# Treex:
# Modular NLP Framework

## Martin Popel

ÚFAL (Institute of Formal and Applied Linguistics)
Charles University in Prague

September 2015, Prague, MT Marathon

# Outline

**Treex**

- Motivation, Treex vs. TectoMT
- Treex architecture
- Treex internals
- Future plans
- Conclusion and examples

# Motivation

**Treex**

## Goals of Treex

- elegant integration of in-house and third-party NLP tools

- modularity, reusability, cooperation

- ability to easily modify and add code in a full-fledged programming language (Perl)

# Treex vs. TectoMT

2005 (Zdeněk Žabokrtský)

**NLP framework**
*TectoMT*

**MT system**
*TectoMT*

lemmatization

tagging

parsing

# Treex vs. TectoMT

**Treex**

**2005** … **2011**

### NLP framework
*TectoMT*

### MT system
*TectoMT*

- lemmatization
- tagging
- parsing

### multi-purpose NLP framework
*Treex*

### MT system
*TectoMT*

- lemmatization
- tagging
- parsing

- coreference
- CzEng analysis
- named entity r.
- SMT preproc.

- PEDT preprocessing
- treebank conversions
- alignment (word,tree)
- etc.

# Treex vs. TectoMT

**Treex**

2005 … 2011

**NLP framework**
*TectoMT*

**MT system**
*TectoMT*

lemmatiza...

ta...g

...ing

**redesigned and reimplemented**
➡ **easier to use**
➡ **more flexible**

**multi-purpose NLP framework**
*Treex*

**MT system**
*TectoMT*

lemmatization

tagging

parsing

coreference

CzEng analysis

named entity r.

SMT preproc.

PEDT preprocessing

treebank conversions

alignment (word,tree)

etc.

# Treex vs. TectoMT

**Treex**

**2005** …

**English**  **Czech**

**NLP framework**
*TectoMT*

**English**  **Czech**

**Russian**  **Tamil**

**NLP framework**

**Polish**  **Esperanto**

**MT system**
*TectoMT*

lemmatizati...

t...g

...ing

**redesigned and reimplemented**
➡ **easier to use**
➡ **more flexible**
➡ **more langs**

**Spanish**  **French**

*TectoMT*

**German**  **Arabic**

coreference  PEDT preprocessing

**Vietnamese**  **Hindi**  ons

named entity r...  alignment (word tree)

**Urdu**  **Finish**

2005

English

English

Czech

Tamil

Framework

Esperanto

French

German

Arabic

coreference  PEDT preprocessing

Vietnamese  Hindi  ons

named entity  alignment (word tree)

Urdu  Finish

M

Tec

TectoMT

reimpl

➡ easie  use
➡ more flexible
➡ more langs

**Special offer
Call now and get
one extra Treex
for free**

# TectoMT

**Treex**

linguistically motivated MT system (English to Czech pilot)

- deep syntactic (tectogrammatical) transfer
- translation process divided to more than 90 "blocks"
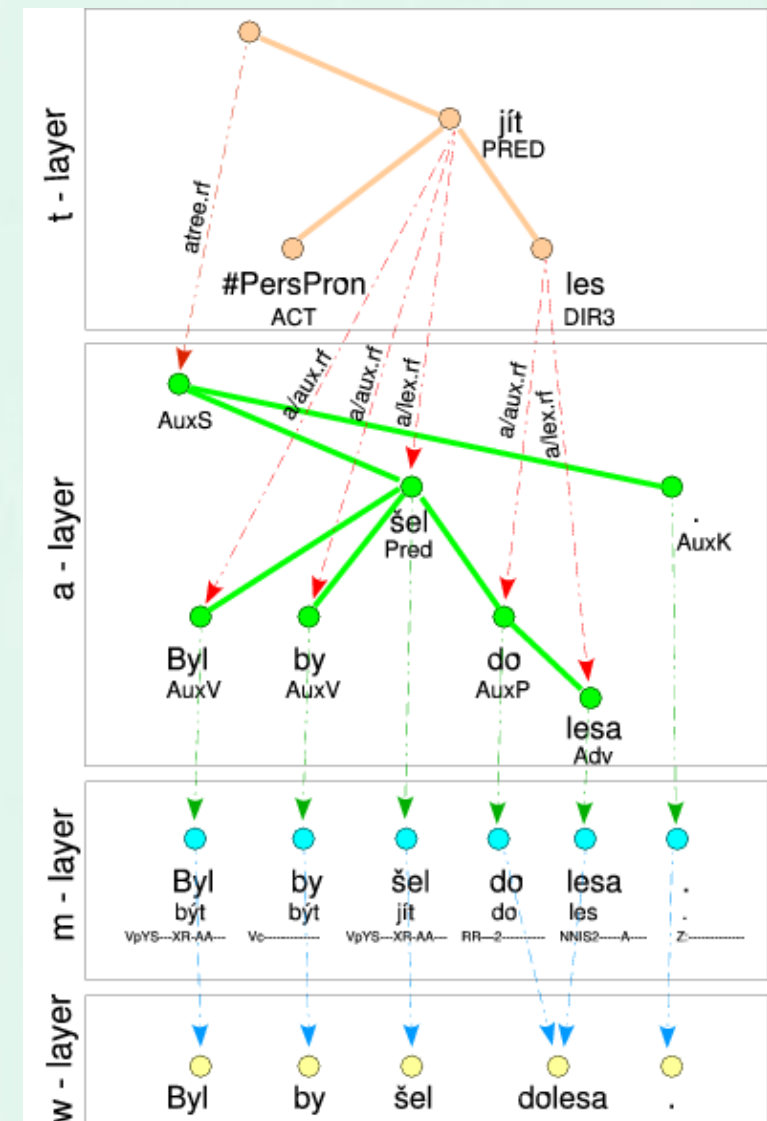- combining  statistical  and  rule based  blocks

**ANALYSIS**     **TRANSFER**     **SYNTHESIS**

**tectogramatical layer** ———————————————— t-layer

fill formems    grammatemes

fill morphological categories

build t-tree

impose agreement

mark edges to contract

add functional words

**analytical layer** ———————————————————————— a-layer

analytical functions

query dictionary    use HMTM

parser (McDonald's MST)

generate wordforms

**morphological layer** —————————————————————— m-layer

tagger (Morce)

concatenate

lemmatization

tokenization

segmentation

source language (English)     target language (Czech)  w-layer

# 4 layers of language description
## implemented in Prague Dependency Treebank (PDT)
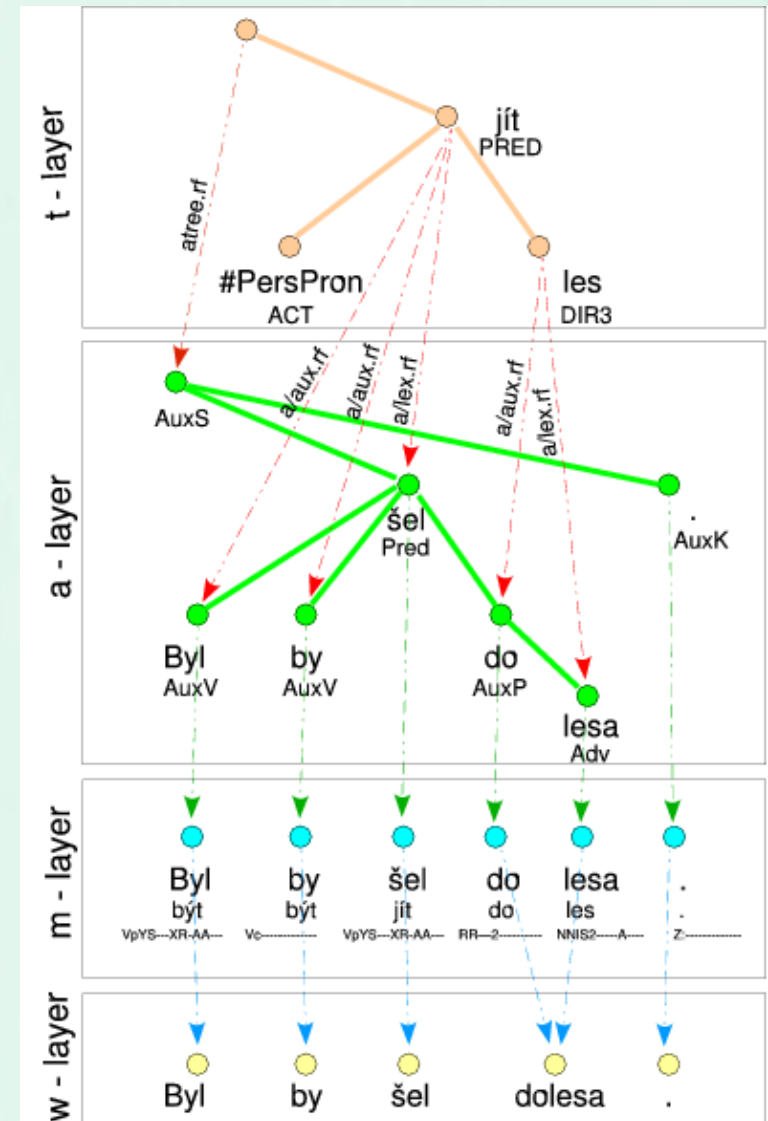
**Treex**

- **tectogrammatical layer**
  deep-syntactic dependency trees

- **analytical layer**
  surface-syntactic dependency trees, labeled edges

- **morphological layer**
  lemma & POS tag for each word

- word layer
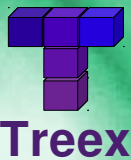  raw (tokenized) text

# 4 layers of language description
## implemented in Prague Dependency Treebank (PDT)

- **tectogrammatical layer**
  deep-syntactic dependency trees

- abstraction from many language-specific phenomena

- autosemantic (meaningful) words
  ~ **nodes**

- functional words (prepositions, auxiliaries)
  ~ **attributes**

- syntactic-semantic relations (dependecies)
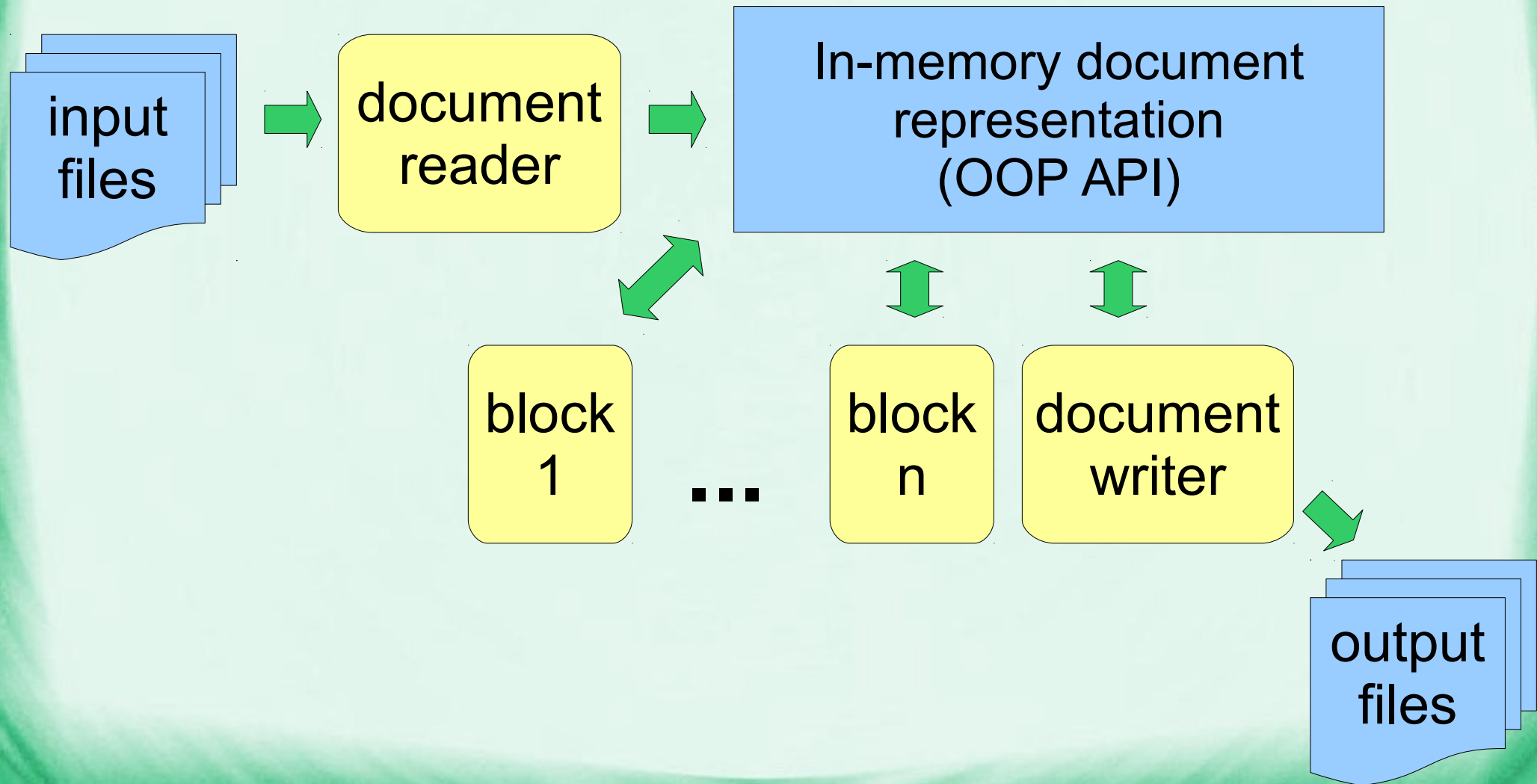  ~ **edges**

- added nodes (e.g. becasue of pro-drop)

- …

- Mostly backward compatible adaptations (adding attributes)
  - **formeme** (n:2, n:k+3, v:že+vfin, v:rc, adj:attr)
  - attributes for clauses, is_passive ($\rightarrow$ diathesis),...
- is_member (for conjuncts on a-layer) is stored with prepositions
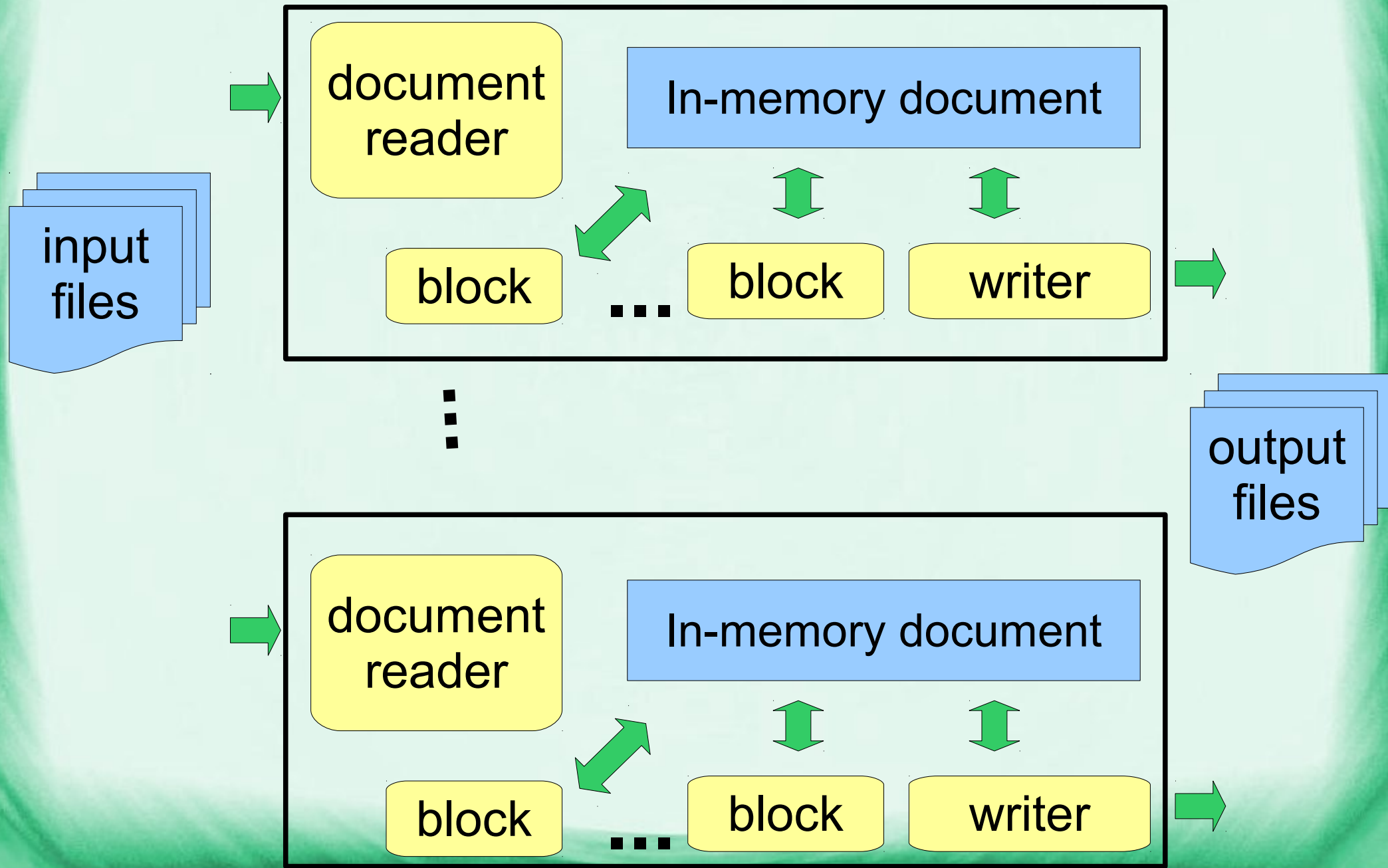
- All layers stored in **one file**
- A-layer and m-layer merged into one
- Two more layers:
  - P-layer phrase-structure trees
  - N-layer named entities

# Treex architecture

Treex

input files → document reader → In-memory document representation (OOP API)

block 1 ... block n | document writer → output files

# Treex architecture
## parallelization (using SGE cluster)

**Treex**

input files

document reader

In-memory document

block · · · block writer

· · ·

document reader

In-memory document

block · · · block writer

output files

# Treex architecture processing units

- **block** – elementary processing unit in Treex
  - corresponding to a given NLP subtask
  - one Perl class (Treex::Block::*), saved in one file
- **scenario** – a sequence of blocks
  - saved in plain text files or a Treex::Scen::* Perl class
  - just a list of the blocks' names and their parameters
- **application** – represents an end-to-end NLP task
  - described by a scenario that
    - starts with a **reader** (input conversion)
    - ends with a **writer** (output conversion)
  - Readers can split the input file into more in-memory docs.
  - There are readers&writers for a number of popular formats: plain text, CoNLL, PDT PML, Penn MRG, Tiger...
    - `*.treex.gz`

# Treex architecture processing units

Blocks can be easily substituted with an alternative solution.

# Treex architecture processing units

**Treex**

Blocks can be easily substituted with an alternative solution.

Scenario A

`W2A::EN::Segment`

`W2A::EN::Tokenize`

`W2A::EN::TagMorce`

`W2A::EN::Lemmatize`

`W2A::EN::ParseMST`

Scenario B

`W2A::SegmentOnNewlines`

`W2A::EN::TagLinguaEn`

`W2A::EN::Lemmatize`

`W2A::EN::ParseMalt`

- **Document**
  - stored in one file
  - sequence of sentences
- **Bundle** ("bundle of trees")
  - corresponds to one sentence
- **Zone**
  - one for each language (Arabic, Czech, English,...)
  - and optionally a variant ("selectors" src, trg, ref,...)
- **Tree**
  - layer of language description: A, T (plus P, N)
  - m-layer is stored with the a-layer in one tree

# Treex architecture data units

# Treex architecture data units

**Treex**

**DOCUMENT**

sentence 1

**BUNDLE**

**Zone en_src**

**W-layer**
*Peter does not love Mary.*

**M-layer**

Peter do n...ove Mary
_NNP _VBZ _...BD _NNP

**A-layer**

Peter do n...ove Mary
Sb AuxV _N...Pred Obj

**T-layer**

Peter love Mary
ACT PRED PAT

**Zone cs_src**

**W-layer**
*Petr nemiluje Marii.*

**...-layer**

Petr m... Marie
NNMS1 _VB-S... -NA NNFS4

**A-layer**

Petr milo... Marie
Sb Pr... Obj

**T-layer**

Petr milovat Marie
ACT PRED PAT

sentence 2

**BUNDLE**

...

...

sentence N

**BUNDLE**

# Treex architecture data units

**Treex**

**DOCUMENT**

sentence 1

## BUNDLE

**Zone en_src**

**W-layer**
*Peter does not love Mary.*

**M-layer**

Peter do n...ove Mary
NNP VBZ ...BD NNP

**A-layer**

Peter do n...ove Mary
Sb AuxV N...Pred Obj

**T-layer**

Peter    love         Mary
ACT      PRED         PAT

**Zone cs_trg**

**W-layer**
*Petr nemiluje Marii.*

**M-layer**

Petr    m...at    Marie
NNMS1 VB- ...P-NA NNFS4

**A-layer**

Petr        mil...    Marie
Sb          Pr...     Obj

**T-layer**

Petr      milovat     Marie
ACT       PRED        PAT

sentence 2

## BUNDLE

...

sentence N

## BUNDLE
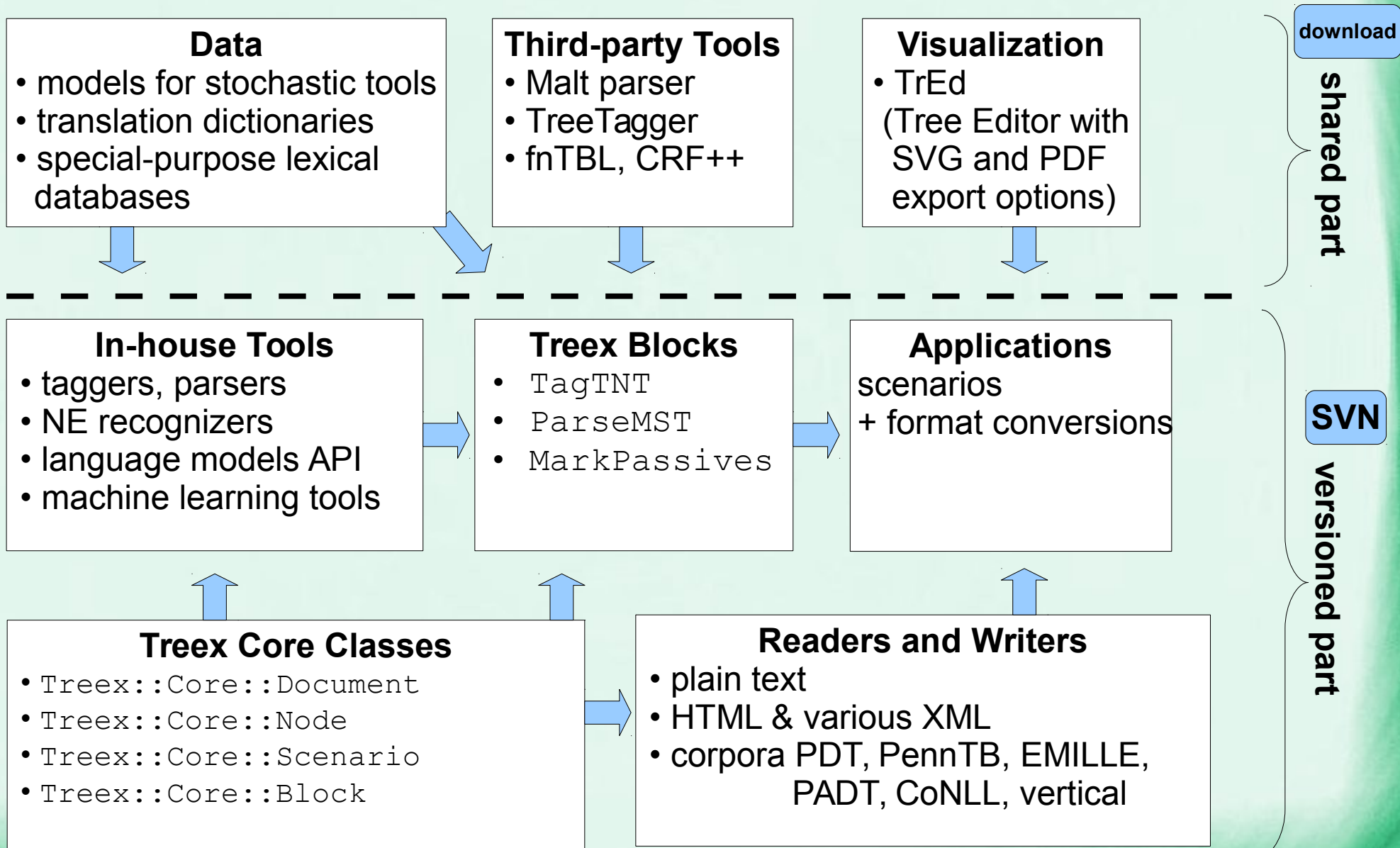
# Internals – Design decisions

- Perl (wrappers for binaries, Java,...)
- Linux (some applications platform-independent)
- OOP (Moose)
- Open source (GNU GPL for the versioned part)
- Neutral w.r.t. methodology (statistical, rule-based)
- Multilingual
- Open standards (Unicode, XML)

# Internals – Components

**Treex**

**shared part**

**Data**
- models for stochastic tools
- translation dictionaries
- special-purpose lexical databases

**Third-party Tools**
- Malt parser
- TreeTagger
- fnTBL, CRF++

**Visualization**
- TrEd
(Tree Editor with SVG and PDF export options)

**SVN**

**versioned part**

**In-house Tools**
- taggers, parsers
- NE recognizers
- language models API
- machine learning tools

**Treex Blocks**
- `TagTNT`
- `ParseMST`
- `MarkPassives`

**Applications**
scenarios
+ format conversions

**Treex Core Classes**
- `Treex::Core::Document`
- `Treex::Core::Node`
- `Treex::Core::Scenario`
- `Treex::Core::Block`

**Readers and Writers**
- plain text
- HTML & various XML
- corpora PDT, PennTB, EMILLE, PADT, CoNLL, vertical

# Internals – Statistics

**Treex**

- Developed since 2005, over ten developers

- Over 400 blocks (140 English, 120 Czech,
  60 English-to-Czech, 30 other languages,
  50 language independent)

- Taggers (5 English, 3 Czech, 1 German and Russian, Tamil)
  Parsers (Dep. 2 English, 3 Czech, 2 German; Const. 2 English)
  Named Entity Recognizers (2 Czech,1 English)

- Speed example: Best version of English-to-Czech MT
  1.2 seconds per sentence plus 90 seconds loading,
  with 20 computers in cluster: 2000 sentences in 4 min

# Conclusion
## Treex main properties

- emphasized efficient development, modular design and reusability

- stratificational approach to the language

- unified object-oriented interface for accessing data structures

- comfortable development
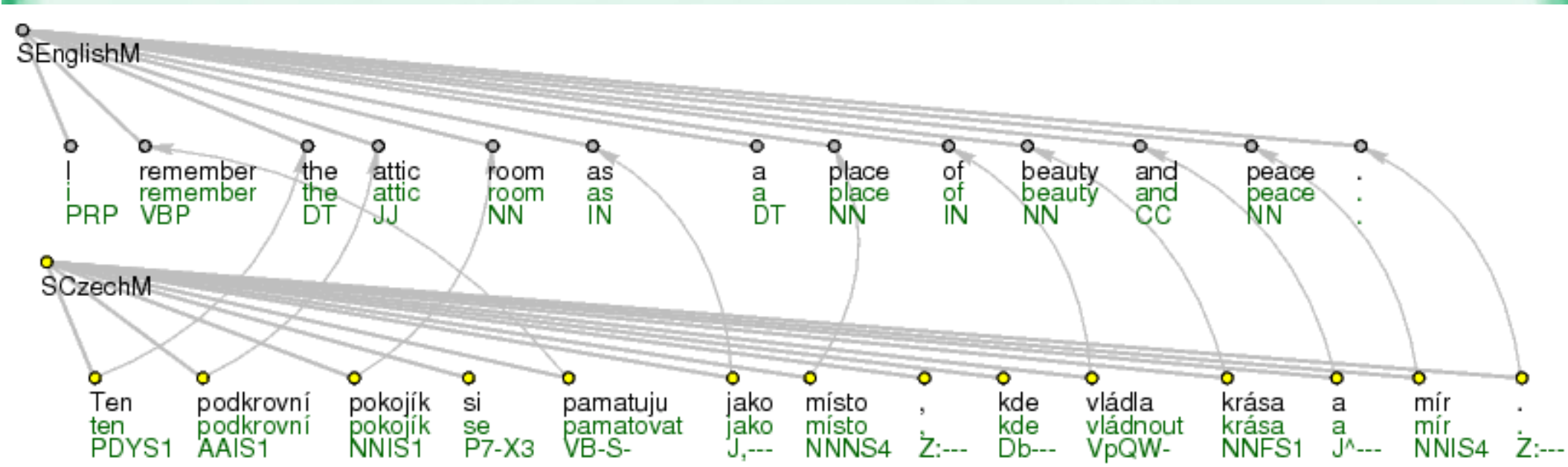
**Treex**
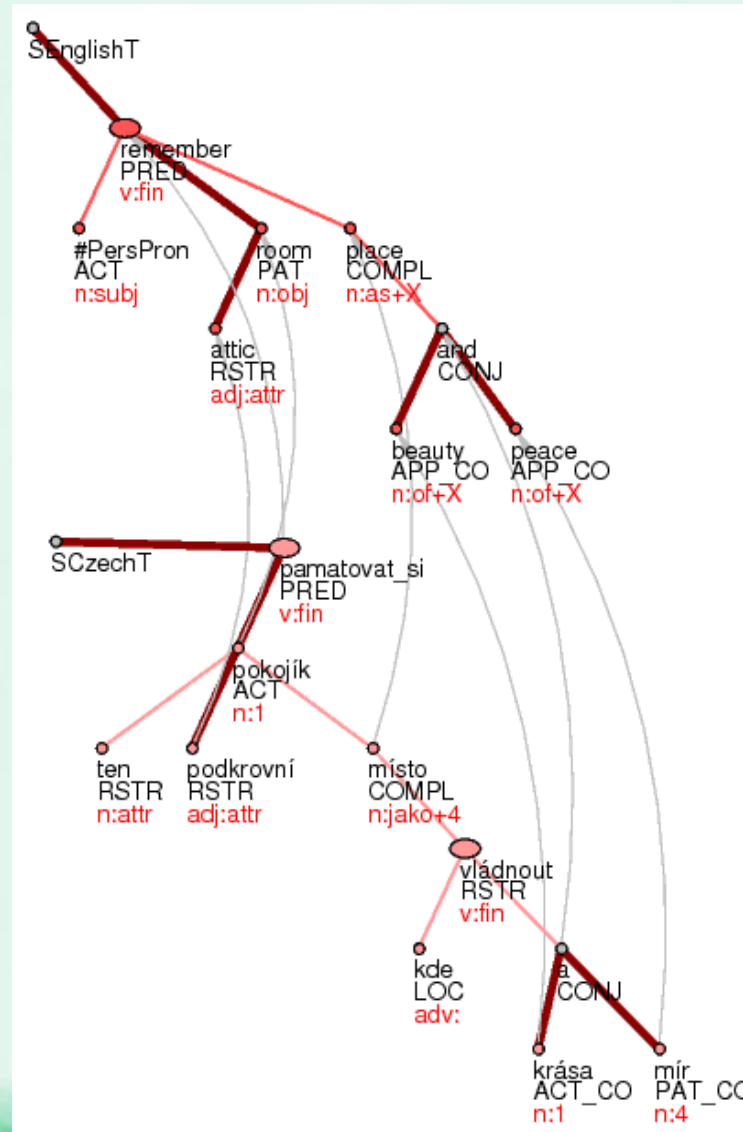
# TrEd visualization

**Treex**

## translation

# TrEd visualization

## word alignment on the morphological layer

# TrEd visualization

word alignment on the tectogrammatical layer
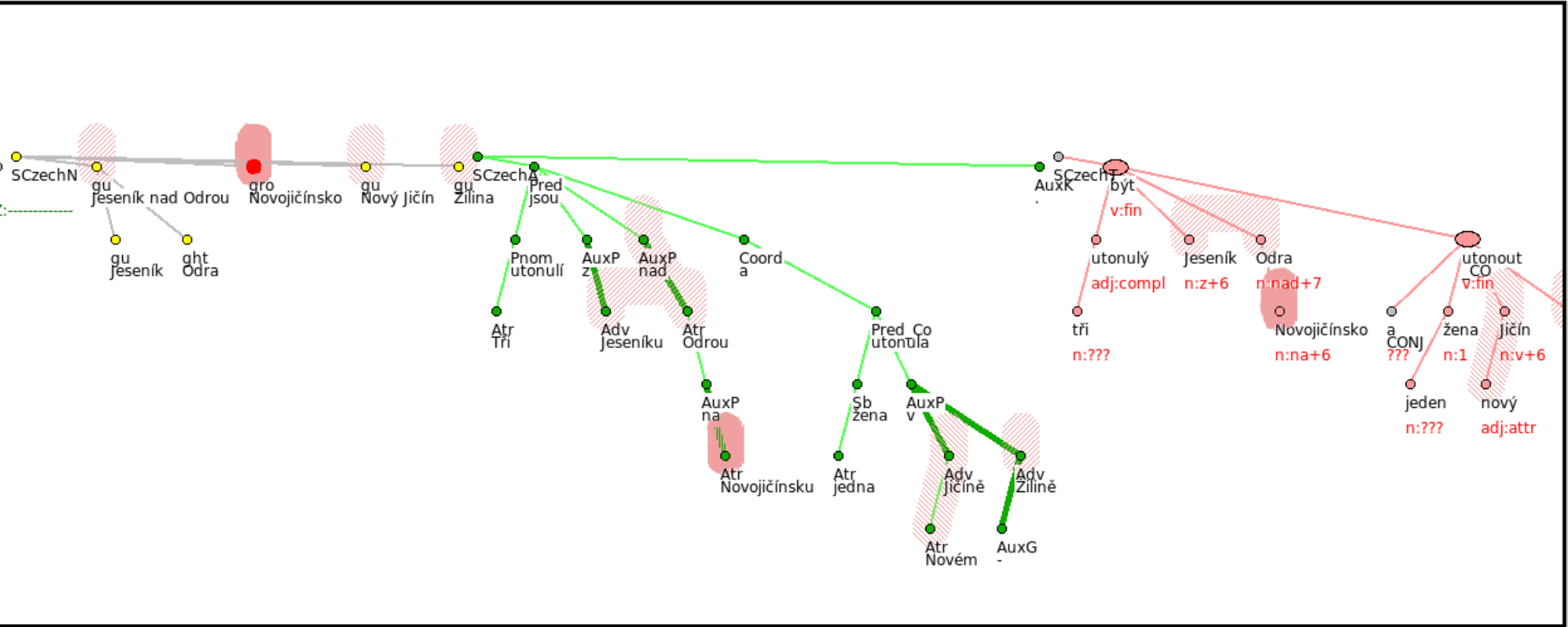
# TrEd visualization

## named entities

# Block example – SVO to SOV code

Treex

```perl
package Tutorial::Svo2SovSolution;
use Moose;
use Treex::Core::Common;
extends 'Treex::Core::Block';


sub process_anode {
  my ( $self, $a_node ) = @_;
  if ( $a_node->tag =~ /^V/ ) {          # verb found
    foreach my $child ( $a_node->get_echildren() ) {
      if ( $child->afun eq 'Obj' ) {     # object found
        # Move the object and its subtree so it precedes the verb
        $child->shift_before_node($a_node);
      }
    }
  }
  return;
}
1;
```

Treex core
Treex convention
**Perl keyword/convention**

# Thank you

**Treex**

Cooperation is welcomed.



http://ufal.mff.cuni.cz/treex

# Thank you

**Treex**

Treex is growing!

http://ufal.mff.cuni.cz/treex