

# Treebanking à la carte<sup>1</sup>

Milena Hnátková Petr Jäger Tomáš Jelínek  
Vladimír Petkevič Alexandr Rosen Hana Skoumalová

Institute of Theoretical and Computational Linguistics  
Faculty of Arts, Charles University, Prague

Dubrovnik, 12th – 14th September 2011

<sup>1</sup>This paper was financially supported by the Grant Agency of the Czech Republic  
No. P406/10/0434

# Outline of the talk

- 1 Introduction
- 2 Main features
- 3 Architecture
- 4 Examples
- 5 Input text processing
- 6 Summary: main assets

# Outline of the talk

- 1 Introduction
- 2 Main features
- 3 Architecture
- 4 Examples
- 5 Input text processing
- 6 Summary: main assets

## Different views of syntactic structure

- Syntax is a discipline of many theories
- Syntactically annotated corpus runs the risk of a theoretical bias
- Theory-specific representations have different appearances but share a large part of content
- Treebank offering different views of a single syntactic annotation is a realistic and appropriate goal

## From scratch? A bad idea!

- Enormous manual efforts went into building treebanks already [Hajič(2006), Hajič et al.(1998), Skut et al.(1997), i.a.]
- Scaling-up possible by automatic tools.

## Improvements over existing annotation schemes

- Potentially underspecified morphological and syntactic core
- Multiple interaction shells, customisable in shape and detail according to the preferences of humans or computer applications
- Accessible to lay users and satisfying experts at the same time

# Outline of the talk

- 1 Introduction
- 2 Main features**
- 3 Architecture
- 4 Examples
- 5 Input text processing
- 6 Summary: main assets

## Syntactic structure

- Internal skeleton structures: constituency-based, with a combination of **binary** and **flat** branching
- Interpretable as **constituency** or **dependency** trees, according to users' specification, visualized with an arbitrary amount of detail, not necessarily by tree graphs
- Surface and deep structure encoded within a single structure: constituents are labelled as **syntactic functions** (including heads as special functions)
- Heads are further specified as **deep** or **surface**
  - ▶ **Deep head**: deep syntactic governor: *bylo by se to povedlo*
  - ▶ **Surface head**: can be identical to the deep head or different: auxiliary *být*, prepositions, subordinate conjunctions, numerals

## Three levels

- Word order and syntactic structure as distinct dimensions, each sentence is represented at three inter-linked levels:
  - ▶ **graphemics** (orthographic words, contractions)
  - ▶ **morphology** (syntactic words, including haplogitized items)
  - ▶ **syntax** (trees, no nodes for pro-dropped subjects)



## Syntactic phenomena

Annotation of:

- **Agreement of various types**
- **Compound periphrastic verbal forms** (passives, conditional structures, future...)
- **Grammatical co-reference** (grammatical control, relative/reflexive pronouns, predicative complements)
- **Multi-word units** (collocations)

## Expressive power

- Expressive enough to accommodate analyses of arbitrary granularity
- Ambiguous or undecidable phenomena represented by **underspecification** and **distributive disjunction**
- Annotation of any kind can be missing, a sentence may be a mere list of words

## Specifications

- Annotation must be licensed by a formal grammar. Words and constituents have their appropriate (potentially **underspecified**) sets of features
- Lexicons are used to index forms, syntactic words and compound forms
- Customizable visualizations are enabled by formal definitions

# Outline of the talk

- 1 Introduction
- 2 Main features
- 3 Architecture**
- 4 Examples
- 5 Input text processing
- 6 Summary: main assets

## Representation layers:

- **text**
- **morphology**
- **syntactic structure**

## Lexicons:

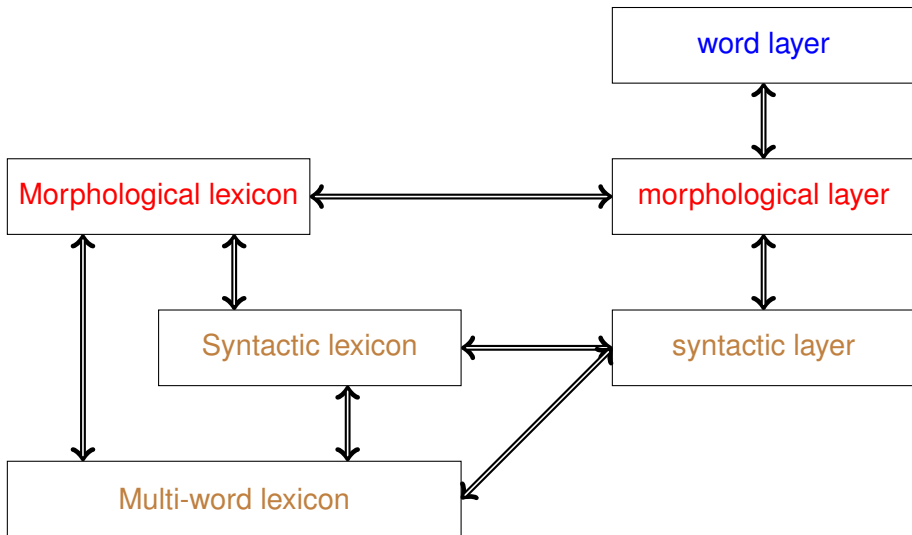
- **morphological**
- **syntactic**
- **multi-word expressions**

## Two-way links between layers, and between layers and lexicons

- to link information across the layers
- to provide lexeme-specific information
- to identify multi-word expressions, including periphrastic forms

## Links within a tree

- **Agreement**
- **Compound (multi-word) verbal predicates**
- **Grammatical coreference**
- ...



## Word layer

- tokenized, including punctuation and MWE: *česko-slovenský*
- contractions left intact (not interpreted): *očs, ses*

## Morphological layer

- morphological analysis and lemmatization of all forms
- contractions split (i.e. interpreted): *očs* → *o co jsi*
- some punctuation marks glued back with word forms:  
*atd., česko-*

## Syntactic layer

- **constituency-based structure**, representable according to user's options/specifications
- punctuation omitted

## Syntactic structure

Syntactic structure is represented by a constituency-based tree where:

- each nonterminal node is assigned a **type** & a **syntactic function**
- each terminal node is assigned a **syntactic function**



## Hierarchy of types

- **TypeHeaded**
- **TypeUnHeaded**
  - ▶ **Coord** – coordination
  - ▶ **Adord** – adordination
  - ▶ **Unspec** – unspecified (for collocations and other)

## Syntactic functions

Special functions for TypeHeaded:

- **SurfHead** – surface head: auxiliary *být/bývat*, prepositions, subordinate conjunctions, numerals in quantified expressions: *pět dětí*
- **DeepHead** – in case it differs from SurfHead (head nouns in PPs, autosemantic verbs in analytical predicates...)
- **Head** – both **SurfHead** & **DeepHead**

## Syntactic functions – continued

Other functions for TypeHeaded:

- **Subj** – subject
- **Attr** – attribute
- **Obj-Advb**
  - ▶ **Obj**
  - ▶ **Advb**
- **VbAttr** – predicative complement
- **RefITant** – reflexive element (*si*, *se*) for inherent reflexives
- **Deagent** – deagentive reflexive
- **Apos** – apposition
- **InDep** – independent syntactic element (parenthesis, vocative syntactic noun...)

## Syntactic functions – continued

Special function for TypeUnHeaded structures:

- **Memb** – member of a TypeUnHeaded structure

## Morphological lexicon

- list of lemmas with inflection paradigms
- a lemma is introduced if two words differ in morphological paradigms—not only in syntactic properties or in semantics:
  - ▶ *travička* ‘little grass’/‘female poisoner’ has only **one lemma**,
  - ▶ *člen* ‘member’ has two lemmas, as it is either **masculine animate** or **masculine inanimate**.

## Syntactic lexicon

- list of lemmas with their syntactic properties
- inherent reflexives have separate entries
- *rozhodnout* and *rozhodnout se* ‘decide’ are two separate entries
- *vidět* ‘see’ is one entry
- **valency frame** entries
- different valency frames listed under one lemma
- ...

## Multi-word lexicon

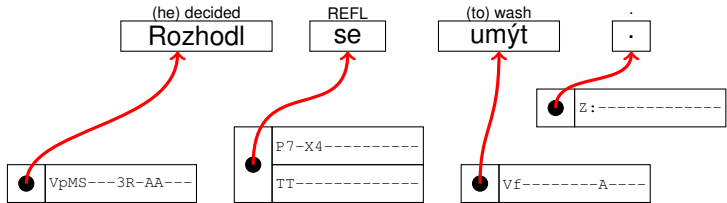
- collocations: *křížem krážem* “in all directions”, *nechat na holičkách* “leave in the lurch”
- types of analytical verb forms: *bych byl přišel* “(I) would have come”, *jsi přišel* “(you) have come”
- inherent reflexives: *usmíváš se* “(you) smile”; *rozhodne se* “(he) decides”
- (reflexive) deagentive construction: *jde se* “let’s go”
- reflexive passives: *bábovka se peče* “the cake is baking/being baked”
- analytical passives: *je čten* “is read”
- nominal predicates: *je velký* “is big”
- agreement patterns: *Lucie ho viděla opilého* “Lucie saw him drunk”

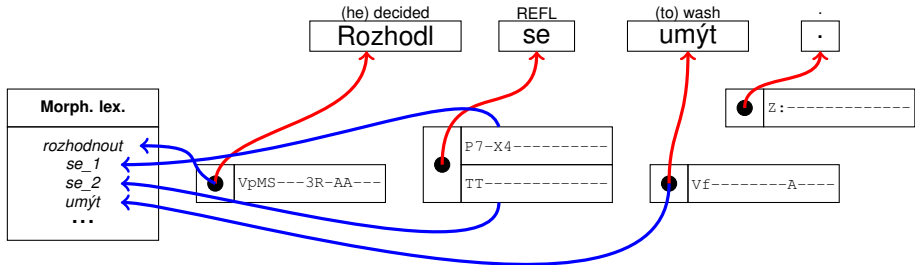
# Outline of the talk

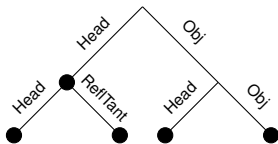
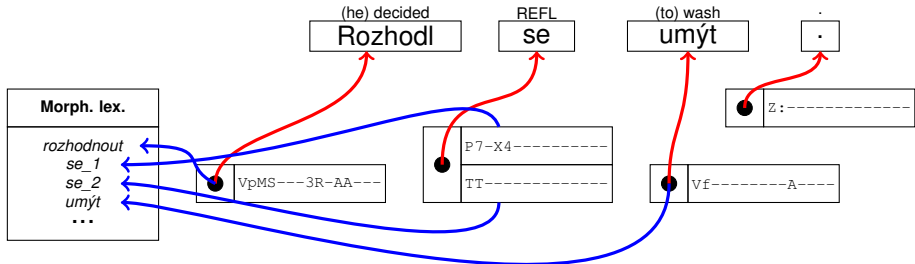
- 1 Introduction
- 2 Main features
- 3 Architecture
- 4 Examples**
- 5 Input text processing
- 6 Summary: main assets

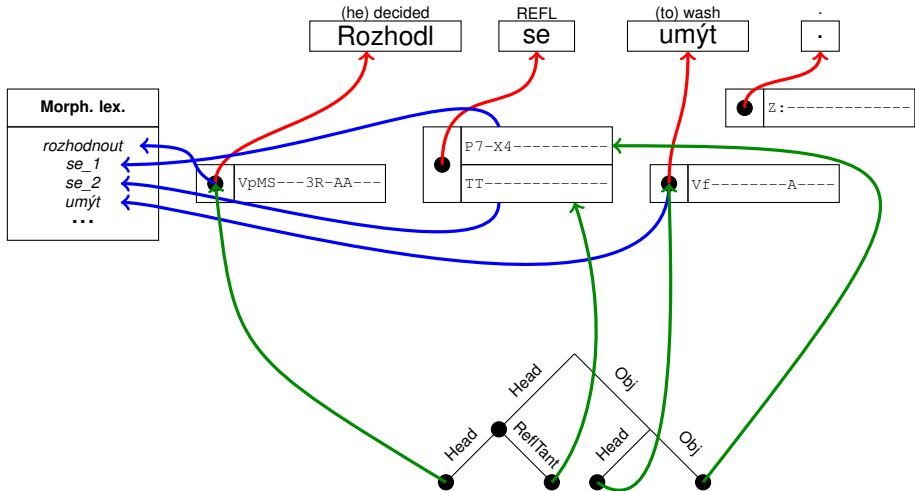


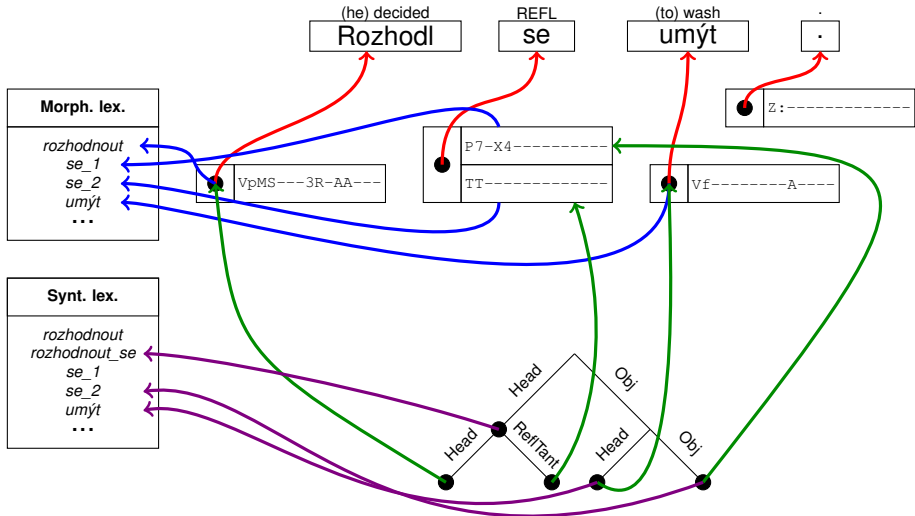
(he) decided	REFL	(to) wash	.
Rozhodl	se	umýt	.

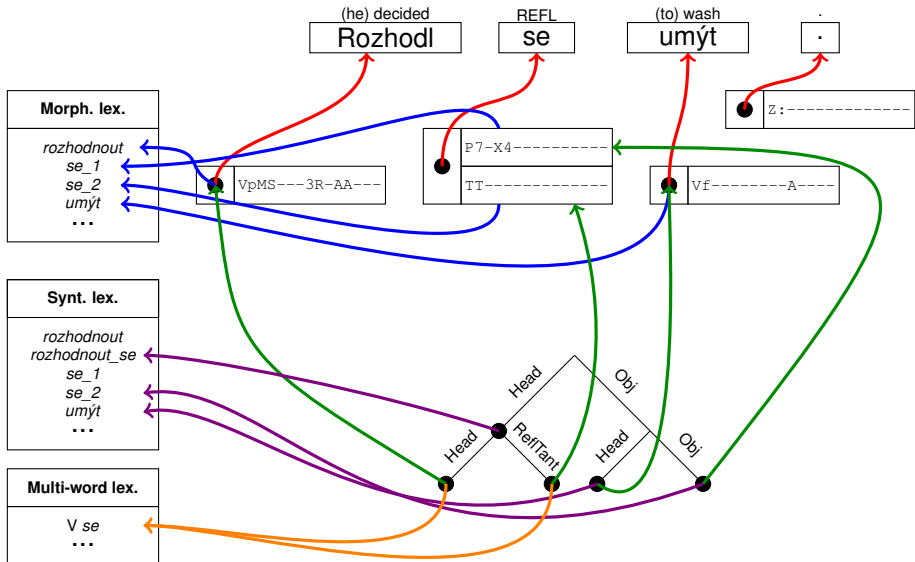


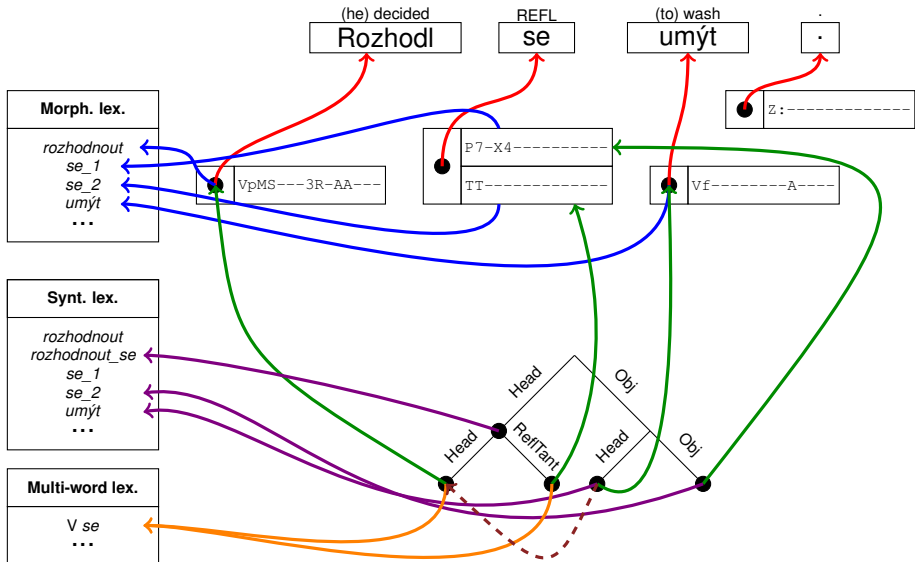












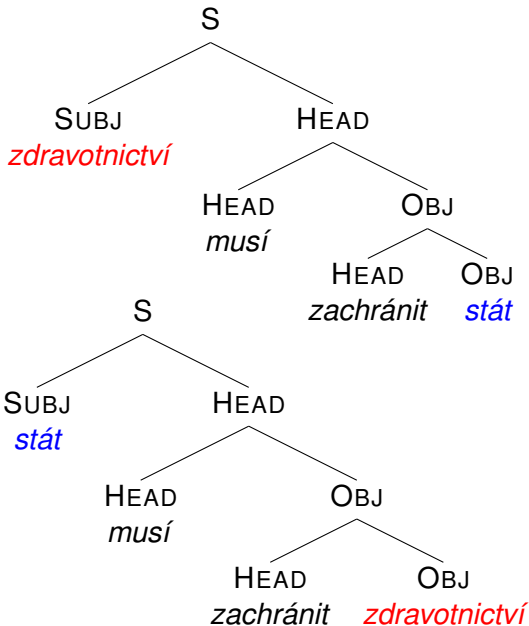


# Subject/object ambiguity

- (1) **Zdravotnictví** musí zachránit **stát**.  
**health service**<sub>nom/acc</sub> must save **state**<sub>nom/acc</sub>

## Two different readings:

- #1 Health service must save the State.
- #2 Health service must be saved by the government.



*Morphological analysis* of (1) with some values unspecified:

- ① zdravotnictví    *noun*, CASE=*X*, NUM=*sg*, GEND=*n*
- ② musí            *verbfin*, PERS=*3*, NUM=*sg*
- ③ zachránit       *verbinf*
- ④ stát             *noun*, CASE=*Y*, NUM=*sg*, GEND=*m*

*Constituents* in one of the two possible syntactic structures of (1), some boxed numbers refer to the forms above:

- ⑤ [ ③zachránit ④stát ]
- ⑥ [ ②musí ⑤ ]
- ⑦ [ ①zdravotnictví ⑥ ]

Two possible structures with constraints on category values and overriding clauses:

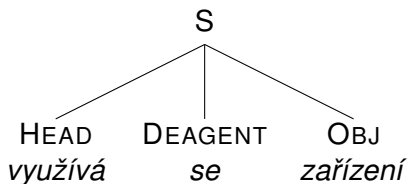
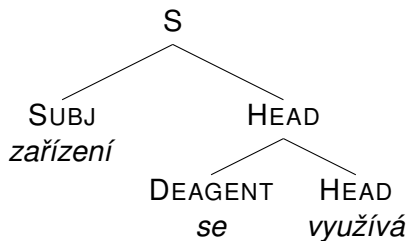
#1 = ⑦, *X=nom*, *Y=acc*

#2 = ⑦, *X=acc*, *Y=nom*, ① → ④, ④ → ①

## Another type of subject/object ambiguity

Reflexive passive:

- (2) Zařízení<sub>Nom/Gen</sub> se využívá.  
 device REFL uses  
 'The device is being used.'

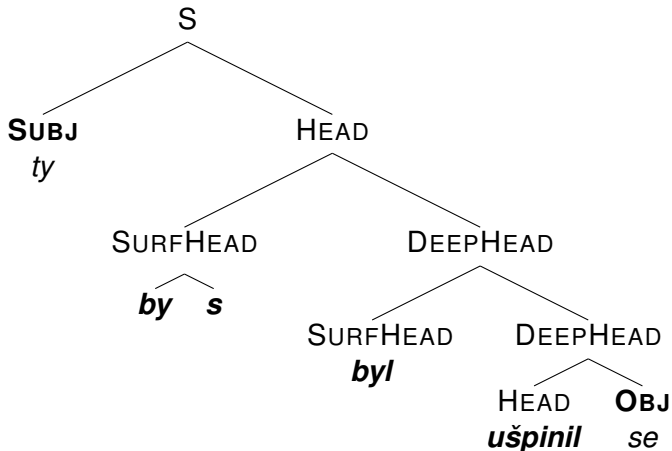


# Treating contractions

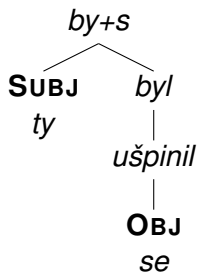
- (3) **Ty** by **ses** byl ušpinil.  
 you would REFL+AUX<sub>2nd,sg</sub> be<sub>ppl</sub>e get dirty<sub>ppl</sub>e  
 ‘You would have got dirty.’

**Ty** by **ses** byl ušpinil.

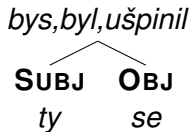
(4)



- (5) **Surface dependency structure** derived from (4)



- (6) **Deep dependency structure** derived from (4)



# Outline of the talk

- 1 Introduction
- 2 Main features
- 3 Architecture
- 4 Examples
- 5 Input text processing**
- 6 Summary: main assets

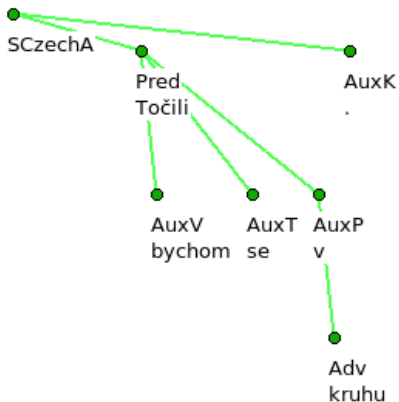


## Input text processing:

- **tokenization** and **sentence segmentation** → sentence on the **graphemic/word level**
- **morphological analysis** and **disambiguation** (rule-based disambiguation, MorČe, collocation module) → sentence on the **morphological level**
- **McDonald et al.**'s parser (or more parsers with a voting scenario) applied to the disambiguated sentence; McDonald et al.'s parser can be parameterized
- **automatic correction of the parse** still in `tectomt` format
- **conversion** of the corrected parse from the `tectomt` format to ours + modifications:
  - ▶ phenomena that in a dependency tree require arbitrary decisions: constructions with auxiliary verbs, coordinated constructions, lists
  - ▶ **disjunction** accounting for structural ambiguities expressed by combined functions **AttrAdv**, **ObjAdv**

# PDT representation

- (7) Točili bychom se v kruhu.  
 Turn would<sub>1st,pl</sub> REFL in circle.



# Internal structure

Ling Viewer

File Transform Help

Točili bychom se v kruhu .

Show

Horizontal

Folder structure

Basic variant

morphology

lemma

functions

variants

structures

DeepHead

SurfHead  
bychom  
být  
Vc-P---1-----I

Head

Head  
Točili  
točit  
VpMP---1R-AA---P

RefTant  
se  
se  
P7-X4-----

Advb

SurfHead  
v  
v  
RR--6-----

DeepHead  
kruhu  
kruh  
NNNS6-----A

# Surface structure

**Ling Viewer**

File Transform Help

Točili bychom se v kruhu .

Show

Horizontal

Surface structure

Basic variant

morphology

lemma

functions

variants

structures

```

graph TD
    Root["bychom  
být  
Vc-P---1-----I"]
    Root --- Node1["Točili  
točit  
VpMP---1R-AA---P"]
    Node1 --- Node2["se  
se  
P7-X4-----"]
    Node1 --- Node3["v  
v  
RR--6-----"]
    Node3 --- Node4["kruhu  
kruh  
NNNS6-----A-----"]
  
```

bychom  
být  
Vc-P---1-----I

DeepHead  
Točili  
točit  
VpMP---1R-AA---P

RefTant  
se  
se  
P7-X4-----

Advb  
v  
v  
RR--6-----

DeepHead  
kruhu  
kruh  
NNNS6-----A-----

# Deep structure

**Ling Viewer** File Transform Help

Točili bychom se v kruhu .

Show

Horizontal

Deep Structure

Basic variant

morphology

lemma

functions

variants

structures

Točili  
točit  
VpMP---1R-AA---P

RefITant  
se  
se  
P7-X4-----

Advb  
kruhu  
kruh  
NNN56-----A-----

SurfHead  
bychom  
být  
Vc-P---1-----I

SurfHead  
v  
v  
RR--6-----

# Outline of the talk




- 1 Introduction
- 2 Main features
- 3 Architecture
- 4 Examples
- 5 Input text processing
- 6 Summary: main assets**

## Main assets

- Exploitation of existing approaches and tools, we cannot develop the treebank from scratch
- multilayer stand-off annotation
- linguistically motivated corrections of the results of the tools used
- entirely automatic processing

## Comparison with analytical layer of PDT

- Single annotation capturing syntactic core, as little theoretical bias as possible but:
  - ▶ various interpretations using underspecification, disjunction
  - ▶ export options into customizable formats
  - ▶ various visualizations of the data

-  Hajič, J. (2006).  
Complex Corpus Annotation: The Prague Dependency Treebank.  
  
In M. Šimková, editor, Insight into the Slovak and Czech Corpus Linguistics, pages 54–73, Bratislava, Slovakia. Veda.
-  Hajič, J., Hajičová, E., Panevová, J., & Sgall, P. (1998).  
Syntax v Českém národním korpusu.  
Slovo a slovesnost, **59**(3), 168–177.
-  Skut, W., Krenn, B., Brants, T., & Uszkoreit, H. (1997).  
An annotation scheme for free word order languages.  
In Proceedings of the Fifth Conference on Applied Natural Language Processing ANLP-97, Washington, DC.