# Introduction to Treex
## Modular NLP Framework

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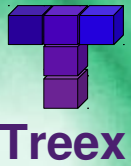
ÚFAL

# Outline

**Treex**

- Motivation, Treex origin (TectoMT)
- Layers of language description
- Treex architecture
- Treex internals
- Overview of tools and applications
- Conclusion and examples

# Motivation

## 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 origin (TectoMT)

2005 (Zdeněk Žabokrtský)

**NLP framework**
*TectoMT*

**MT system**
*TectoMT*

lemmatization

tagging

parsing

# Treex origin (TectoMT)

**Treex**

## 2005 … 2011

### NLP framework
### *TectoMT*

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

| MT system *TectoMT* | lemmatization |
| | tagging |
| | parsing |

| MT system *TectoMT* | lemmatization |
| | tagging |
| | parsing |

| coreference | PEDT preprocessing |
| CzEng analysis | treebank conversions |
| named entity r. | alignment (word,tree) |
| SMT preproc. | etc. |

# Treex origin (TectoMT)

**Treex**

## 2005 ... 2011

### NLP framework
*TectoMT*

**MT system**
*TectoMT*

lemmatizati...

t...ng

...ing

**Now not only tectogrammatics and not only MT**

➡ **renamed**

### multi-purpose NLP framework
*Treex*

**MT system**
*TectoMT*

lemmatization

tagging

parsing

coreference

CzEng analysis

named entity r.

SMT preproc.

PEDT preprocessing

treebank conversions
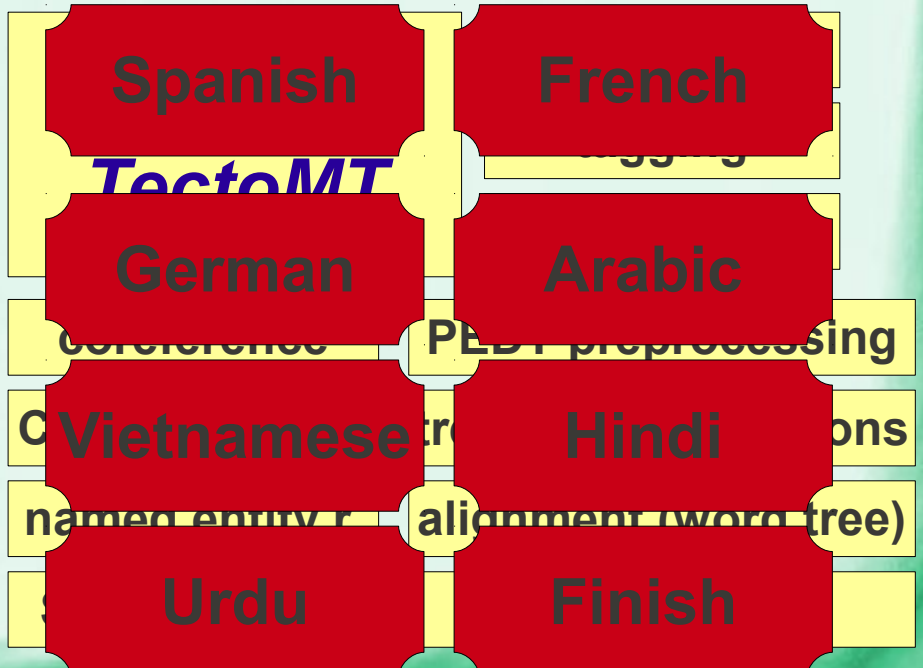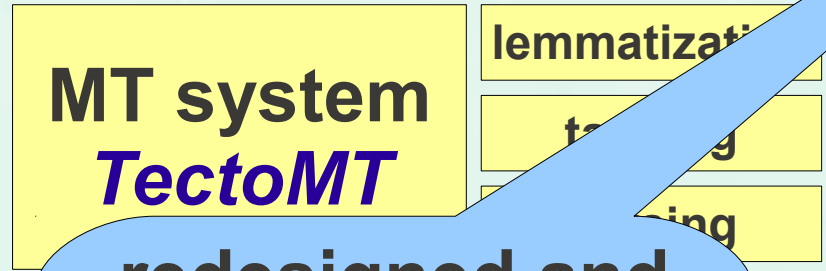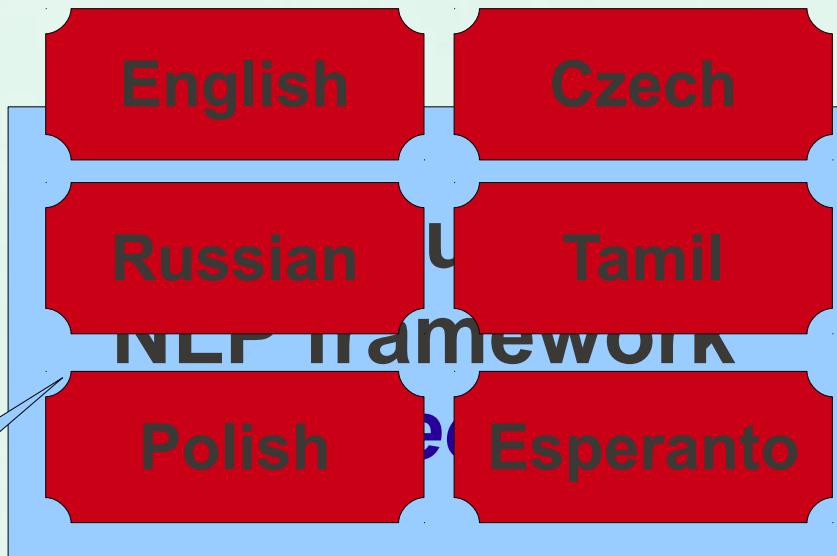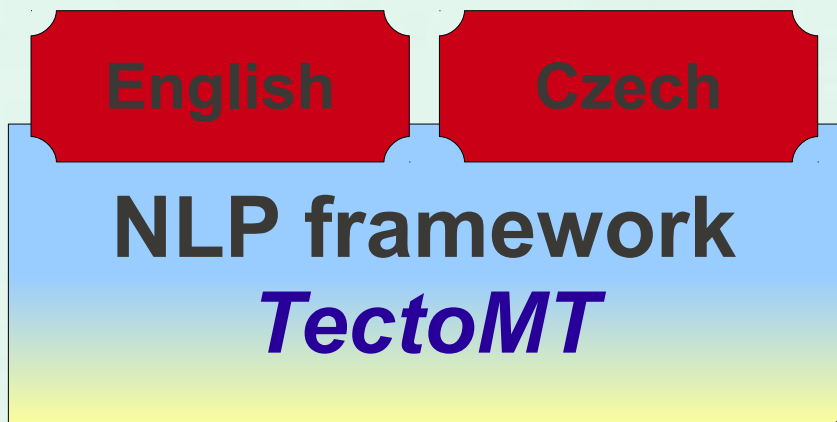
alignment (word,tree)

etc.

# Treex origin (TectoMT)

**Treex**

**2005**  ...  **2011**

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

**NLP framework** *TectoMT*

**MT system** *TectoMT*

| | |
|---|---|
| lemmatizati | |
| t...ng | |
| ...ing | |

**redesigned and reimplemented**

➡ **easier to use**
➡ **more flexible**

**MT system** *TectoMT*

| | |
|---|---|
| lemmatization | |
| tagging | |
| parsing | |

| | |
|---|---|
| coreference | PEDT preprocessing |
| CzEng analysis | treebank conversions |
| named entity r. | alignment (word,tree) |
| SMT preproc. | etc. |

# Treex origin (TectoMT)

**Treex**

2005 …

**English** **Czech**

**NLP framework**
*TectoMT*

**MT system**
*TectoMT*

lemmatiza...

t...g

...ing

**redesigned and reimplemented**

➡ **easier to use**

➡ **more flexible**

➡ **more langs**

**English** **Czech**

**Russian** **Tamil**

**NLP framework**

**Polish** **Esperanto**

**Spanish** **French**

*TectoMT*

**German** **Arabic**

coreference PEDT preprocessing

C **Vietnamese** Hindi ons

named entity alignment (word tree)

**Urdu** **Finish**

# Treex origin (TectoMT)



**2005**

English

Czech

Tamil

Framework

Esperanto

French

German

Arabic

coreference  PEDT preprocessing

Vietnamese  Hindi

named entity  alignment (word tree)

Urdu  Finish

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

TectoMT

reimplemented

➡ **easier to use**
➡ **more flexible**
➡ **more langs**

# TectoMT

**Treex**

## linguistically motivated MT system (English to Czech pilot)

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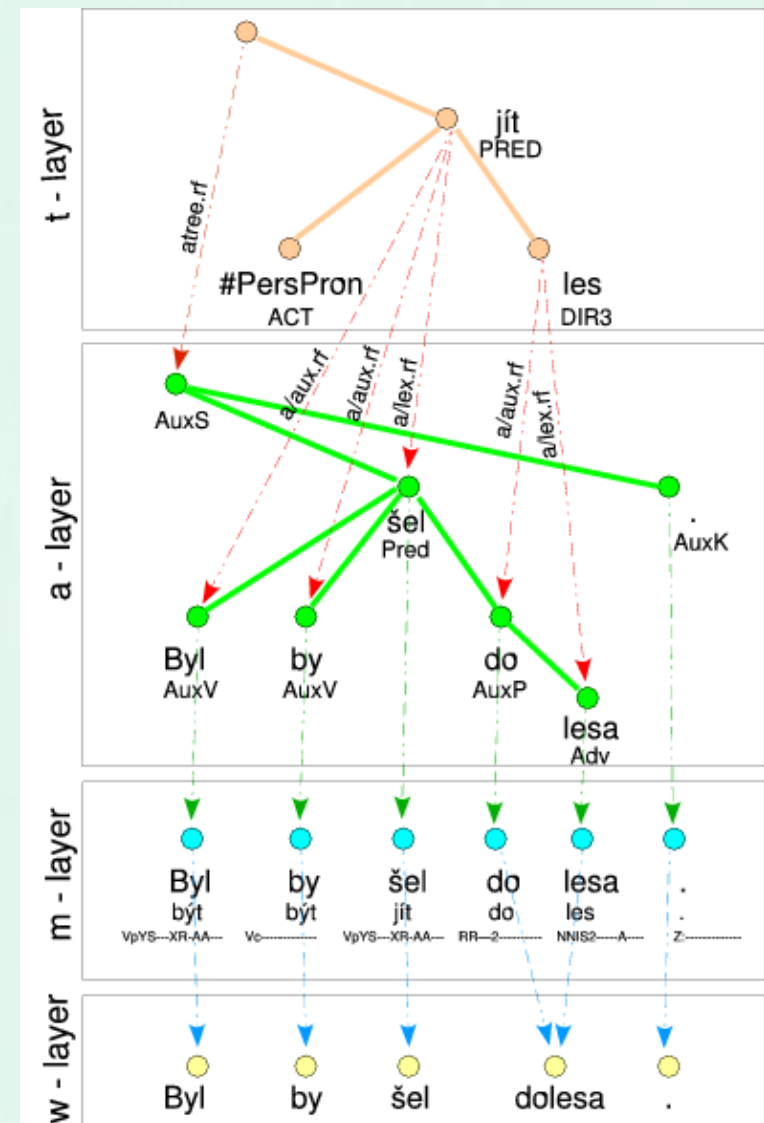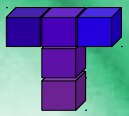- combining  statistical  and  rule based  blocks

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

**tectogramatical layer** ──────────────────────  t-layer

fill formems    grammatemes

query dictionary    use HMTM

fill morphological categories

build t-tree

impose agreement

mark edges to contract

add functional words

**analytical layer** ──────────────────────  a-layer

analytical functions

generate wordforms

parser (McDonald's MST)

**morphological layer** ──────────────────────  m-layer

tagger (Morce)

concatenate

lemmatization

tokenization

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

segmentation

# 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**
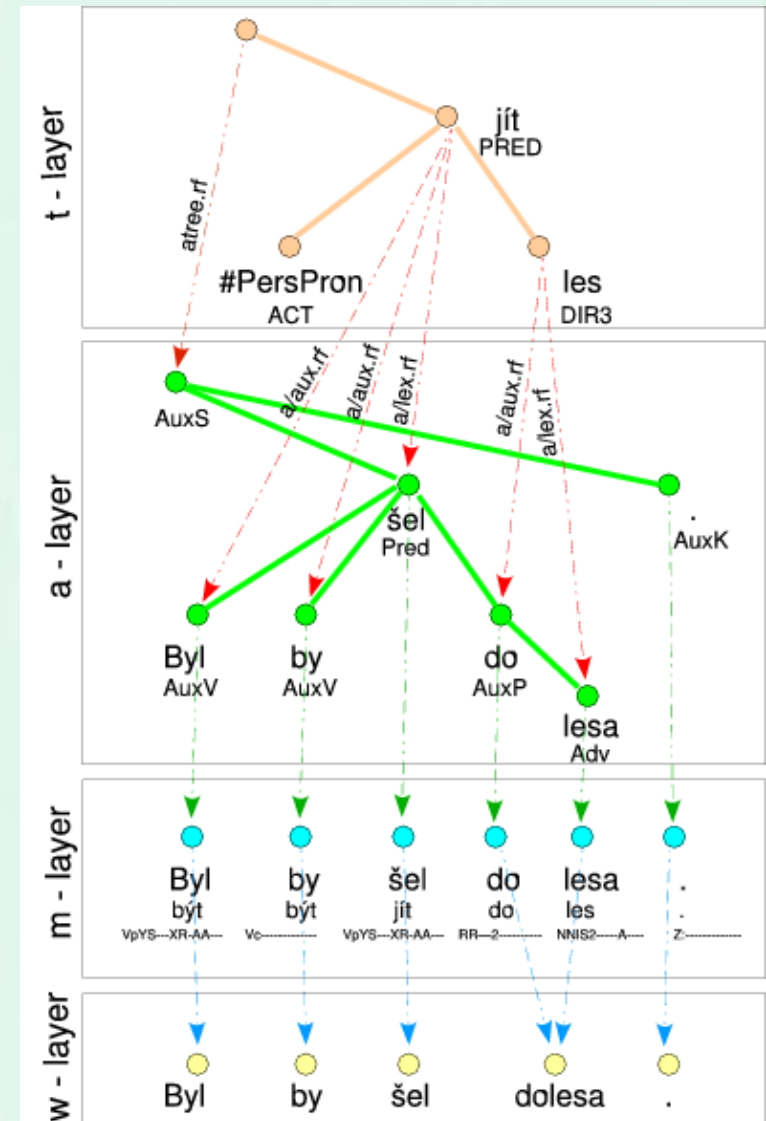  lemma & POS tag for each word

- 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)
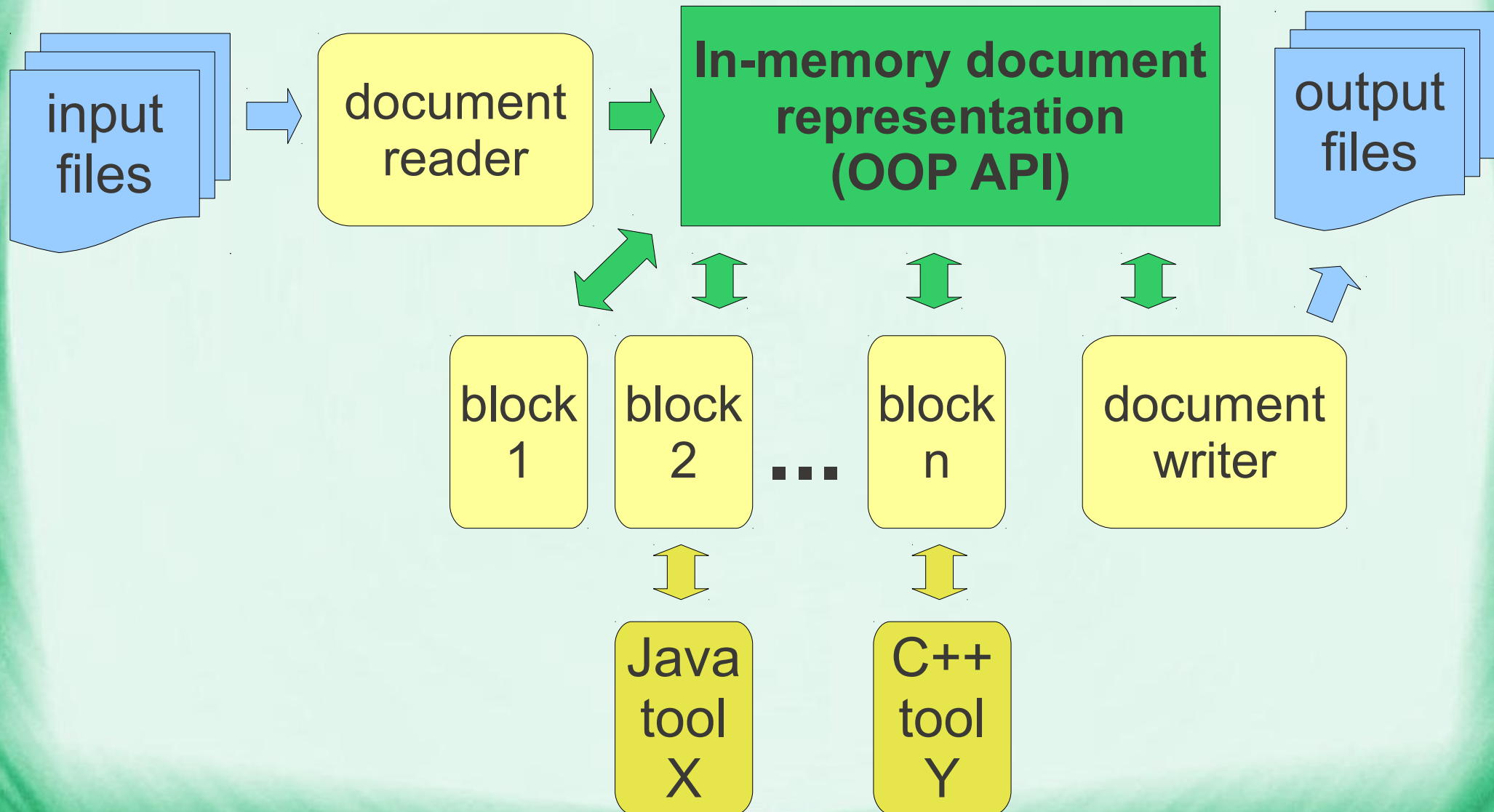
- ...

# layers of language description
## implemented in Treex

- 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),...
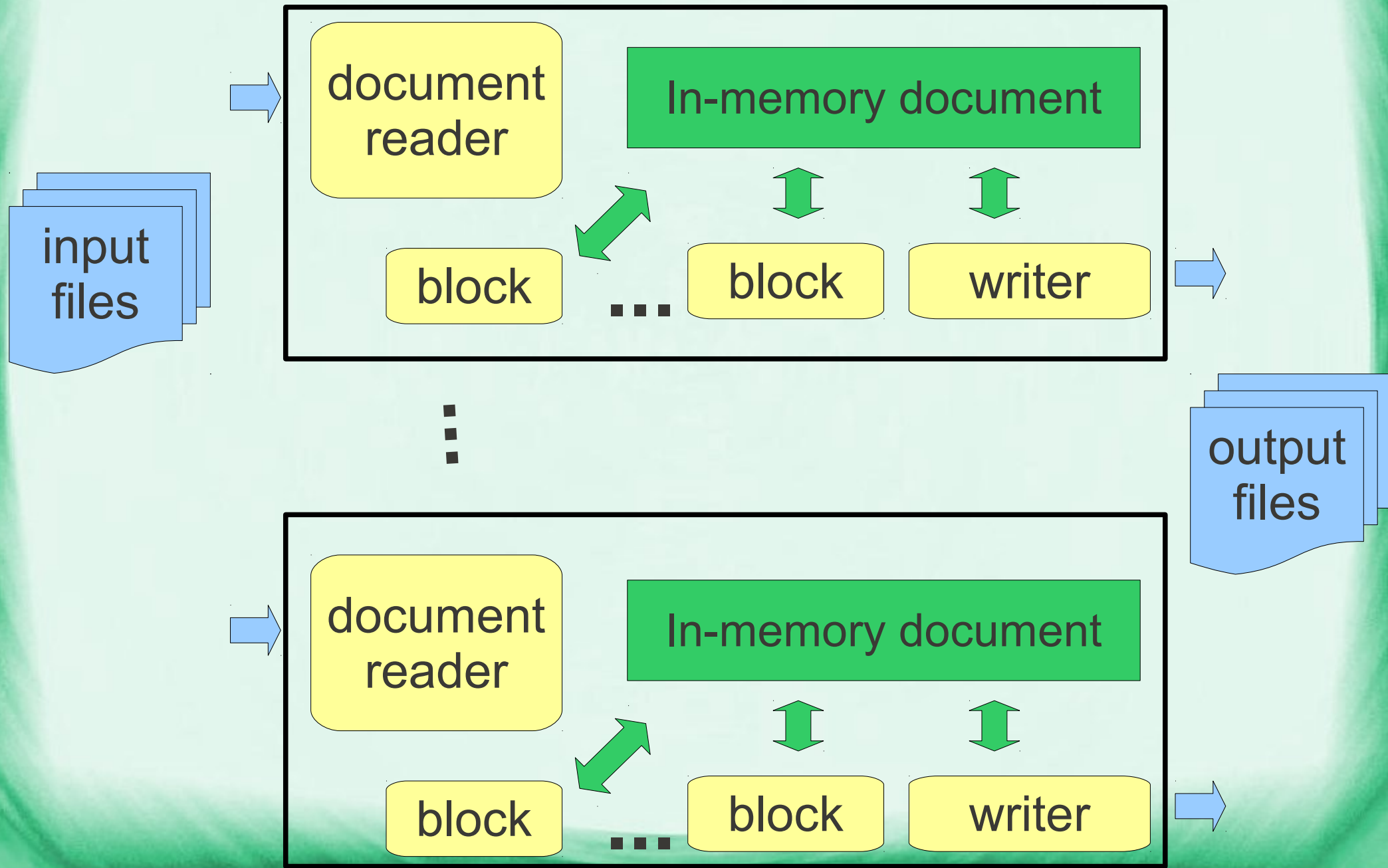- is_member (for conjuncts on a-layer) is stored with prepositions

- 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**

# Treex architecture

**Treex**

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

block 1    block 2    ...    block n    document writer

Java tool X    C++ tool Y

# Treex architecture
## parallelization (using SGE cluster)
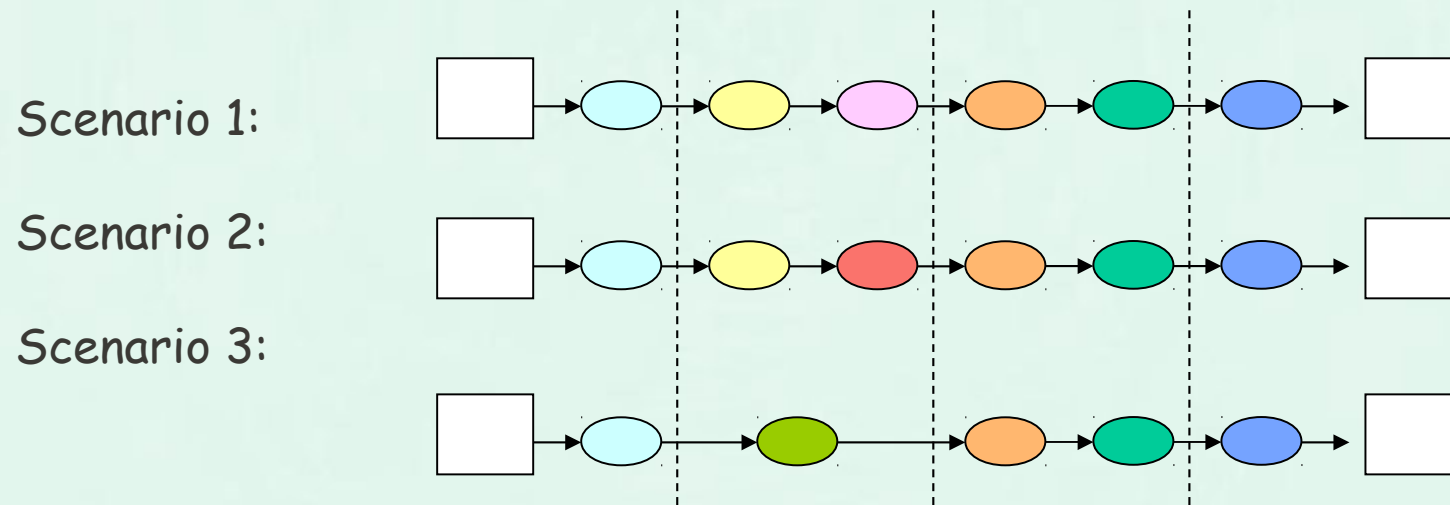
**Treex**

# Treex architecture processing units

- **block** – elementary processing unit in Treex
  - corresponding to a given NLP subtask
  - one Perl class, saved in one file
- **scenario** – a sequence of blocks
  - can be saved in plain text *.scen files
  - 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.

Scenario 1:

Scenario 2:

Scenario 3:

# Treex architecture processing units

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
```
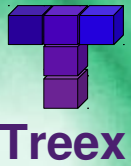
Scenario B

```
W2A::SegmentOnNewlines

W2A::EN::TagLinguaEn

W2A::EN::Lemmatize
W2A::EN::ParseMalt
```
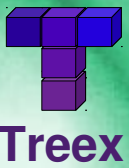
# Treex architecture
## data units

- **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, trans, 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

## DOCUMENT

sentence 1        sentence 2    ...    sentence N
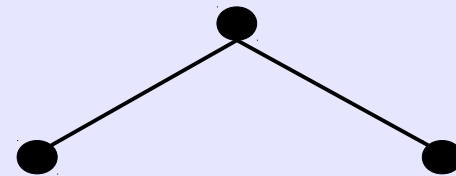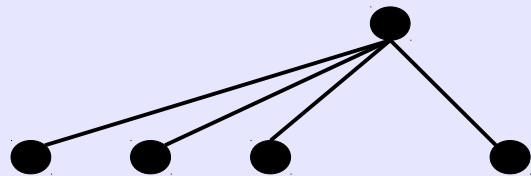
### BUNDLE

| Zone en_src | Zone cs_src |
|---|---|
| **W-layer** *Peter does not love Mary.* | **W-layer** *Petr nemiluje Marii.* |

**A-layer**

| Peter | does | not | love | Mary |
|---|---|---|---|---|
| Peter | do | not | love | Mary |
| Sb | AuxV | Neg | Pred | Obj |
| NNP | VBZ | RB | VBD | NNP |

**A-layer**

| Petr | nemiluje | Marii |
|---|---|---|
| Petr | milovat | Marie |
| Sb | Pred | Obj |
| NNMS1 | VB-S—3P-NA | NNFS4 |

**T-layer**

| Peter | love | Mary |
|---|---|---|
| ACT | PRED | PAT |
| n:subj | v:fin | n:obj |
| sg,... | neg:1,sim,... | sg,... |

**T-layer**

| Petr | milovat | Marie |
|---|---|---|
| ACT | PRED | PAT |
| n:1 | v:fin | n:4 |
| sg,... | neg:1,sim,... | sg,... |

### BUNDLE

### BUNDLE

...

# Treex architecture
# data units

**Treex**

## DOCUMENT

sentence 1                     sentence 2        ...        sentence N

### BUNDLE

| **Zone en_src** | **Zone cs_src** |
|---|---|
| **W-layer** *Peter does not love Mary.* | **W-layer** *Petr nemiluje Marii.* |

**A-layer**



| Peter | does | not | love | Mary |
|---|---|---|---|---|
| Peter | do | not | love | Mary |
| Sb | AuxV | Neg | Pred | Obj |
| NNP | VBZ | RB | VBD | NNP |

**A-layer**

| Petr | nemiluje | Marii |
|---|---|---|
| Petr | milovat | Marie |
| Sb | Pred | Obj |
| NNMS1 | VB-S—3P-NA | NNFS4 |

**T-layer**

| Peter | love | Mary |
|---|---|---|
| ACT | PRED | PAT |
| n:subj | v:fin | n:obj |
| sg,... | neg:1,sim,... | sg,... |

**T-layer**

| Petr | milovat | Marie |
|---|---|---|
| ACT | PRED | PAT |
| n:1 | v:fin | n:4 |
| sg,... | neg:1,sim,... | sg,... |

### BUNDLE

...

form
lemma
afun (edge label)
morphological tag (PoS)

t_lemma
functor (semantic role)
formeme (morphosyntactic abstraction)
grammatemes for number, negation, tense

### BUNDLE

# Treex architecture
# data units

**Treex**

## DOCUMENT

sentence 1                      sentence 2    ...    sentence N

### BUNDLE                      ### BUNDLE          ### BUNDLE

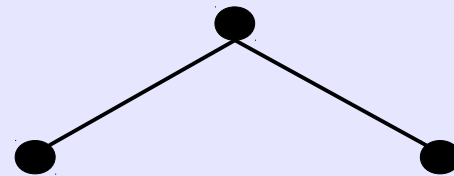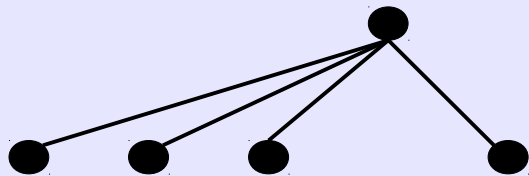**Zone en_src**                 **Zone cs_src**

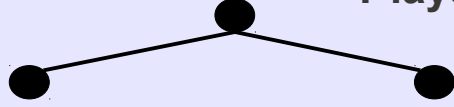**W-layer**                     **W-layer**
*Peter does not love Mary.*     *Petr nemiluje Marii.*

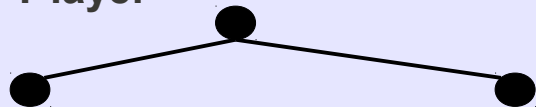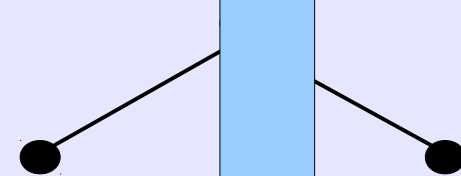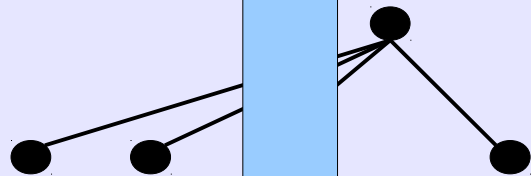**A-layer**                     **A-layer**

| Peter | does | | love | Mary | | Petr | nem... | | Marii |
|-------|------|---|------|------|---|------|--------|---|-------|
| Peter | do | | love | Mary | | Petr | milo... | | Marie |
| Sb | AuxV | | Pred | Obj | | Sb | Pred... | | Obj |
| NNP | VBZ | | VBD | NNP | | NNMS1 | VB-S... | NA | NNFS4 |

form
lemma
afun (edge label)
morphological tag (PoS)

**T-layer**                     **T-layer**

| Peter | love | Mary | | Petr | milovat | Marie |
|-------|------|------|---|------|---------|-------|
| ACT | PRED | PAT | | ACT | PRED | PAT |
| n:subj | v:fin | n:obj | | n:1 | v:fin | n:4 |
| sg,... | neg:1,sim,... | sg,... | | sg,... | neg:1,sim,... | sg,... |

t_lemma
functor (semantic role)
formeme (morphosyntactic abstraction)
grammatemes for number, negation, tense

# Treex architecture data units

**Treex**

## DOCUMENT

sentence 1          sentence 2          ...          sentence N

### BUNDLE

| **Zone en_src** | **Zone cs_trans** |
|---|---|

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

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

**A-layer**

| Peter | does | n | love | Mary | Petr | nem | e | Marii |
|---|---|---|---|---|---|---|---|---|
| Peter | do | r | love | Mary | Petr | mil | | Marie |
| Sb | AuxV | N | Pred | Obj | Sb | Pre | | Obj |
| NNP | VBZ | R | VBD | NNP | NNMS1 | VB-S | P-NA | NNFS4 |

**A-layer**

**T-layer**

| Peter | love | Mary | Petr | milovat | Marie |
|---|---|---|---|---|---|
| ACT | PRED | PAT | ACT | PRED | PAT |
| n:subj | v:fin | n:obj | n:1 | v:fin | n:4 |
| sg,... | neg:1,sim,... | sg,... | sg,... | neg:1,sim,... | sg,... |

**T-layer**

### BUNDLE

...

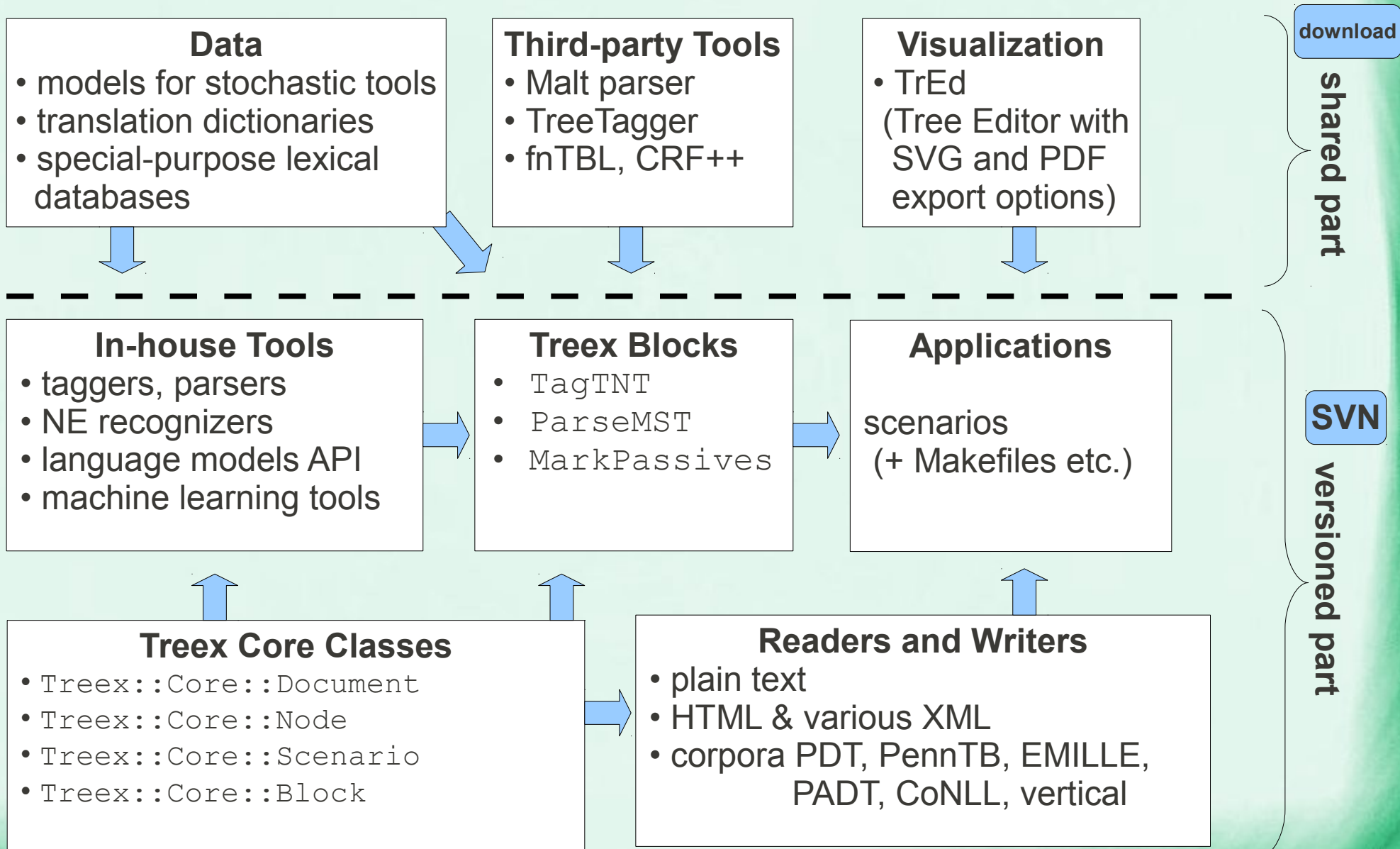form
lemma
afun (edge label)
morphological tag (PoS)

t_lemma
functor (semantic role)
formeme (morphosyntactic abstraction)
grammatemes for number, negation, tense

### BUNDLE

# Internals – Design decisions

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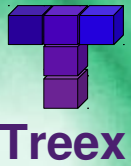- Open standards (Unicode, XML)

# Internals – Components

**Treex**

**shared part**

| | | |
|---|---|---|
| **Data**<br>• models for stochastic tools<br>• translation dictionaries<br>• special-purpose lexical databases | **Third-party Tools**<br>• Malt parser<br>• TreeTagger<br>• fnTBL, CRF++ | **Visualization**<br>• TrEd<br>(Tree Editor with SVG and PDF export options) |

**SVN**

**versioned part**

| | | |
|---|---|---|
| **In-house Tools**<br>• taggers, parsers<br>• NE recognizers<br>• language models API<br>• machine learning tools | **Treex Blocks**<br>• `TagTNT`<br>• `ParseMST`<br>• `MarkPassives` | **Applications**<br><br>scenarios<br>(+ Makefiles etc.) |

| | |
|---|---|
| **Treex Core Classes**<br>• `Treex::Core::Document`<br>• `Treex::Core::Node`<br>• `Treex::Core::Scenario`<br>• `Treex::Core::Block` | **Readers and Writers**<br>• plain text<br>• HTML & various XML<br>• 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
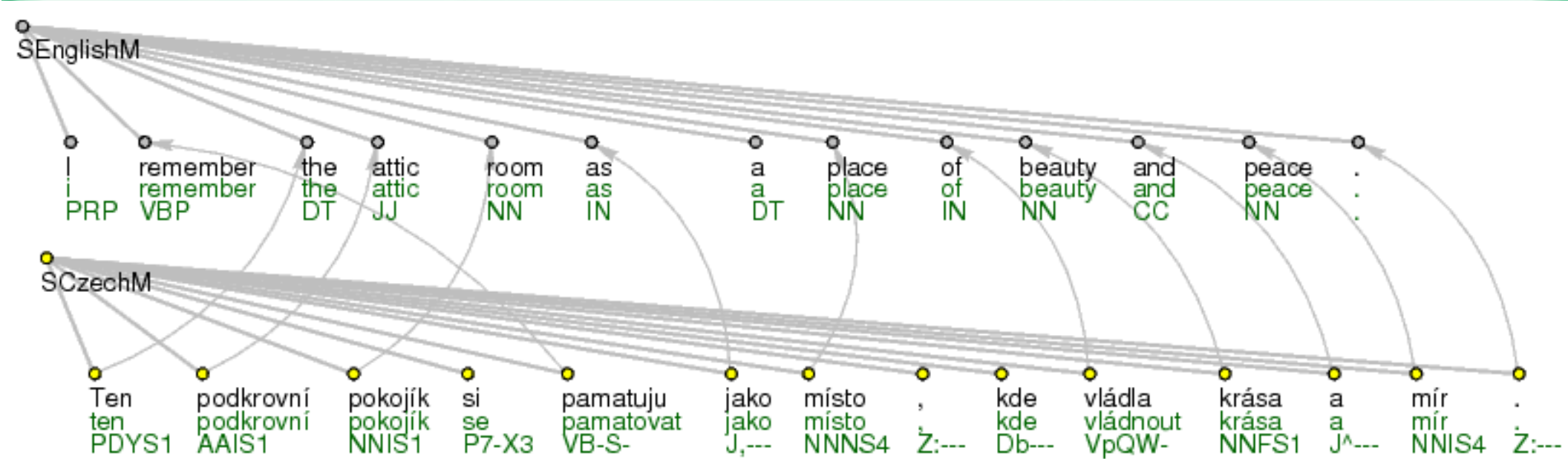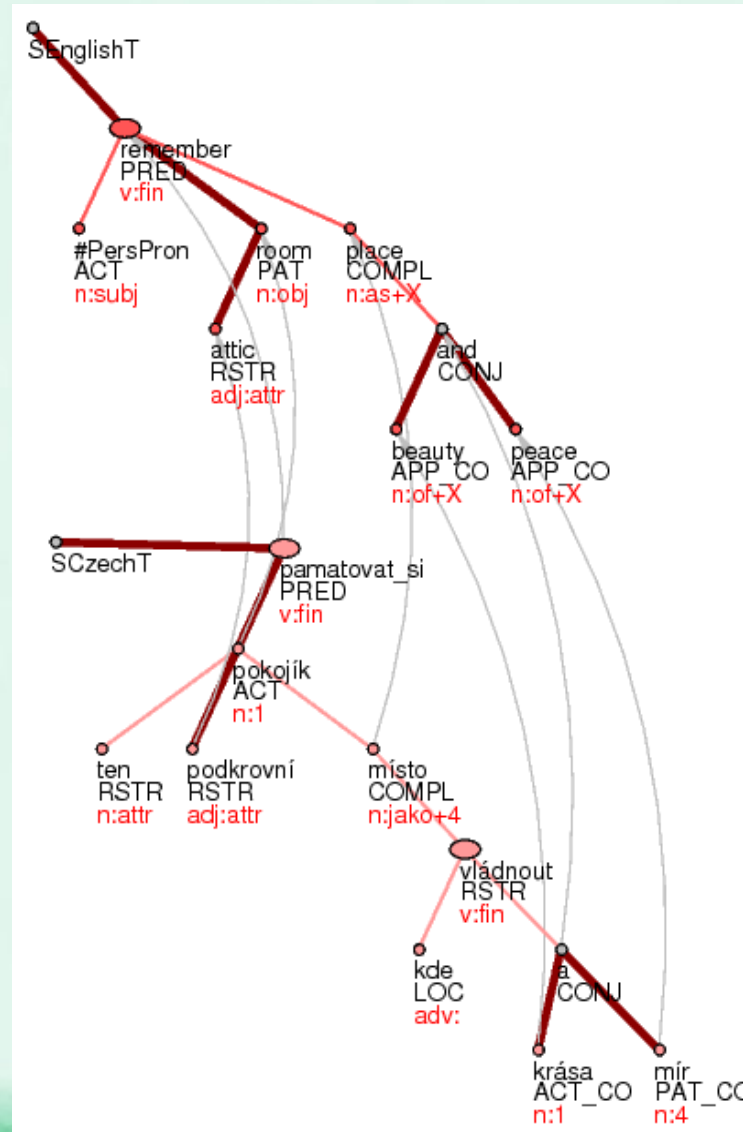
Treex

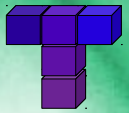# TrEd visualization

## translation

# TrEd visualization

## word alignment on the morphological layer

# TrEd visualization

**Treex**

word alignment on the tectogrammatical layer

# TrEd visualization

## named entities

# Block example – SVO to SOV code

Treex

```perl
package Tutorial::Solution::Svo2Sov;
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;
```
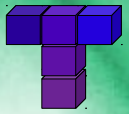
Treex core
Treex convention
Perl keyword/convention

# Thank you

**Treex**

Cooperation is welcomed.



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

# Thank you

**Treex**

Treex is growing!



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