# STYX

## Prague Dependency Treebank as an exercise book of Czech

**Ondřej Kučera**

# Contents

1. **Introduction** (motivation, PDT, implementation)

2. **Filtering sentences**

3. **Transformations of trees**

4. **STYX:** FilterSentences, Charon, Styx

# Motivation

- children of today use computers regularly
  - games, web surfing, chatting, writing, drawing
- why couldn't they parse sentences or determine parts of speech?

# Building an exercise book

## Manually

- extremely hard
  - choose (make up) the sentences
  - annotate them
- considerably limited number of sentences
- often too simple sentences not reflecting the real usage of the language

# Building an exercise book

## Automatically

- if we have annotated data

- the work of choosing the sentences and annotating them is already done

- the data in corpus reflect the real usage of the language

- the number of sentences corresponds to the size of the corpus

- PDT

# Prague Dependency Treebank

- annotated on four layers (word, morphological, analytical, tectogrammatical)
- inner data format: PML (Prague Markup Language) – based on XML

# PDT vs. school syntax

- annotation rules of PDT allow to process any sentence
  ⇨ filtering sentences

- Analytical layer of PDT differs from the school syntax in many ways
  ⇨ transformations of analytical trees

# Filtering sentences

## Filtering in numbers

- nine different filters
- starting number of sentences: 49,442
- after application of the filters: 11,705
- about 23.7% sentences kept

# Transformation of trees

- three basic transformations
- particular transformations consist of
  - combining of the three basic transformations
  - rules for modification of the syntactic functions

# Transformations of trees

## Example

# Implementation

## Java

- high-level language with number of mechanisms protecting programmers "against themselves"

- portability

- presence of SWT library

## SWT

- Standard Widget Toolkit

- provides native look and feel of graphical user interface

- speed

# Implementation

# Implementation

# STYX:
## FilterSentences, Charon, Styx

m-layer

a-layer

t-layer

candidate set

**FilterSentences**

transformations

exercise book

**Charon**

exercise   exercise   exercise

**Styx**

# FilterSentences

- used for applying the filters
- reads data in PML format
- each sentence is tested by a filter
- output data contains the sentences that the filter kept
- output again in PML format

# Charon

- "administrative" program
- loads all sentences available
- the user selects sentences that he or she wants to have in the exercise
- in the end the user saves the exercise

# Charon

# Charon

# Styx

- exercise book itself
- user loads an exercise previously created and saved in Charon
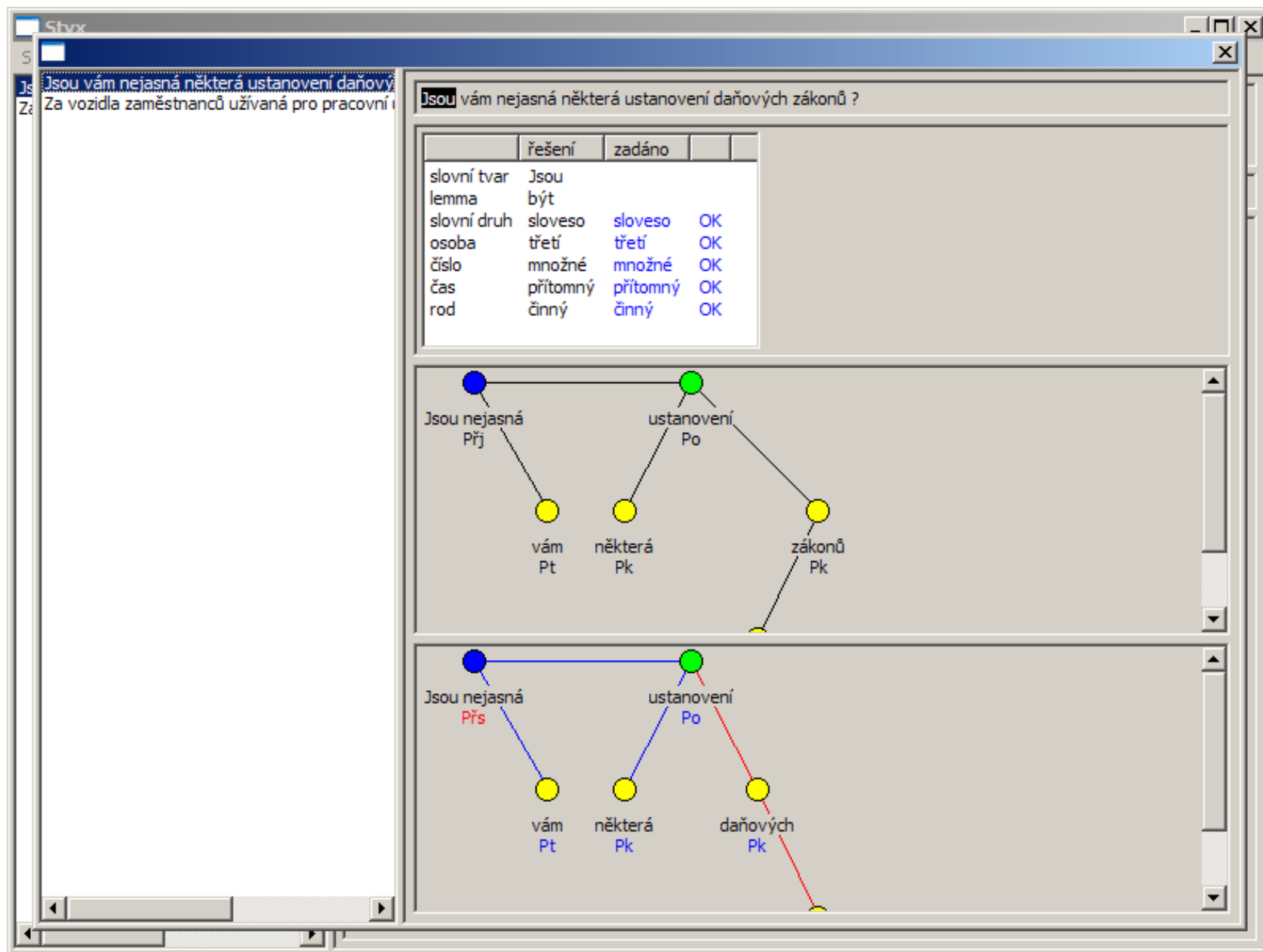
# Styx

# Styx

# Styx

# Questions

and perhaps some answers…