# Annotation of Multiword Expressions in the Prague Dependency Treebank
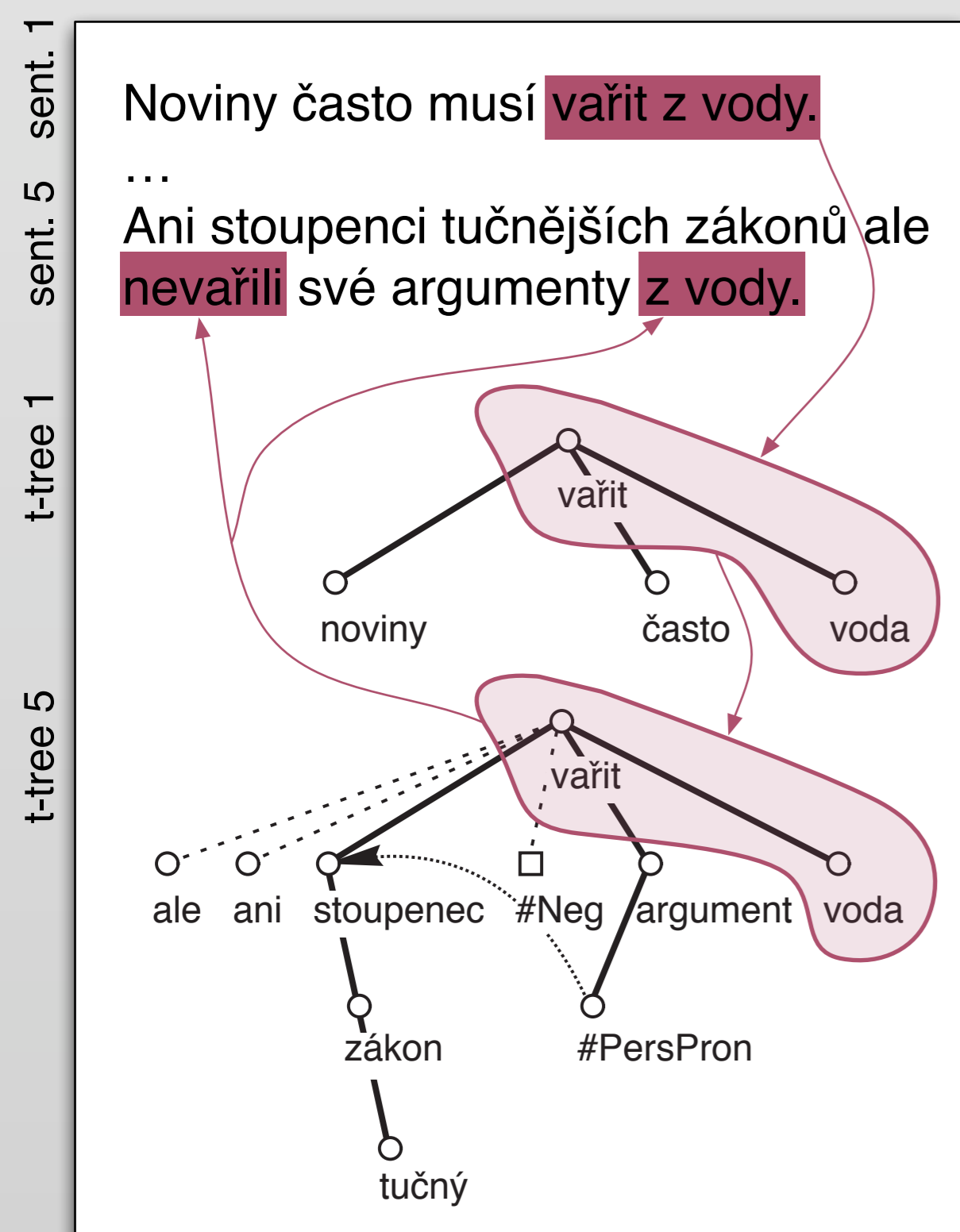
Eduard Bejček and Pavel Straňák

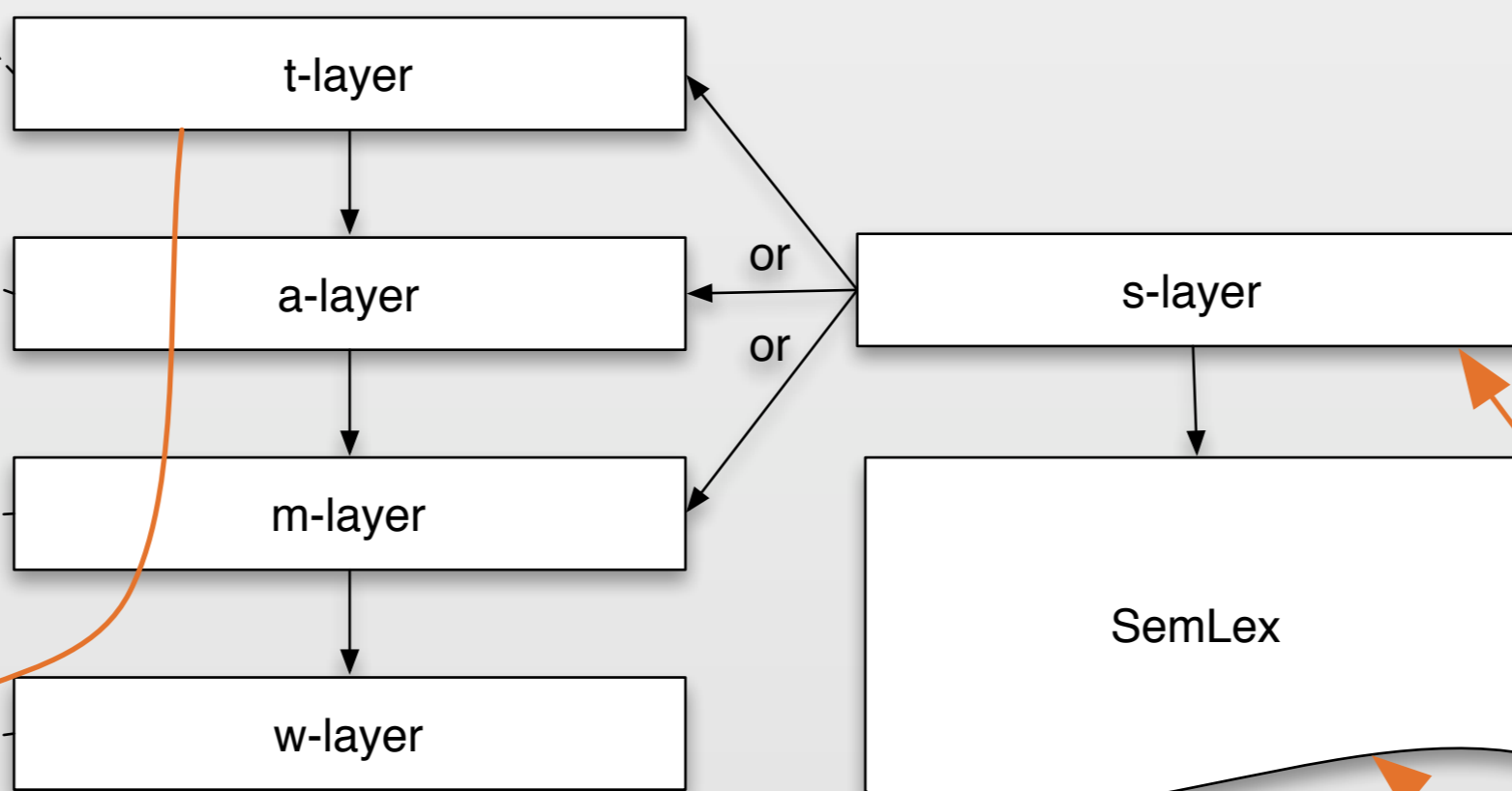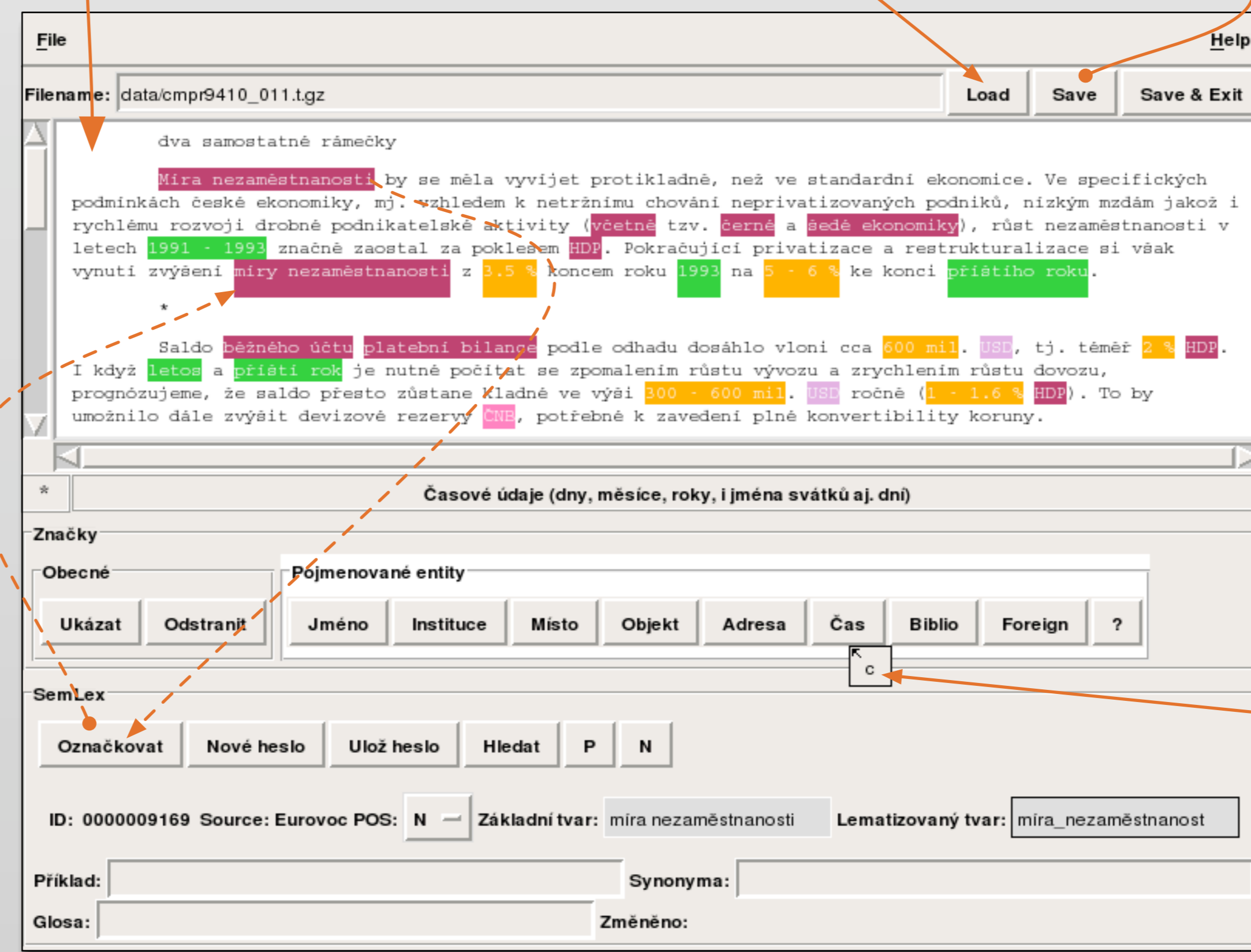We annotate multiword expressions (MWEs) in Prague Dependency Treebank (❶).

We use PML scheme to capture the annotation of whole PDT as well as of our MWEs (❷).

Our annotators work with an annotation tool (❺) and mark occurences in sentences, which are generated from the deep syntactic layer (t-layer) of PDT: see ❶, ❷, ❹.

Our aim is to improve a representation of MWEs (incl. named entities) in t-trees (❸), because current approach is considered to be insufficient (❻).
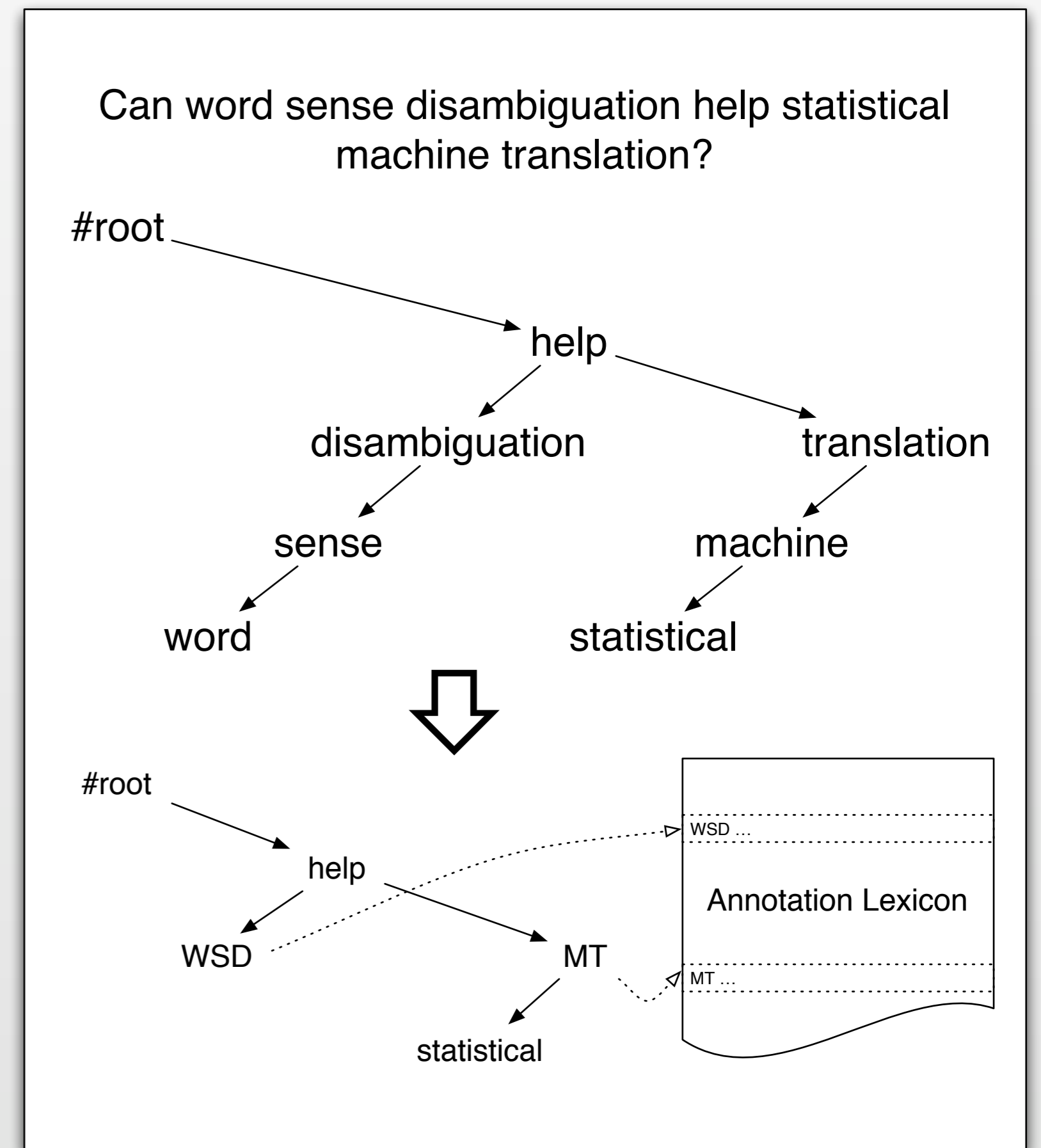


**❶ PDT 2.0**



**❷ PML files and SemLex**



Can word sense disambiguation help statistical machine translation?

**❸ Scheme of changes**



Noviny často musí **vařit z vody**.
…
Ani stoupenci tučnějších zákonů ale **nevařili** své argumenty **z vody**.

Annotate a *new* MWE + automatic pre-annotation of all other occurrences (see ❹)

**❹ Use of t-trees in the pre-annotation**



Load + automatic preannotation (from SemLex)

one-key shortcuts for annotation

searchable SemLex browser and editor

**❺ Annotation interface**



**❻ t-layer in PDT 2.0**

## Inter-annotator Agreement

Each t-node may be: (i) annotated with a SemLex entry (ii) one of nine types of named entities or (iii) not annotated. This yields a scale between full agreement and none. Each type of agreement is assigned a weight according to approximate amount of information it provides (❼).

Then we use the slightly modified pi measure (❽) on these weighted values to compute inter-annotator agreement.

|  | Agreement | | | | Disagreement |
|---|---|---|---|---|---|
|  | Annotated | | | Not annot. | |
|  | Agreement on NE / lexia | |  |  |  |
|  | Full agr. | Disagr. |  |  |  |
| class $c$ | 1 | 2 | 3 | 4 | 5 |
| # of t-nodes $n$ | 10,527 | 2,365 | 389 | 83,287 | 3,988 |
| weight $w$ | 1 | 0.5 | 0.25 | 0.052 | 0 |

**❼ Annotated t-nodes**

## Current State

- We use tectogrammatical tree structures of MWEs for automatic pre-annotation.
  - The richer the tectogrammatical annotation the better the possibilities for automatic pre-annotation, which minimises human errors
- Weighted measure that accounts for partial agreement as well as estimation of maximal agreement
  - The resulting $\pi_w = 0.676$ is statistically significant
  - Agreement should gradually improve as:
    - we clean up the annotation lexicon
    - more entries are pre-annotated automatically
    - and further types of pre-annotation are employed.

$$\pi_w = \frac{A_o - A_e}{\hat{U} - A_e}$$

$$\hat{U} = \frac{n_{A \cup B}}{N} + 0.052 \cdot \frac{N - n_{A \cup B}}{N} = 0.215$$

$$\pi_w = \frac{A_o - A_e}{\hat{U} - A_e} = \frac{0.160 - 0.047}{0.215 - 0.047} = 0.676$$

**❽ Agreement measure**