# **Towards Semantic Annotation Supported by Dependency Linguistics and ILP**

Jan Dědek

Department of Software Engineering, Faculty of Mathematics and Physics,
Charles University in Prague, Czech Republic

Seminář informatiky I2, 2. 11. 2010, MFF UK, Praha
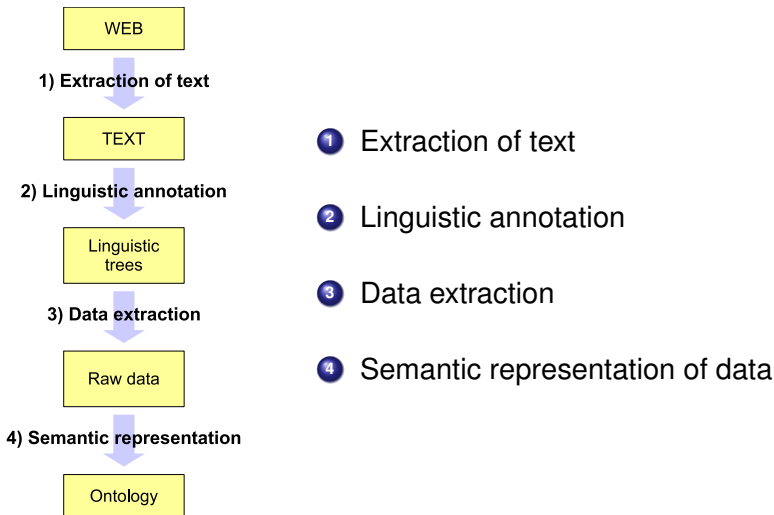Doktorský projekt Res Informatica

**Outline**

Our Information Extraction System
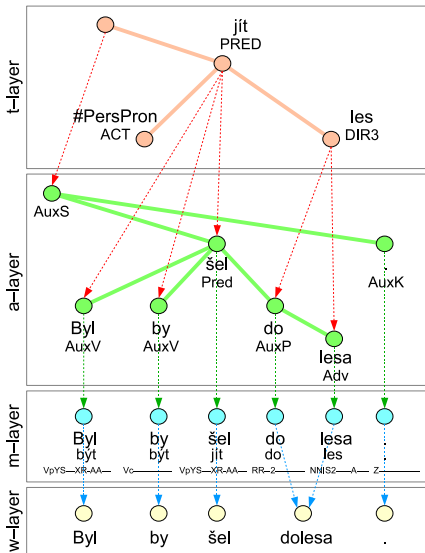
## Introduction to the Presented Work

- Extraction of semantic information from texts.
    - In Czech language.
    - Coming from web pages.
- Using of Semantic Web ontologies.
    - RDF, OWL
- Exploiting of linguistic tools.
    - **Prague Dependency Treebank** project.
    - **TectoMT** project (ÚFAL MFF UK).
    - **GATE** project (The University of Sheffield).
    - Experiments with the **Czech WordNet**.
- Rule based extraction method.
    - Extraction rules $\approx$ linguistic tree queries
    - ILP learning of extraction rules

Our Information Extraction System

# Schema of the extraction process

```
WEB
```

**1) Extraction of text**

```
TEXT
```

**2) Linguistic annotation**

```
Linguistic
trees
```

**3) Data extraction**

```
Raw data
```

**4) Semantic representation**

```
Ontology
```

**❶** Extraction of text

**❷** Linguistic annotation

**❸** Data extraction

**❹** Semantic representation of data

Linguistics we have used

## Layers of linguistic annotation in PDT



- Tectogrammatical layer
- Analytical layer
- Morphological layer

- PDT 2.0 on-line:

http://ufal.mff.cuni.cz/pdt2.0/

*Sentence:*

Byl by šel dolesa.
He-was would went toforest.

## Tools for machine linguistic annotation

1. Segmentation and tokenization
2. Morphological analysis
3. Morphological tagging
4. McDonnald's Maximum Spanning Tree parser
   – Czech adaptation
5. Analytical function assignment
6. Tectogrammatical analysis
   – Developed by Václav Klimeš

- Available within the TectoMT[1] project

---

[1] http://ufal.mff.cuni.cz/tectomt/

Introduction
○○○○●○○○○

Our Information Extraction Method
○○○○○○○○○○○○○

Conclusion
○○

Linguistics we have used

## Example of tectogrammatical tree



- Lemmas
- Functors
- Semantic parts of speech

*Sentence:*

Ve zdemolovaném trabantu na místě zemřeli dva muži – 82letý senior a další muž, jehož totožnost zjišťují policisté.

Two men died on the spot in demolished trabant – . . .

Domain of fire-department articles

# Example of the web-page with a report of a fire department

Domain of fire-department articles

## Domain of our experiments

- Fire-department articles
- From The Ministry of Interior of the Czech Republic[2]
- Extensive experiments
  - More than 800 articles
  - 1.2 MB of textual data
  - Extracting information about injured and killed people
  - 470 matches of the extraction rule,
    200 numeric values of quantity (described later)
- Intensive experiments
  - 50 articles
  - Precisely manually tagged
  - Used of the evaluation of the learning procedure

---

[2]`http://www.mvcr.cz/rss/regionhzs.html`

Domain of fire-department articles

**Example of processed text**



- Information to be extracted is decorated.
- See the last sentence on the next slide.

Domain of fire-department articles

# Example of a linguistic tree



…, škodu vyšetřovatel předběžně vyčíslil na osm tisíc korun.

…, investigating officer preliminarily reckoned the damage to be 8 000 CZK.

- Our IE method uses tree queries (tree patterns)

Introduction
○○○○○○○○○

Our Information Extraction Method
●○○○○○○○○○○○○○

Conclusion
○○

Manually created rules

T-jihomoravsky49640.txt-001-p1s4
root

zemřít
PRED
v

#PersPron
ACT
n.pron.def.pers

Trabant
LOC.basic
n.denot

#Dash
APPS
coap

zdemolovaný
RSTR
adj.denot

místo
LOC.basic
n.denot

muž
ACT
n.denot

a
CONJ
coap

dva
RSTR
adj.quant.def

senior
DENOM
n.denot

muž
DENOM
n.denot

82letý
RSTR
adj.denot

další
RSTR
adj.denot

zjišťovat
RSTR
v

totožnost
ACT
n.denot.neg

policista
ACT
n.denot

který
APP
n.pron.indef

... two ...

- How to extract the information about two dead people?

Manually created rules

# Extraction rules – Netgraph queries



- Tree patterns on shape and nodes (on node attributes).
- Evaluation gives actual matches of particular nodes.
- Names of nodes allow use of references.

## Raw data extraction output

```
<QueryMatches>
  <Match root_id="T-vysocina63466.txt-001-p1s4" match_string="2:0,7:3,8:4,11:2">
    <Sentence>
      Při požáru byla jedna osoba lehce zraněna - jednalo se
      o majitele domu, který si vykloubil rameno.
    </Sentence>
    <Data>
      <Value variable_name="action_type" attribute_name="t_lemma">zranit</Value>
      <Value variable_name="injury_manner" attribute_name="t_lemma">lehký</Value>
      <Value variable_name="participant" attribute_name="t_lemma">osoba</Value>
      <Value variable_name="quantity" attribute_name="t_lemma">jeden</Value>
    </Data>
  </Match>
  <Match root_id="T-jihomoravsky49640.txt-001-p1s4" match_string="1:0,13:3,14:4">
    <Sentence>
      Ve zdemolovaném trabantu na místě zemřeli dva muži - 82letý senior
      a další muž, jehož totožnost zjišťují policisté.
    </Sentence>
    <Data>
      <Value variable_name="action_type" attribute_name="t_lemma">zemřít</Value>
      <Value variable_name="participant" attribute_name="t_lemma">muž</Value>
      <Value variable_name="quantity" attribute_name="t_lemma">dva</Value>
    </Data>
  </Match>
  <Match root_id="T-jihomoravsky49736.txt-001-p4s3" match_string="1:0,3:3,7:1">
    <Sentence>Čtyřiatřicetiletý řidič nebyl zraněn.</Sentence>
    <Data>
      <Value variable_name="action_type" attribute_name="t_lemma">zranit</Value>
      <Value variable_name="a-negation" attribute_name="m/tag">VpYS---XR-(N)A---
      </Value>
      <Value variable_name="participant" attribute_name="t_lemma">řidič</Value>
    </Data>
  </Match>
</QueryMatches>
```
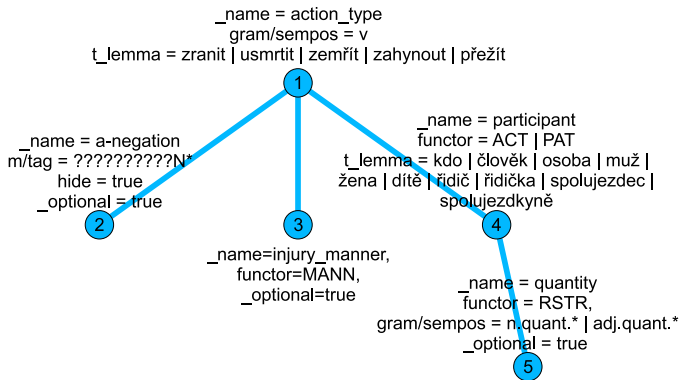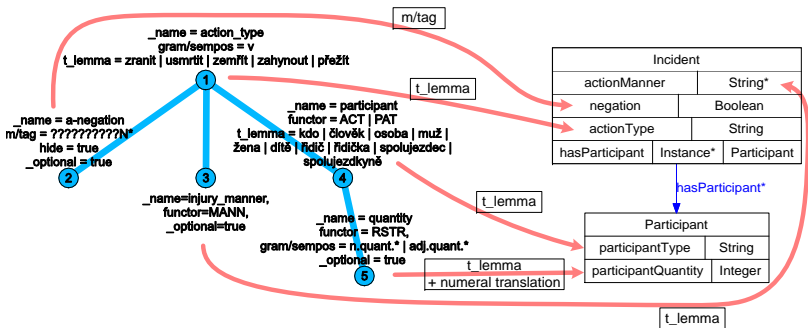
**SELECT** action_type.t_lemma, a-negation.mtag, injury_manner.t_lemma,

participant.t_lemma, quantity.t_lemma **FROM** *\*\*\*extraction rule\*\*\**

# Semantic interpretation of extraction rules



- Determines how particular values of attributes are used.
- Gives semantics to extraction rule.
- Gives semantics to extracted data.

Introduction
○○○○○○○○○

Our Information Extraction Method
○○○○○○○●○○○○○○

Conclusion
○○

Learning of rules

Learning of rules

# Integration of ILP in our extraction process



- Transformation of trees to logic representation.
- Today: just first promising experiments.

Learning of rules

# Logic representation of linguistic trees



```
tree_root(node0_0). node(node0_0).
id(node0_0, t_jihomoravsky49640_txt_001_p1s4).
$$$$$$$$ node0_1 $$$$$$$$$$$$$$$$$
node(node0_1).
functor(node0_1, pred).
gram_sempos(node0_1, v).
t_lemma(node0_1, zemrit).
$$$$$$$$ node0_2 $$$$$$$$$$$$$$$$$
node(node0_2).
functor(node0_2, act).
gram_sempos(node0_2, n_pron_def_pers).
t_lemma(node0_2, x_perspron).
$$$$$$$$ node0_3 $$$$$$$$$$$$$$$$$
node(node0_3). id(node0_3,
functor(node0_3, loc).
gram_sempos(node0_3, n_denot).
t_lemma(node0_3, trabant).
...

edge(node0_0, node0_1). edge(node0_1, node0_2).
edge(node0_1, node0_3). edge(node0_3, node0_4).
edge(node0_4, node0_5). edge(node0_3, node0_6).
edge(node0_3, node0_7). edge(node0_3, node0_8).
...
```

Logic representation

Source web page

... two ...

Linguistic trees

Introduction
ooooooooooo

Our Information Extraction Method
oooooooooooooo

Conclusion
oo

Learning of rules

# Root/Subtree Preprocessing/Postprocessing (Chunk learning)



reckon

thousand

Root

CZK

Sub-tree

PRED
vyčíslit
v

EFF
#PersPron
n.pron.def.pers

PAT
škoda
n.denot

ACT
vyšetřovatel
n.denot

MANN
předběžně
adj.denot

CPR
tisíc
n.quant.def

damage

investigating officer

eight

RSTR
osm
adj.quant.def

MAT
koruna
n.denot

…, škodu vyšetřovatel předběžně vyčíslil na osm tisíc korun.

…, investigating officer preliminarily reckoned the damage to be eight thousand Crowns (CZK).

**Examples of learned rules, Czech words are translated.**

### Example

[Rule 1] [Pos cover = 14 Neg cover = 0]
damage_root(A) :- lex_rf(B,A), has_sempos(B,'n.quant.def'),
   tDependency(C,B), tDependency(C,D),
   has_t_lemma(D,'investigator').

[Rule 2] [Pos cover = 13 Neg cover = 0]
damage_root(A) :- lex_rf(B,A), has_functor(B,'TOWH'),
   tDependency(C,B), tDependency(C,D), has_t_lemma(D,'damage').

[Rule 1] [Pos cover = 7 Neg cover = 0]
injuries(A) :- lex_rf(B,A), has_functor(B,'PAT'),
   has_gender(B,anim), tDependency(B,C), has_t_lemma(C,'injured').

[Rule 8] [Pos cover = 6 Neg cover = 0]
injuries(A) :- lex_rf(B,A), has_gender(B,anim), tDependency(C,B),
   has_t_lemma(C,'injure'), has_negation(C,neg0).

Introduction
00000000

Our Information Extraction Method
000000000000●

Conclusion
00

Evaluation

**Evaluation results**

| task/method | matching | missing | excess | overlap | prec.% | recall% | F1.0% |
|---|---|---|---|---|---|---|---|
| **damage/ILP** | 14 | 0 | 7 | 6 | 51.85 | 70.00 | 59.57 |
| **damage/ILP – lenient measures** | | | | | 74.07 | 100.00 | 85.11 |
| **dam./ILP-roots** | 16 | 4 | 2 | 0 | 88.89 | 80.00 | 84.21 |
| **damage/Paum** | 20 | 0 | 6 | 0 | 76.92 | 100.00 | 86.96 |
| **injuries/ILP** | 15 | 18 | 11 | 0 | 57.69 | 45.45 | 50.85 |
| **injuries/Paum** | 25 | 8 | 54 | 0 | 31.65 | 75.76 | 44.64 |
| **inj./Paum-afun** | 24 | 9 | 38 | 0 | 38.71 | 72.73 | 50.53 |

- 10-fold cross validation
- Two tasks: 'damage' and 'injuries'
- Root/subtree preprocessing/postprocessing used for 'damage' task

## Summary

- Proposed a system for extraction of semantic information
- Based on third party linguistic tools (TectoMT[3])
- Extraction rules adopted from Netgraph[4] application.
- ILP used for learning rules.
- All methods integrated inside GATE[5].

- Our future research will concentrate on:
    - Extension of the method with WordNet technology.
    - Adaptation of this method on other languages.
    - Evaluation of the method on other datasets.

---

[3] http://ufal.mff.cuni.cz/tectomt/
[4] http://quest.ms.mff.cuni.cz/netgraph/
[5] http://gate.ac.uk/

## Inter-project cooperation

- Mainly with **I-3 Matematická lingvistika** group
- Directly with David Mareček
- Indirectly with all participants of the linguistic projects

- Making all the linguistic tools working, updates
- Planned feedback in the future
  (comparison of linguistic tools)
- Porting PDT formalism and TMT tools to the GATE platform
  (side effect of our work, but may be usefull in the future).