

# Introduction to TectoMT

Zdeněk Žabokrtský, Martin Popel

Institute of Formal and Applied Linguistics  
Charles University in Prague

# Outline

---

- PART 1
  - What is TectoMT?
  - TectoMT's architecture
  - Overview of TectoMT's tools and applications
  
- PART 2 - demo

# What is TectoMT?

- multi-purpose NLP software framework
- created at UFAL since 2005
  
- main linguistic features
  - layered language representation
  - linguistic data structures adopted from the Prague Dependency Treebank
  
- main technical features
  - highly modular, open-source
  - numerous NLP tools already integrated (both existing and new)
  - all tools communicating via a uniform OO infrastructure
  - Linux + Perl
  - reuse of PDT technology (tree editor TrEd, XML...)

# Why “TectoMT” ?

- Tecto..
  - refers the (Praguian) tectogrammar
  - deep-syntactic dependency-oriented sentence representation
  - developed by Petr Sgall and his colleagues since 1960s
  - large scale application in the Prague Dependency Treebank
- .....MT
  - the main application of TectoMT is Machine Translation
- however, not only “tecto” and not only “MT” !!!
- re-branding planned for 2011: TectoMT → Treex

# What is not TectoMT?

---

- TectoMT (as a whole) is not an end-user application
  - it is rather an experimental lab for NLP researchers
- however, releasing of single-purpose stand-alone applications is possible

# Motivation for creating TectoMT

## ■ First, technical reasons:

- Want to make use of more than two NLP tools in your experiment? Be ready for endless data conversions, need for other people's source code tweaking, incompatibility of source code and model versions...
- ⇒ unified software infrastructure might help us in many aspects.

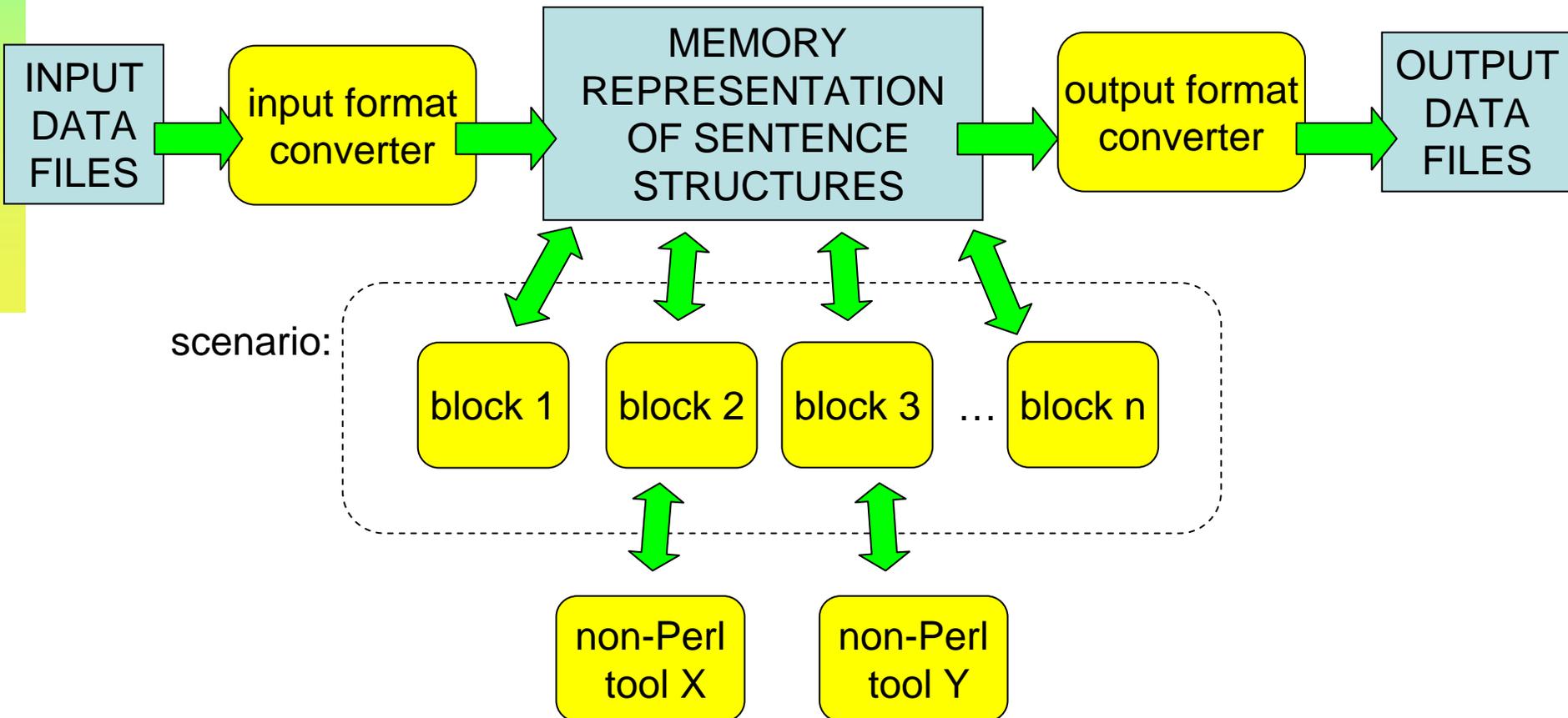
## ■ Second, our long-term MT plan:

- We believe that tectogrammar (deep syntax) as implemented in Prague Dependency Treebank might help to (1) **reduce data sparseness**, and (2) find and **employ structural similarities** revealed by tectogrammar even between typologically different languages.

# Main Design Decisions

- Linux
- Perl as the core language
- set of well-defined, linguistically relevant layers of language representation
- neutral w.r.t. chosen methodology ("rules vs. statistics")
- emphasis on modularity
  - each task implemented by a sequence of blocks
  - each block corresponds to a well-defined NLP subtask
  - reusability and substitutability of blocks
- support for distributed processing

# Data Flow Diagram in a typical application in TectoMT



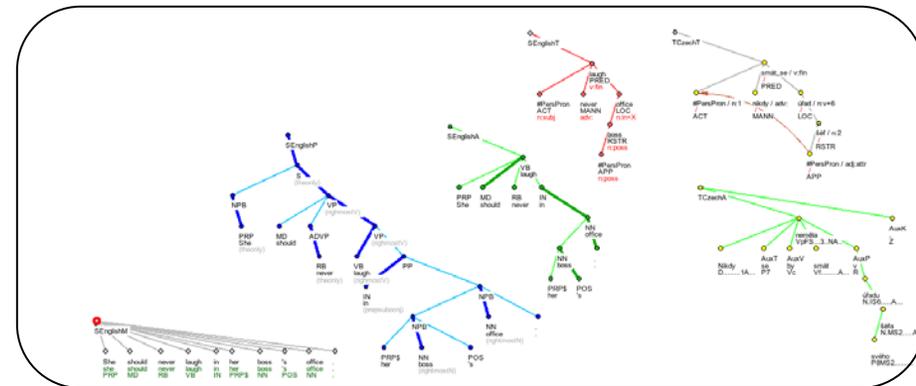
# Hierarchy of data-structure units

## ■ document

- the smallest independently storable unit (~ xml file)
- represents a text as a sequence of bundles, each representing one sentence (or sentence tuples in the case of parallel documents)

## ■ bundle

- set of tree representations of a given sentence



## ■ tree

- representation of a sentence on a given layer of linguistic description

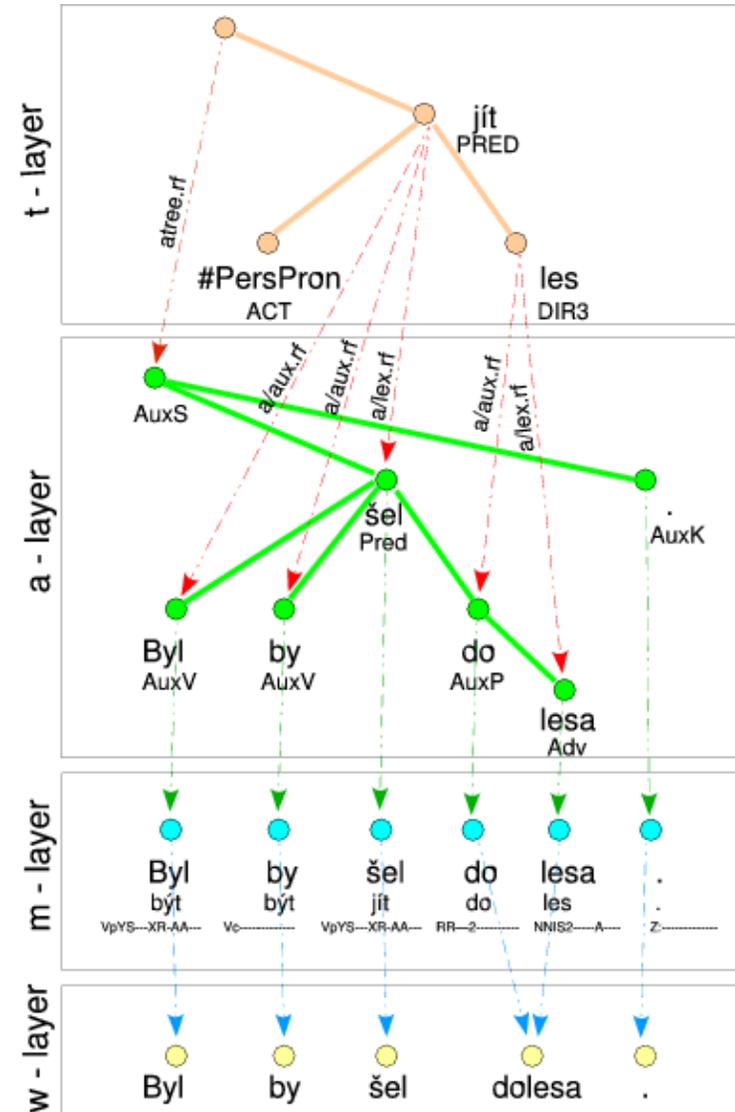
## ■ node

## ■ attribute

- document's, node's, or bundle's name-value pairs

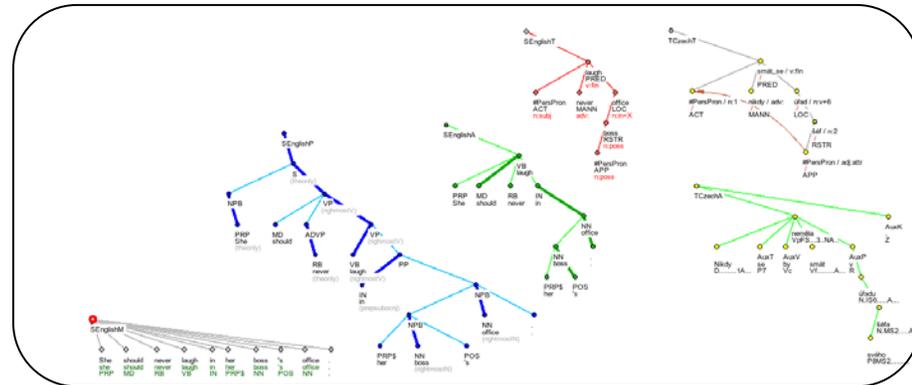
# Tree types adopted from PDT

- tectogrammatical layer
  - deep-syntactic dependency tree
- analytical layer
  - surface-syntactic dependency tree
  - 1 word (or punct.) ~ 1 node
- morphological layer
  - sequence of tokens with their lemmas and morphological tags



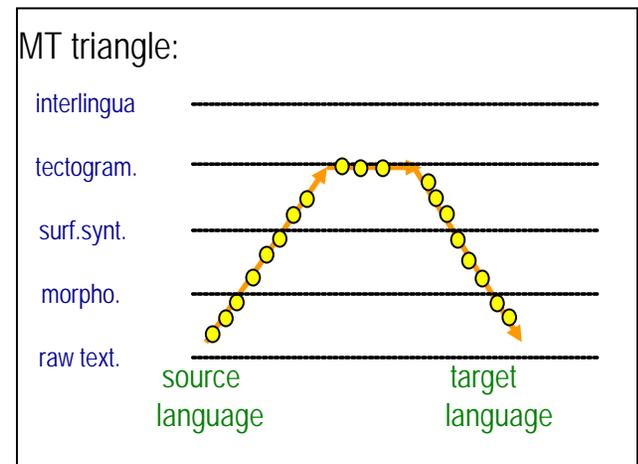
# Trees in a bundle

- in each bundle, there can be at most one tree for each "layer"
  - set of possible layers =  $\{S,T\} \times \{\text{English,Czech,...}\} \times \{M,A,T,P, N\}$
  - S - source, T-target (analysis vs. synthesis, MT perspective)
  - M - morphological analysis
  - P - phrase-structure tree
  - A - analytical tree
  - T - tectogrammatical tree
  - N - instances of named entities
- 
- Example: SEnglishA - tectogrammatical analysis of an English sentence on the source-language side



# Hierarchy of processing units

- block
  - the smallest individually executable unit
  - with well-defined input and output
  - block parametrization possible (e.g. model size choice)
- scenario
  - sequence of blocks, applied one after another on given documents
- application
  - typically 3 steps:
    1. conversion from the input format
    2. applying the scenario on the data
    3. conversion into the output format



# Blocks

- technically, Perl classes derived from `TectoMT::Block`
- either method `process_bundle` (if sentences are processed independently) or method `process_document` must be defined
- several hundreds blocks in TectoMT now, for various purposes:
  - blocks for analysis/transfer/synthesis, e.g.
    - `SEnglishW_to_SEnglishM::Lemmatize_mtree`
    - `SEnglishP_to_SEnglishA::Mark_heads`
    - `TCzechT_to_TCzechA::Vocalize_prepositions`
  - blocks for alignment, evaluation, feature extraction, etc.
- some of them only implement simple rules, some of them call complex probabilistic tools
- English-Czech tecto-based translation currently composes of roughly 140 blocks

# Tools available as TectoMT blocks

- to integrate a stand-alone NLP tool into TectoMT means to provide it with the standardized block interface
- already integrated tools:
  - taggers
    - Hajič's tagger, Raab&Spoustová Morče tagger, Rathnaparkhi MXPOST tagger, Brants's TnT tager, Schmid's Tree tagger, Coburn's `Lingua::EN::Tagger`
  - parsers
    - Collins' phrase structure parser, McDonalds dependency parser, Malt parser, ZŽ's dependency parser
  - named-entity recognizer
    - Stanford Named Entity Recognizer, Kravalová's SVM-based NE recognizer
  - miscel.
    - Klimeš's semantic role labeller, ZŽ's C5-based afun labeller, Ptáček's C5-based Czech preposition vocalizer, ...

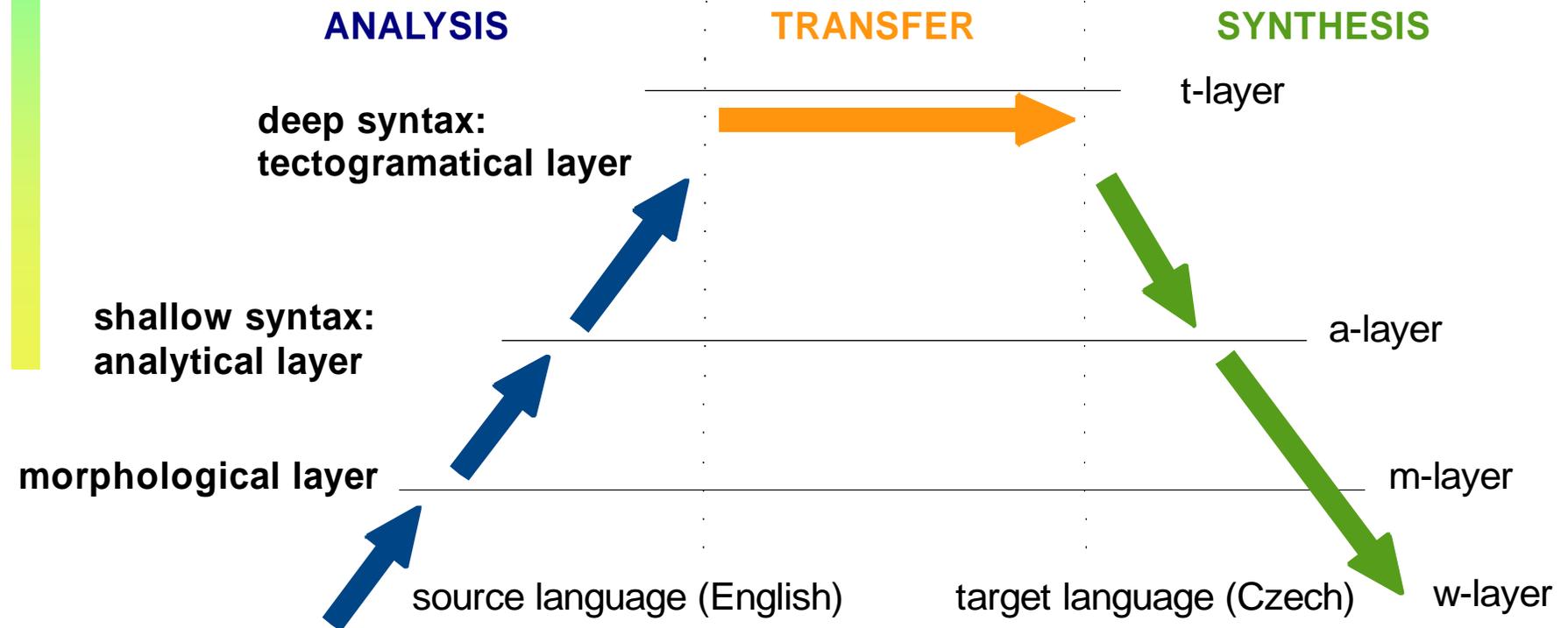
# Other TectoMT components

- "core" - Perl libraries forming the core of TectoMT infrastructure, esp. for memory representation of (and interface to) to the data structures
- numerous file-format converters (e.g. from PDT, Penn treebank, Czeg corpus, WMT shared task data etc. to our xml format)
- TectoMT-customized Pajas' tree editor TrEd
- tools for parallelized processing (Bojar)
- data, esp. trained models for the individual tools, morphological dictionaries, probabilistic translation dictionaries...
- tools for testing (regular daily tests), documentation...

# Languages in TectoMT

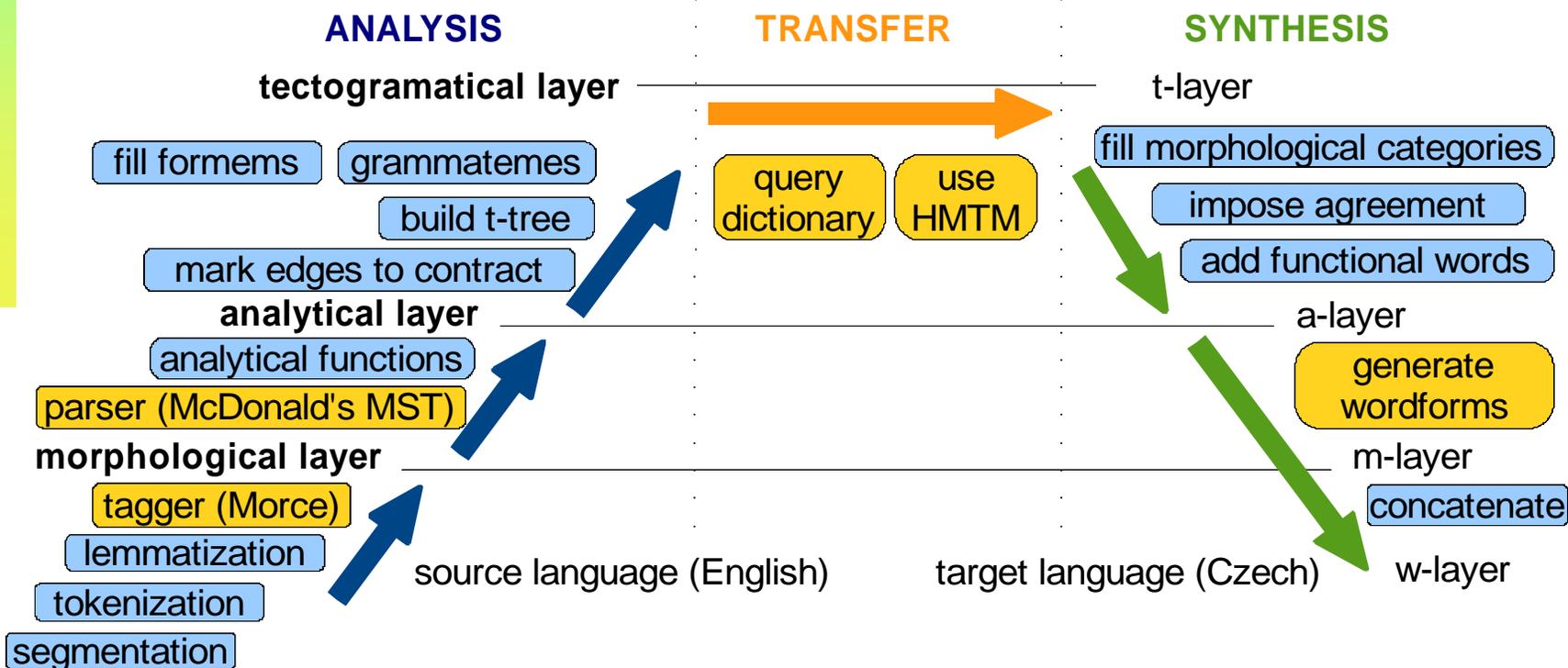
- full-fledged sentence PDT-style analysis/transfer/synthesis for English and Czech
  - using state-of-the-art tools
- prototype implementations of PDT-style analyses for a number of other languages
  - mostly created by students
  - Polish, French, German, Tamil, Spanish, Esperanto...

# English-Czech translation in TectoMT



# English-Czech translation in TectoMT

rule based & statistical blocks



# Real Translation Scenario

## **SEnglishW\_to\_SEnglishM::**

### **Tokenization**

Normalize\_forms

Fix\_tokenization

### **TagMorce**

Fix\_mtags

### **Lemmatize\_mtree**

## **SEnglishM\_to\_SEnglishN::**

Stanford\_named\_entities

Distinguish\_personal\_names

## **SEnglishM\_to\_SEnglishA::**

### **McD\_parser**

Fill\_is\_member\_from\_deprel

Fix\_tags\_after\_parse

McD\_parser REPARSE=1

Fill\_is\_member\_from\_deprel

Fix\_McD\_topology

Fix\_nominal\_groups

Fix\_is\_member

Fix\_atree

Fix\_multiword\_prep\_and\_conj

Fix\_dicendi\_verbs

Fill\_afun\_AuxCP\_Coord

### **Fill\_afun**

## **SEnglishA\_to\_SEnglishT::**

### **Mark\_edges\_to\_collapse**

Mark\_edges\_to\_collapse\_neg

### **Build\_tree**

Fill\_is\_member

Move\_aux\_from\_coord-  
\_to\_members

Fix\_tlemmas

Assign\_coap\_funcctors

Fix\_either\_or

Fix\_is\_member

Mark\_clause\_heads

Mark\_passives

Assign\_funcctors

Mark\_infin

Mark\_relclause\_heads

Mark\_relclause\_coref

Mark\_dsp\_root

Mark\_parentheses

Recompute\_deepord

Assign\_nodetype

### **Assign\_grammatemes**

### **Detect\_formeme**

Rehang\_shared\_attr

Detect\_voice

Fix\_imperatives

Fill\_is\_name\_of\_person

Fill\_gender\_of\_person

Add\_cor\_act

Find\_text\_coref

## **SEnglishT\_to\_TCzechT::**

### **Clone\_tree**

Translate\_LF\_phrases

Translate\_LF\_joint\_static

Delete\_superfluous\_tnodes

Translate\_F\_try\_rules

### **Translate\_F\_add\_variants**

Translate\_F\_rerank

Translate\_L\_try\_rules

### **Translate\_L\_add\_variants**

Translate\_LF\_numerals\_by\_rules

Translate\_L\_filter\_aspect

Transform\_passive\_constructions

Prune\_personal\_name\_variants

Remove\_unpassivizable\_variants

Translate\_LF\_compounds

Cut\_variants

Rehang\_to\_eff\_parents

### **Translate\_LF\_tree\_Viterbi**

Rehang\_to\_orig\_parents

Fix\_transfer\_choices

Translate\_L\_female\_surnames

Add\_noun\_gender

Add\_relpron\_below\_rc

Change\_Cor\_to\_PersPron

Add\_PersPron\_below\_vfin

Add\_verb\_aspect

Fix\_date\_time

Fix\_grammatemes\_after\_transfer

Fix\_negation

Move\_adjectives\_before\_nouns

Move\_genitives\_to\_postposit

Move\_relclause\_to\_postposit

Move\_dicendi\_closer\_to\_dsp

Move\_PersPron\_next\_to\_verb

Move\_enough\_before\_adj

Fix\_money

Recompute\_deepord

Find\_gram\_coref\_for\_refl\_pron

Neut\_PersPron\_gender\_from\_antec

Override\_pp\_with\_phrase\_translation

Valency\_related\_rules

Fill\_clause\_number

Turn\_text\_coref\_to\_gram\_coref

## **TCzechT\_to\_TCzechA::**

### **Clone\_atree**

Distinguish\_homonymous\_mlemmas

Reverse\_number\_noun\_dependency

### **Init\_morphcat**

Fix\_possessive\_adjectives

Mark\_subject

Impose\_pron\_z\_agr

Impose\_rel\_pron\_agr

Impose\_subjpred\_agr

Impose\_attr\_agr

Impose\_compl\_agr

Drop\_subj\_pers\_prons

Add\_prepositions

Add\_subconjs

Add\_reflex\_particles

Add\_auxverb\_compound\_passive

Add\_auxverb\_modal

Add\_auxverb\_compound\_future

Add\_auxverb\_conditional

Add\_auxverb\_compound\_past

Add\_clausal\_expletive\_pronouns

Resolve\_verbs

Project\_clause\_number

Add\_parentheses

Add\_sent\_final\_punct

Add\_subord\_clause\_punct

Add\_coord\_punct

Add\_apposition\_punct

Choose\_mlemma\_for\_PersPron

### **Generate\_wordforms**

### **Move\_clitics\_to\_wackemagel**

Recompute\_ordering

Delete\_superfluous\_prepos

Delete\_empty\_nouns

Vocalize\_prepositions

Capitalize\_sent\_start

Capitalize\_named\_entities

## **TCzechA\_to\_TCzechW::**

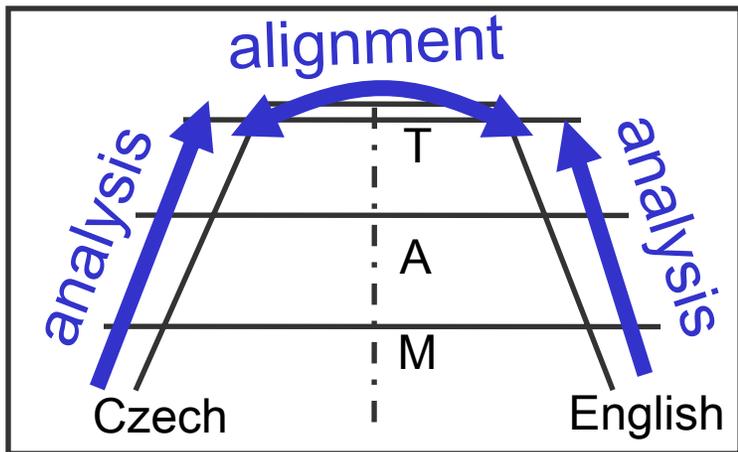
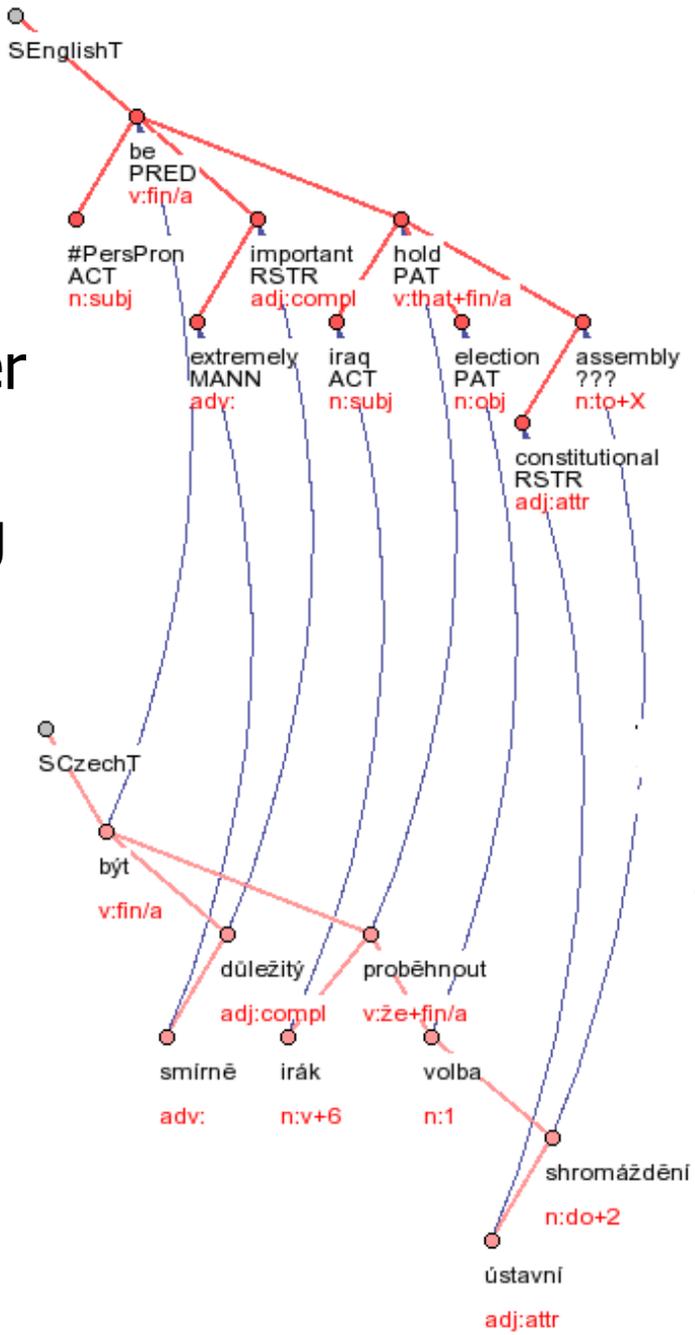
### **Concatenate\_tokens**

Ascii\_quotes

Remove\_repeated\_tokens

# Parallel analysis

- data needed for training the transfer phase models
- Czech-English parallel corpus CzEng
- 8 mil. pairs of sentences with automatic PDT-style analyses and alignment



# Summary of Part I

---

- TectoMT (→Treeex)
  - environment for NLP experiments
  - multipurpose, multilingual
  - PDT-style linguistic structures
  - Linux+Perl, open-source
  - modular architecture (several hundreds of modules)
  - capable of processing massive data
  - will be released at CPAN