# English-Czech Machine Translation Using TectoMT

## Martin Popel
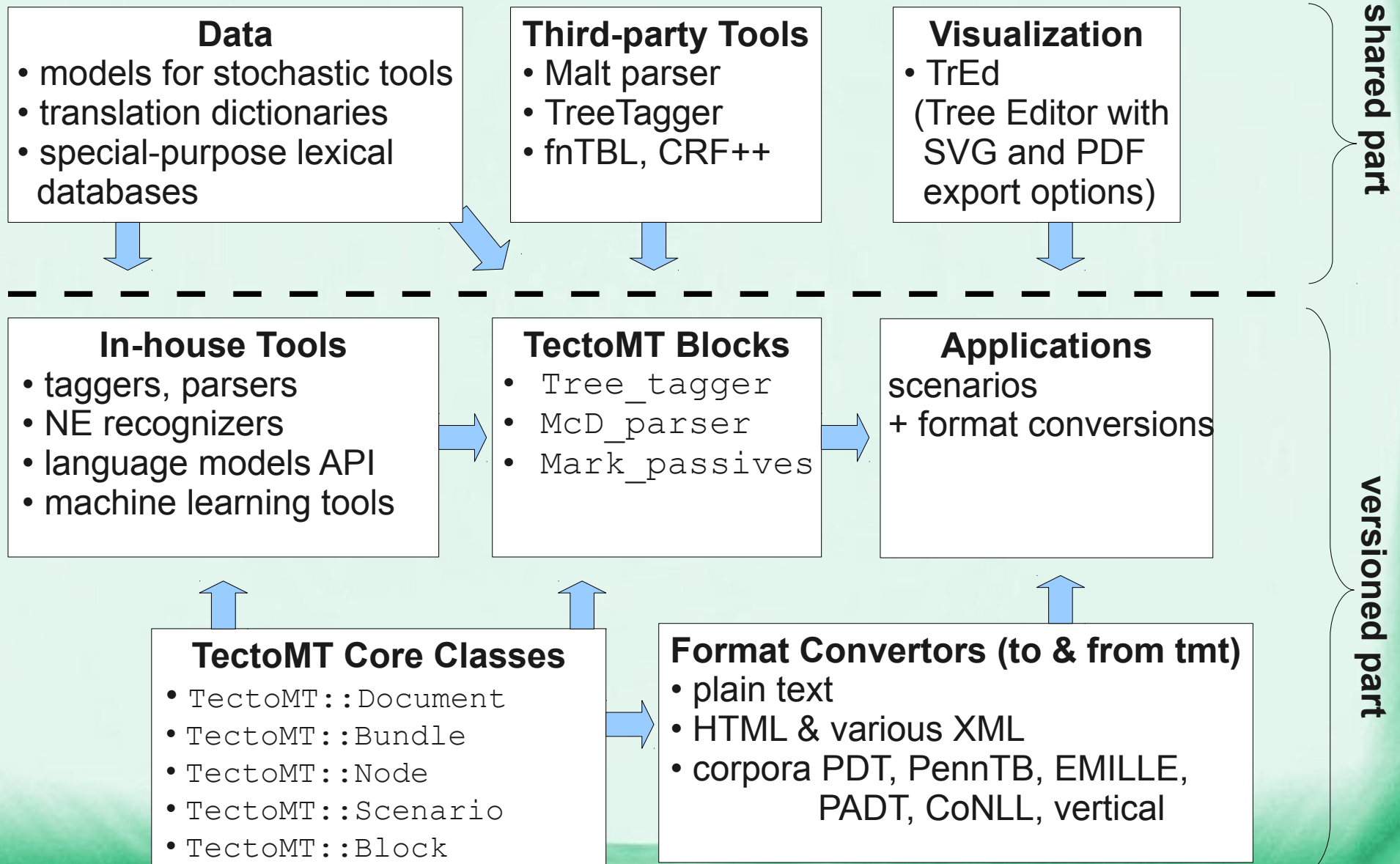ÚFAL, Charles University in Prague

19th Week of Doctoral Students, June 3, 2010

# Outline

- TectoMT – NLP framework and MT system

- Demo translation step by step

- Hidden Markov Tree Models (HMTM)

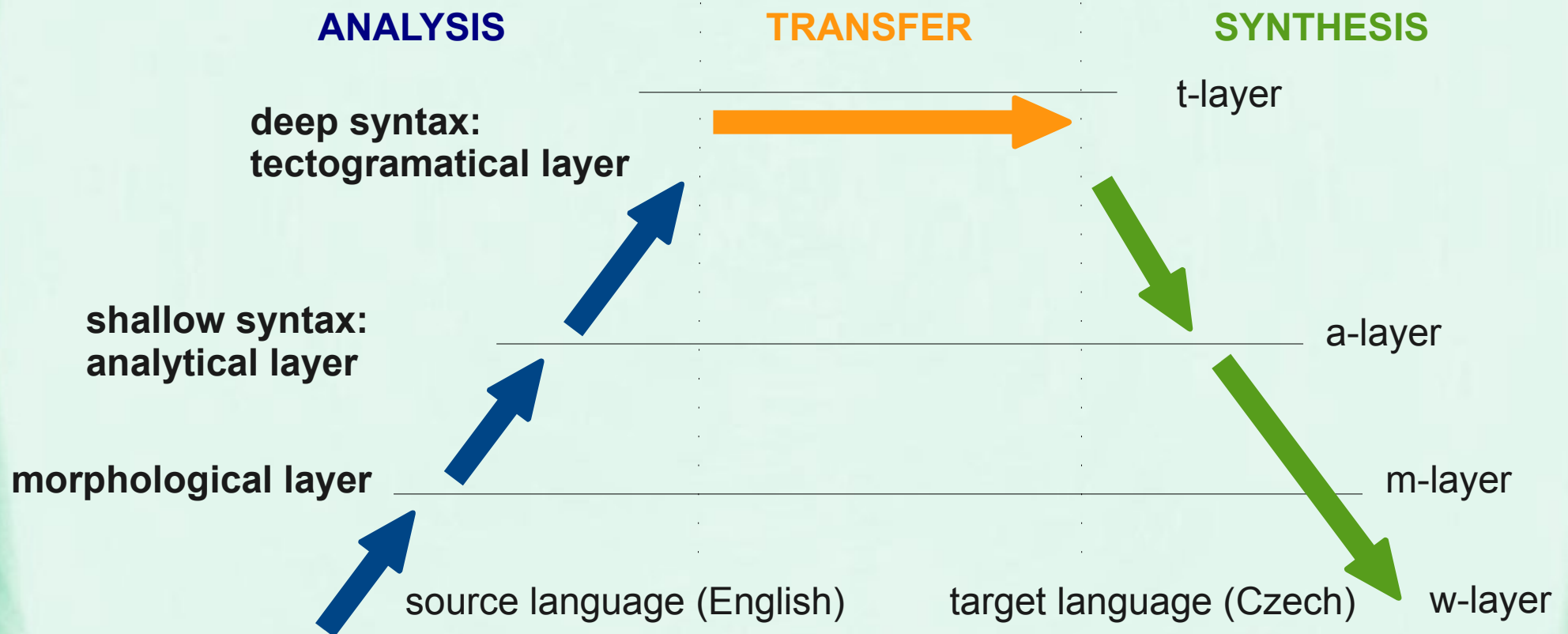- Results – three metrics of translation quality

# TectoMT – framework for NLP

**TectoMT**
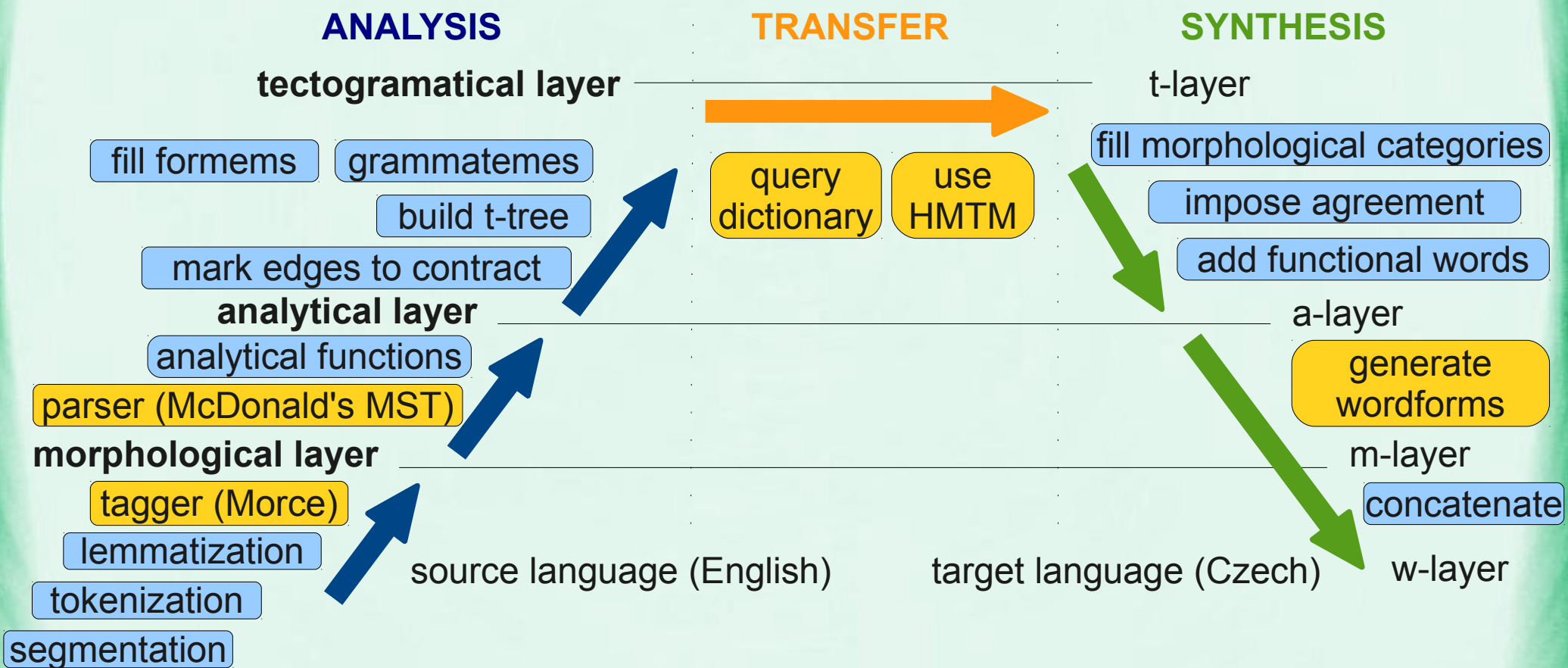
## modular, open source, Perl, Linux, OOP-style

**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)

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

**TectoMT Blocks**
- `Tree_tagger`
- `McD_parser`
- `Mark_passives`

**Applications**
scenarios
+ format conversions

**TectoMT Core Classes**
- `TectoMT::Document`
- `TectoMT::Bundle`
- `TectoMT::Node`
- `TectoMT::Scenario`
- `TectoMT::Block`

**Format Convertors (to & from tmt)**
- plain text
- HTML & various XML
- corpora PDT, PennTB, EMILLE, PADT, CoNLL, vertical

# TectoMT – MT system

## transfer over the tectogrammatical layer

# TectoMT – MT system
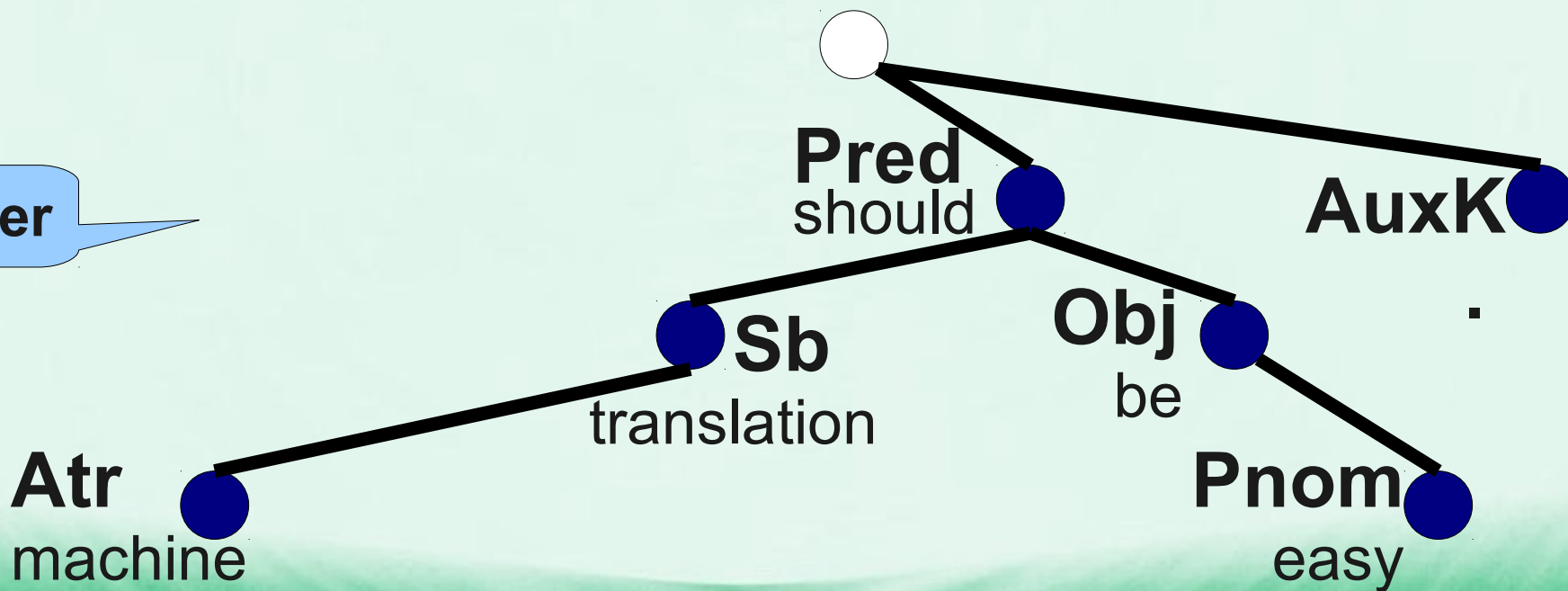
TectoMT

rule based & statistical blocks

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

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

query dictionary   use HMTM

fill formems   grammatemes

build t-tree

fill morphological categories

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)

lemmatization

concatenate

tokenization

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

segmentation

# Demo Translation – Analysis

**TectoMT**

**raw text** — Machine translation should be easy.

**m-layer**

machine    translation    should    be    easy  .

**NN**    **NN**    **MD**    **VB**    **JJ**    **.**

**a-layer**

**Pred**
should

**AuxK**

**Sb**
translation

**Obj**
be

**Atr**
machine

**Pnom**
easy

# Demo Translation – Analysis

TectoMT

**raw text** Machine translation should be easy.

**m-layer**

machine **NN** translation **NN** should **MD** be **VB** easy **JJ** . **.**

Mark functional words

**a-layer**

Pred
should

AuxK

Sb
translation

Obj
be

Atr
machine

Pnom
easy

# Demo Translation – Analysis

**TectoMT**

**raw text**

Machine translation should be easy.

**m-layer**

machine translation should be easy .
**NN** **NN** **MD** **VB** **JJ** **.**

Build t-tree (backbone)

**t-layer**

be

translation easy

machine

# Demo Translation – Analysis

TectoMT

**raw text**

Machine translation should be easy.

**m-layer**

machine **NN**  translation **NN**  should **MD**  be **VB**  easy **JJ**  . **.**

Fill formems

**t-layer**

be **v:fin**

translation **n:subj**

easy **adj:compl**

machine **n:attr**

# Demo Translation – Analysis

TectoMT

**raw text**

Machine translation should be easy.

**m-layer**

machine **NN**   translation **NN**   should **MD**   be **VB**   easy **JJ**   . **.**

Fill grammatemes

**t-layer**

**tense = simple, modality, conditional**

be **v:fin**

**n:attr** machine

translation **n:subj**
**number = singular**

easy **adj:compl**
**degree of comparison = positive**

# Demo Translation – Transfer

**Build target t-tree by cloning**

source t-layer

be **v:fin**

translation **n:subj**

easy **adj:compl**

machine **n:attr**

target t-layer

be **v:fin**

translation **n:subj**

easy **adj:compl**

machine **n:attr**

# Demo Translation – Transfer

TectoMT

Get translation variants
for lemmas and formems

source t-layer

be **v:fin**

translation **n:subj**

easy **adj:compl**

machine **n:attr**

target t-layer

být
mít
**v:fin**
**v:inf**

překlad
převod
**n:1**

snadný
jednoduchý
**adj:compl**
**n:1**
**adv:**

počítač
stroj
strojový
**n:2**
**n:attr**
**adj:attr**

# Demo Translation – Transfer

# Demo Translation – Synthesis

TectoMT

Build target a-layer by cloning

target t-layer

být **v:fin**

překlad **n:1**

snadný **adj:compl**

strojový **adj:attr**

target a-layer

být

překlad

snadný

strojový

# Demo Translation – Synthesis

TectoMT

Fill morphological categories

target t-layer

být **v:fin**

překlad **n:1**

snadný **adj:compl**

strojový **adj:attr**

target a-layer

být

překlad

snadný

**degree = positive**

strojový
**degree = positive**

**number = singular**
**gender = masc. inanim.**
**case = nominative**

# Demo Translation – Synthesis

TectoMT

Impose agreement

target t-layer

být — v:fin

překlad — n:1

snadný — adj:compl

strojový — adj:attr

target a-layer

number = singular
case = nominative
gender = masc. inanim.

strojový
degree = positive

překlad
number = singular
gender = masc. inanim.
case = nominative

být
number = singular
gender = masc. Inanim.

snadný
degree = positive
number = singular
case = nominative
gender = masc. inanim.

# Demo Translation – Synthesis

TectoMT

Add functional words

target t-layer

být **v:fin**

překlad **n:1**

snadný **adj:compl**

strojový **adj:attr**

target a-layer

mít

překlad

strojový

by být snadný

# Demo Translation – Synthesis

TectoMT

Reorder clitics

target t-layer

být  v:fin

překlad  n:1

snadný  adj:compl

strojový  adj:attr

target a-layer

mít

překlad

by

být  snadný

strojový

# Demo Translation – Synthesis

TectoMT

Generate wordforms

target t-layer

být  v:fin

překlad  n:1

snadný  adj:compl

strojový  adj:attr

target a-layer

měl

překlad

strojový

by

být  snadný

# Demo Translation – Synthesis

Concatenate tokens for output

target t-layer

být **v:fin**

překlad **n:1**

snadný **adj:compl**

strojový **adj:attr**

target a-layer

měl

překlad

strojový

by

být snadný

**Strojový překlad by měl být snadný.**

# Demo Translation – Real Scenario

TectoMT

**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_ttree
Fill_is_member
Move_aux_from_coord-
_to_members
Fix_tlemmas
Assign_coap_functors
Fix_either_or
Fix_is_member

Mark_clause_heads
Mark_passives
Assign_functors
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_ttree
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_wackernagel
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

# HMTM – Motivation

# HMTM – Motivation

**TectoMT**

Select the best label for each node

source t-layer

be **v:fin**

translation **n:subj**

easy **adj:compl**

machine **n:attr**

target t-layer

být|**v:fin**, být|**v:inf**, mít|**v:fin**, mít|**v:inf**

překlad|**n:1**, převod|**n:1**

snadný|**adj:compl**, jednoduchý|**adj:compl**, ...

počítač|**n:2**, počítač|**n:attr**, strojový|adj:attr, ...

# HMTM in MT

**TectoMT**

Source tree (Czech) ──TRANSFER──▶ Target tree (English)

$$P(\text{optimal\_tree}) = P_E(\text{strojový} \mid \text{machine}) \cdot P_T(\text{machine} \mid \text{translation}) \cdot$$
$$P_E(\text{překlad} \mid \text{translation}) \cdot P_T(\text{translation} \mid \text{be}) \cdot$$
$$P_E(\text{snadný} \mid \text{easy}) \cdot P_T(\text{easy} \mid \text{be}) \cdot$$
$$P_E(\text{být} \mid \text{be}) \cdot P_T(\text{be} \mid \text{ROOT})$$

ROOT

být

$P_E(\text{být} \mid \text{have}) = \textbf{0.01}$

$P_E(\text{být} \mid \text{be}) = \textbf{0.8}$

překlad   snadný

$P_E(\text{překlad} \mid \text{arcade}) = \textbf{0.7}$

$P_E(\text{překlad} \mid \text{translation}) = \textbf{0.6}$

strojový

$P_T(\text{machine} \mid \text{translation}) = \textbf{0.02}$

$P_E(\text{strojový} \mid \text{machine}) = \textbf{0.4}$

$P_E(\text{strojový} \mid \text{engine}) = \textbf{0.5}$

**ROOT**

0.01

**be**   have

0.0001   1×10⁻³   0.002   0.001

**translation**   arcade

1×10⁻³   1×10⁻⁸

**machine**   engine

**easy**   simple

ANALYSIS

SYNTHESIS

Source sentence:
*Strojový překlad by měl být snadný.*

Target sentence:
*Machine translation should be easy.*

$P_E(\text{source} \mid \text{target})$ … emission probabilities … **translation model**

$P_T(\text{dependent} \mid \text{governing})$ … transition probabilities … **target-language tree model**

# Results – BLEU
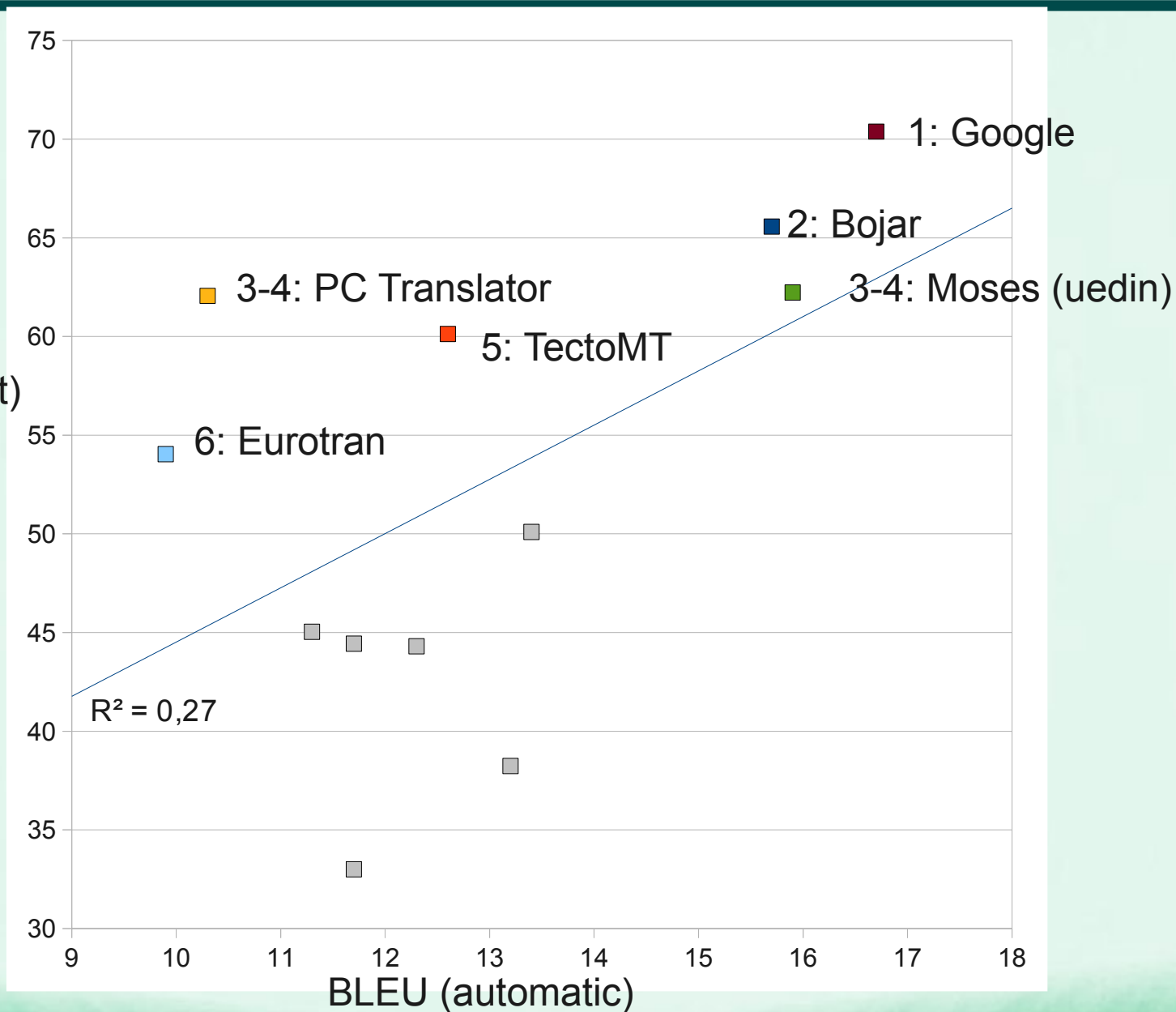


WMT = Workshop on Statistical Machine Translation

Legend:
- bojar
- tectomt
- pctrans
- uedin
- google
- eurotran

TectoMT:

| Year | | BLEU |
|------|------|------|
| 2008 | WMT | 6,9 |
| 2009 | WMT | 7,3 |
| 2009/09 | | 10,2 |
| 2010/01 | | 10,4 |
| 2010/02 | | 11,3 |
| 2010/03 | WMT | 12,6 |

# Results – BLEU vs. Ranks

# Results – task-based evaluation

## Example

- **Translated text:**

  *Bílý dům náčelník štábu Rahm Emanuel je naplánovaný k pojídaný vsedě v budově Kongresu se sněmovnou diskutující Nancy Pelosi v 9 p.m.*

- **Yes/No Statements:**

  - *Nancy Pelosi je zaměstnankyní Bílého domu.*

  - *Schůzka se koná dopoledne.*

  - *Schůzka se koná v Bílem domě.*

## Results

| TectoMT | PC Translator | Google | Moses |
|---------|---------------|--------|-------|
| 80.56   | 80.23         | 78.68  | 73.58 |

Jan Berka, Martin Černý, 2010: 18 annotators, 1905 * 3 statements, 4 systems evaluated

[https://svn.ms.mff.cuni.cz/projects/NPFL087/export/598/projects/task-based-evaluation/tbe.pdf]
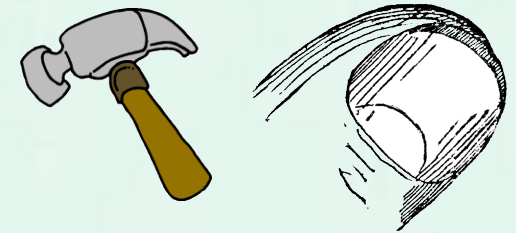
# Examples of Translation

A miss by an inch is a miss by a mile.

Slečna palec je slečna miliónu.

I'd rather be a hammer than a nail.

Spíše bych byl kladivo než nehet.

A bird in the hand is worth two in the bush.

Pták v ruce je cenný dvakrát v Bushovi.