchronous processing using `asyncio` module – already included
- Multiprocessing using `multiprocessing` module – to be included

Big data testing & alternative tagger

- We are getting testing data from the cooperating news company
- Currently the normalized input is being made for the big data testing to be started
- Optional use of different morphological analysis tools to achieve the best performance
- MorphoDiTa is currently being implemented
- We will surely try the new neural tagger presented earlier by Marie Stará



Thank you for your attention!

This work was supported by the project of specific research *Čeština v jednotě synchronie a diachronie* (Czech language in unity of synchrony and diachrony; project no. MUNI/A/1061/2018) and by the Technology Agency of the Czech Republic under the project TL02000146.