

TectoMT: Machine Translation System

Zdeněk Žabokrtský, Martin Popel

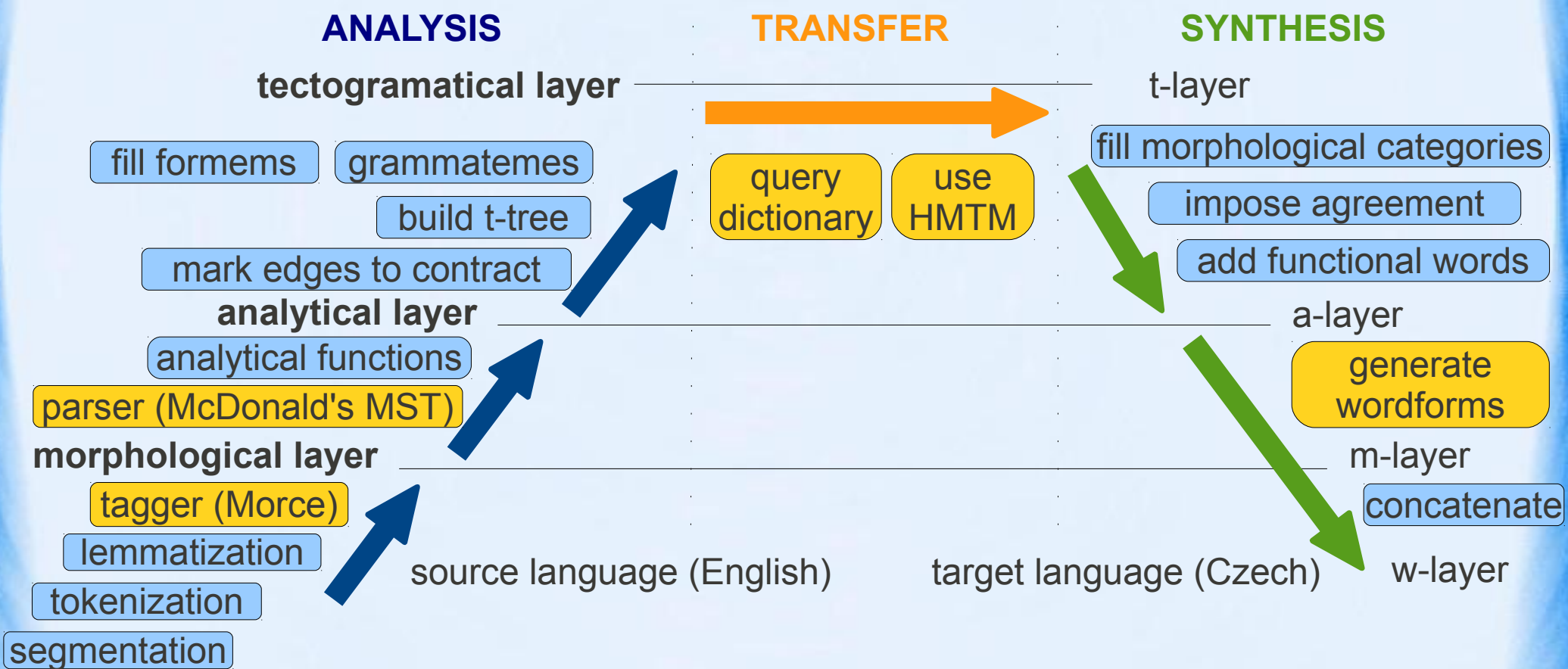
ÚFAL (Institute of Formal and Applied Linguistics)
Charles University in Prague



META-NET Course on Advanced MT Resources,
Prague Dec. 17, 2010

Translation scheme

rule based & statistical blocks



Demo Translation – Analysis

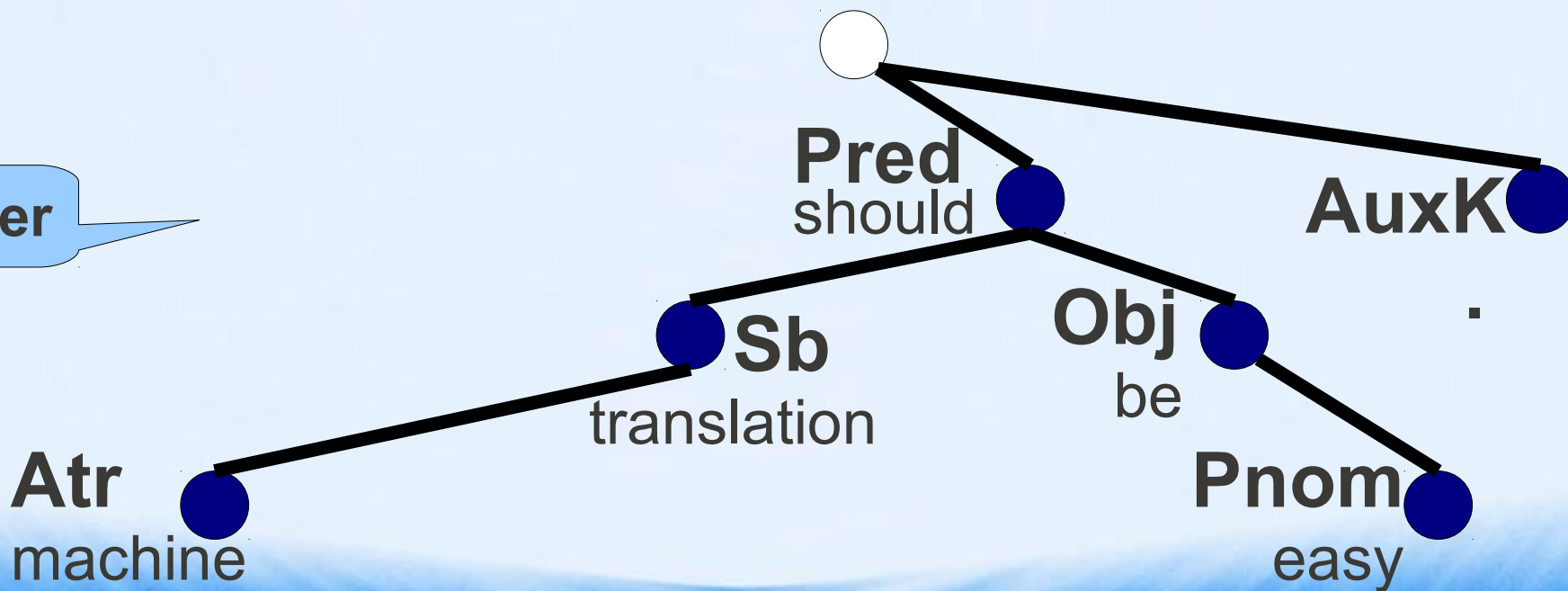
raw text

Machine translation should be easy.

m-layer

| | | | | | |
|---------|-------------|--------|----|------|---|
| ● | ● | ● | ● | ● | ● |
| machine | translation | should | be | easy | . |
| NN | NN | MD | VB | JJ | . |

a-layer



Demo Translation – Analysis

raw text

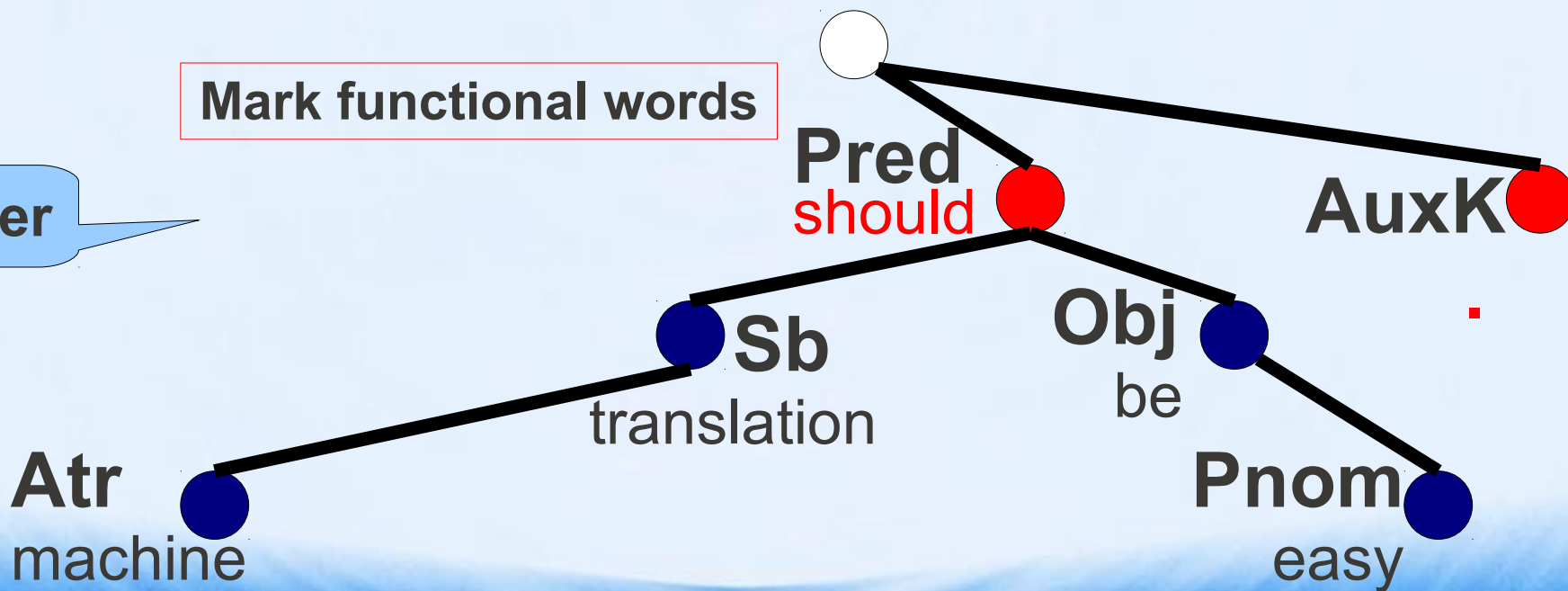
Machine translation should be easy.

m-layer

| | | | | | |
|---------|-------------|--------|----|------|---|
| ● | ● | ● | ● | ● | ● |
| machine | translation | should | be | easy | . |
| NN | NN | MD | VB | JJ | . |

Mark functional words

a-layer



Demo Translation – Analysis

raw text

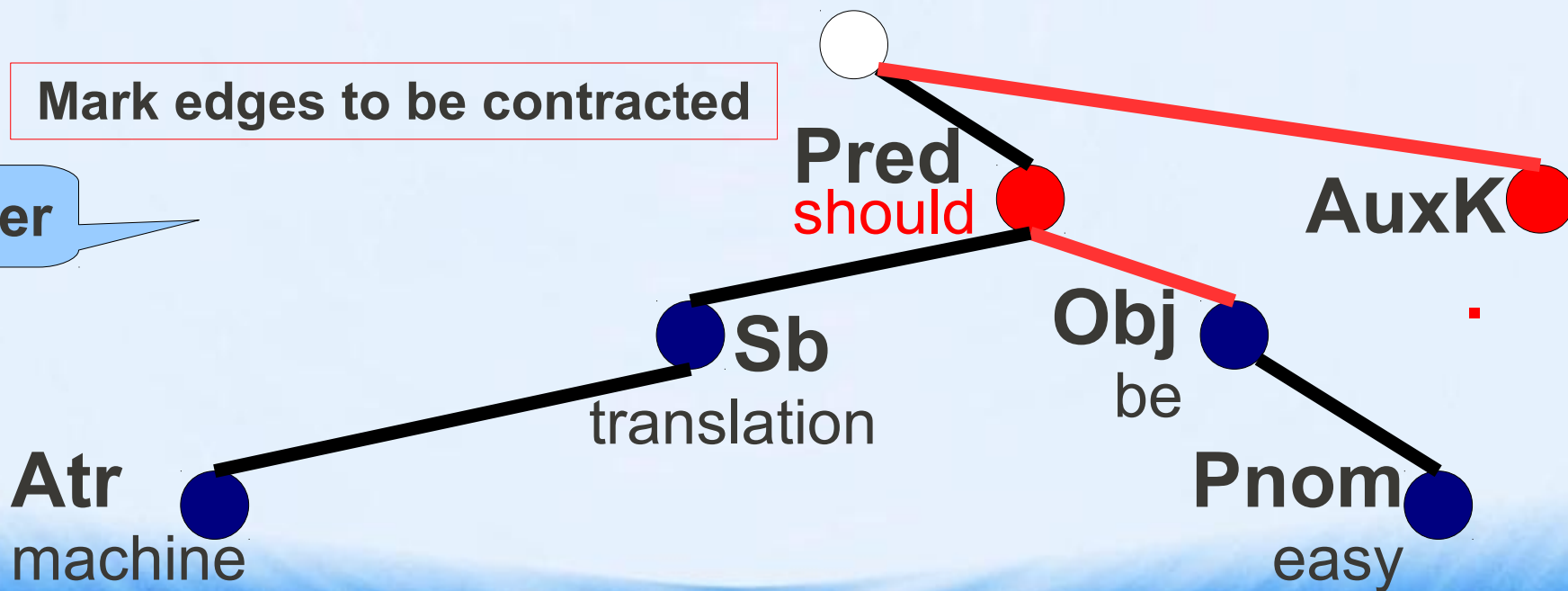
Machine translation should be easy.

m-layer

| | | | | | |
|---------|-------------|--------|----|------|---|
| ● | ● | ● | ● | ● | ● |
| machine | translation | should | be | easy | . |
| NN | NN | MD | VB | JJ | . |

Mark edges to be contracted

a-layer



Demo Translation – Analysis

raw text

Machine translation should be easy.

m-layer

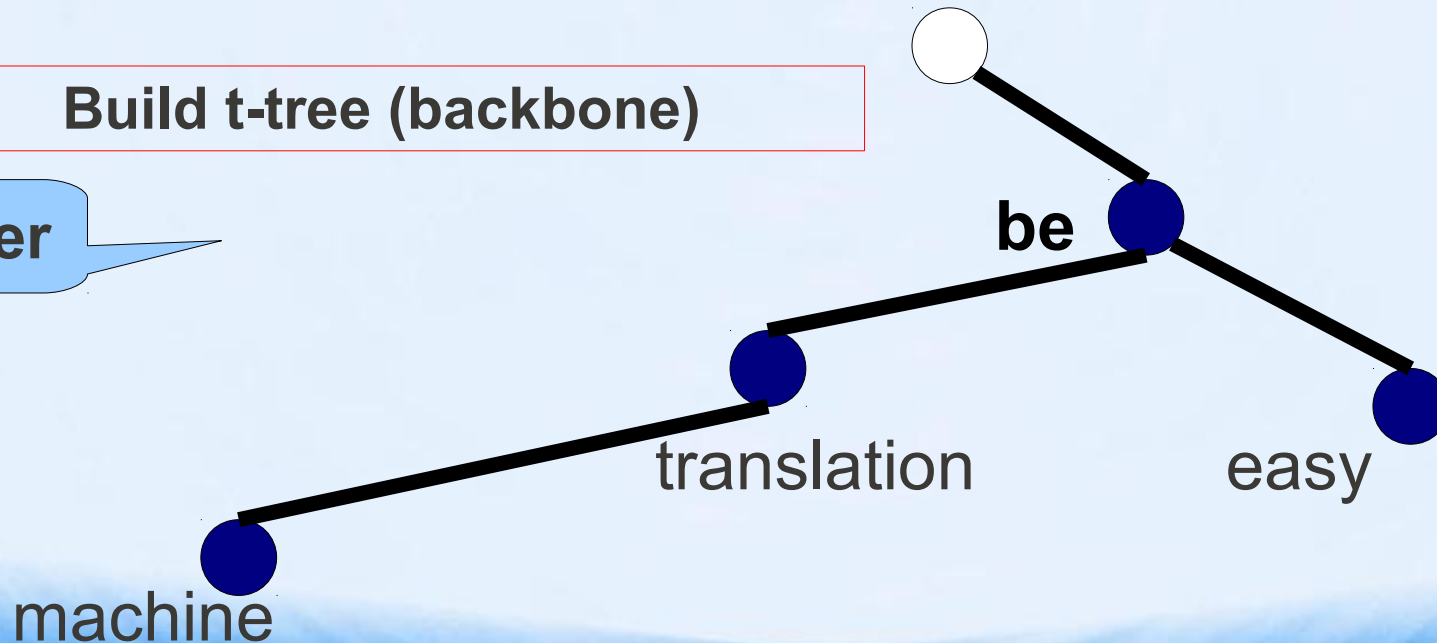
● ● ● ● ● ●

machine translation should be easy .

 NN NN MD VB JJ .

Build t-tree (backbone)

t-layer



Demo Translation – Analysis

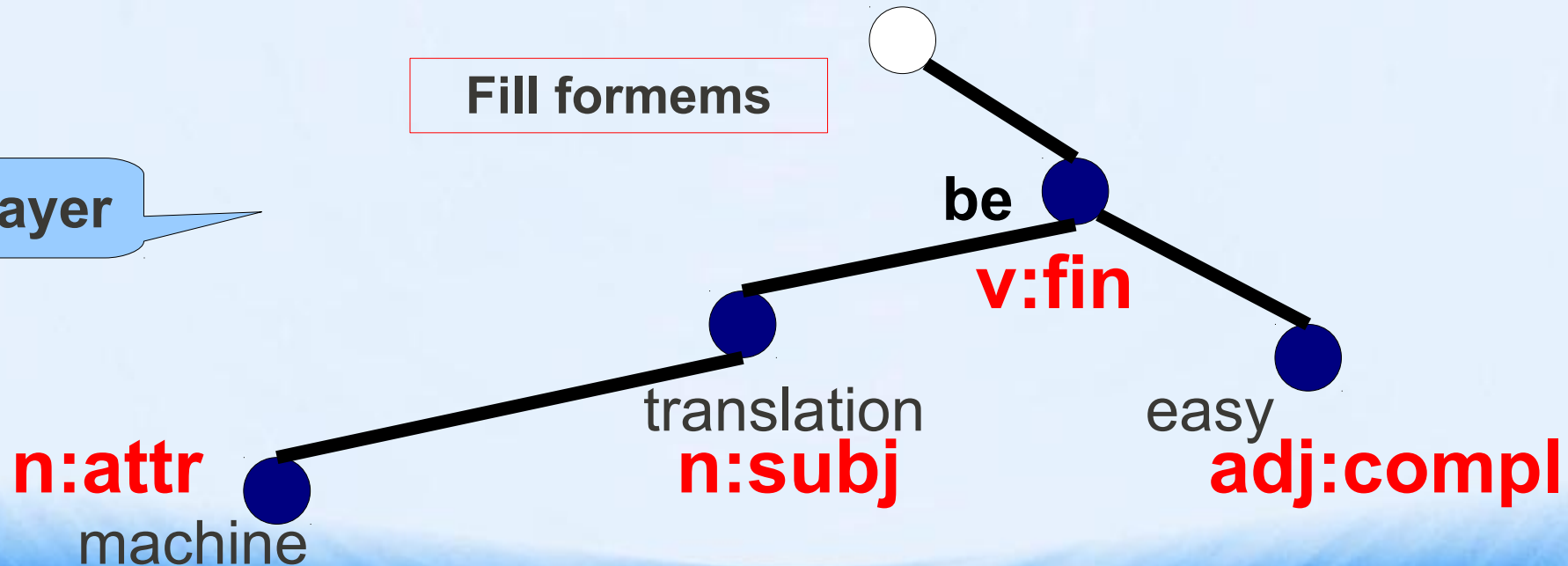
raw text

Machine translation should be easy.

m-layer

| | | | | | |
|---------|-------------|--------|----|------|---|
| ● | ● | ● | ● | ● | ● |
| machine | translation | should | be | easy | . |
| NN | NN | MD | VB | JJ | . |

t-layer



Demo Translation – Analysis

raw text

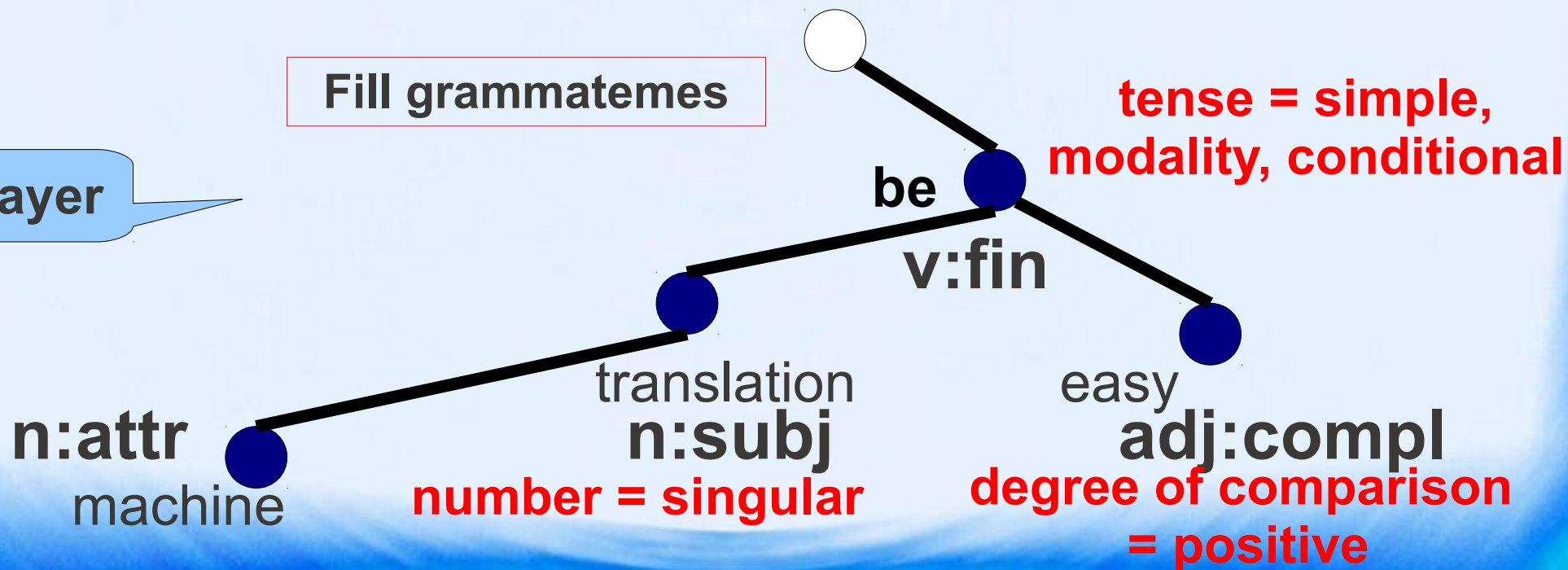
Machine translation should be easy.

m-layer

| | | | | | |
|---------|-------------|--------|----|------|---|
| ● | ● | ● | ● | ● | ● |
| machine | translation | should | be | easy | . |
| NN | NN | MD | VB | JJ | . |

Fill grammatememes

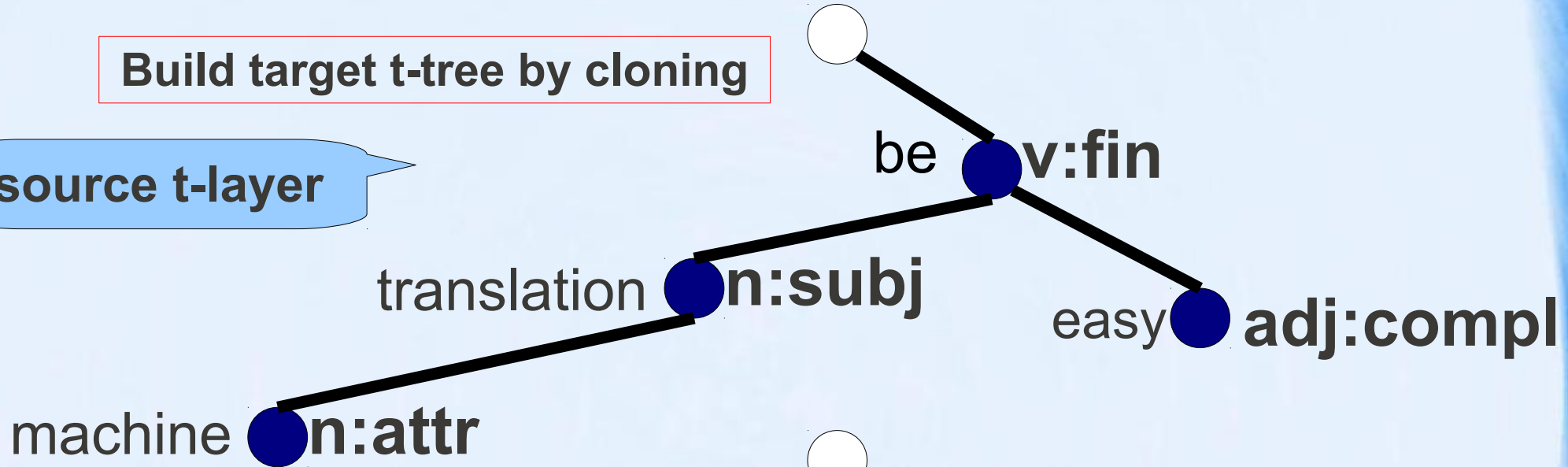
t-layer



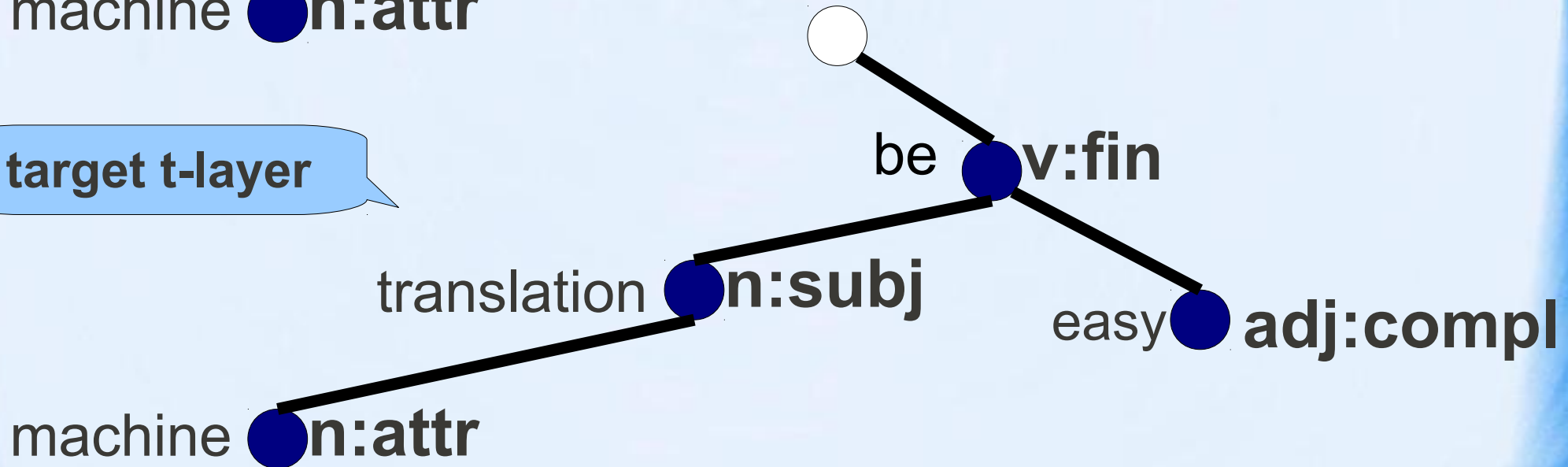
Demo Translation – Transfer

Build target t-tree by cloning

source t-layer



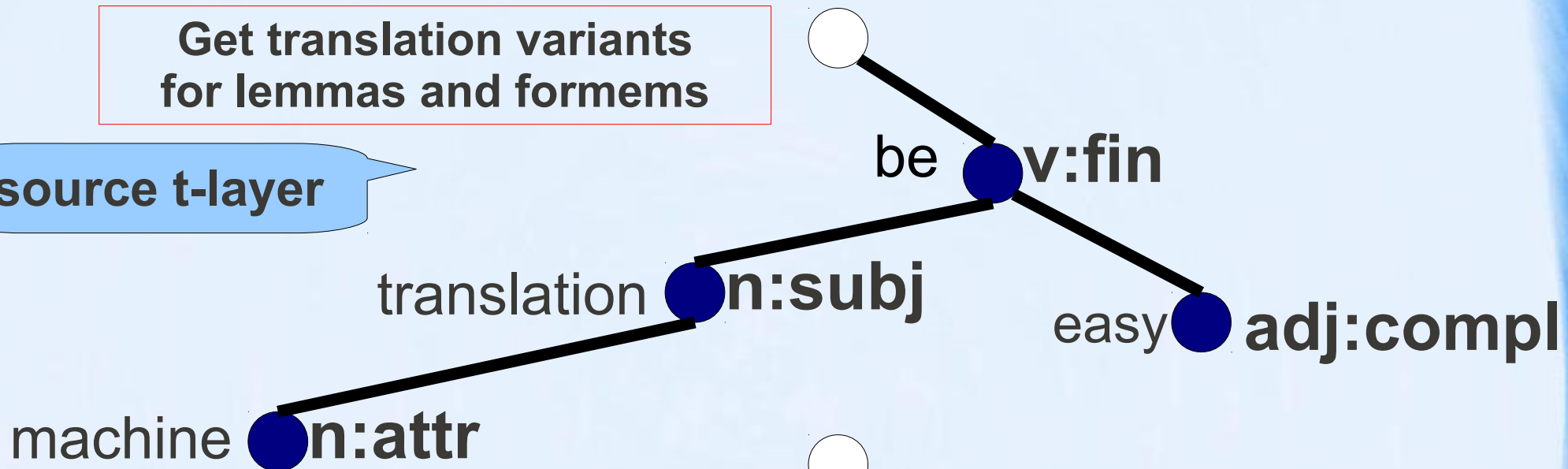
target t-layer



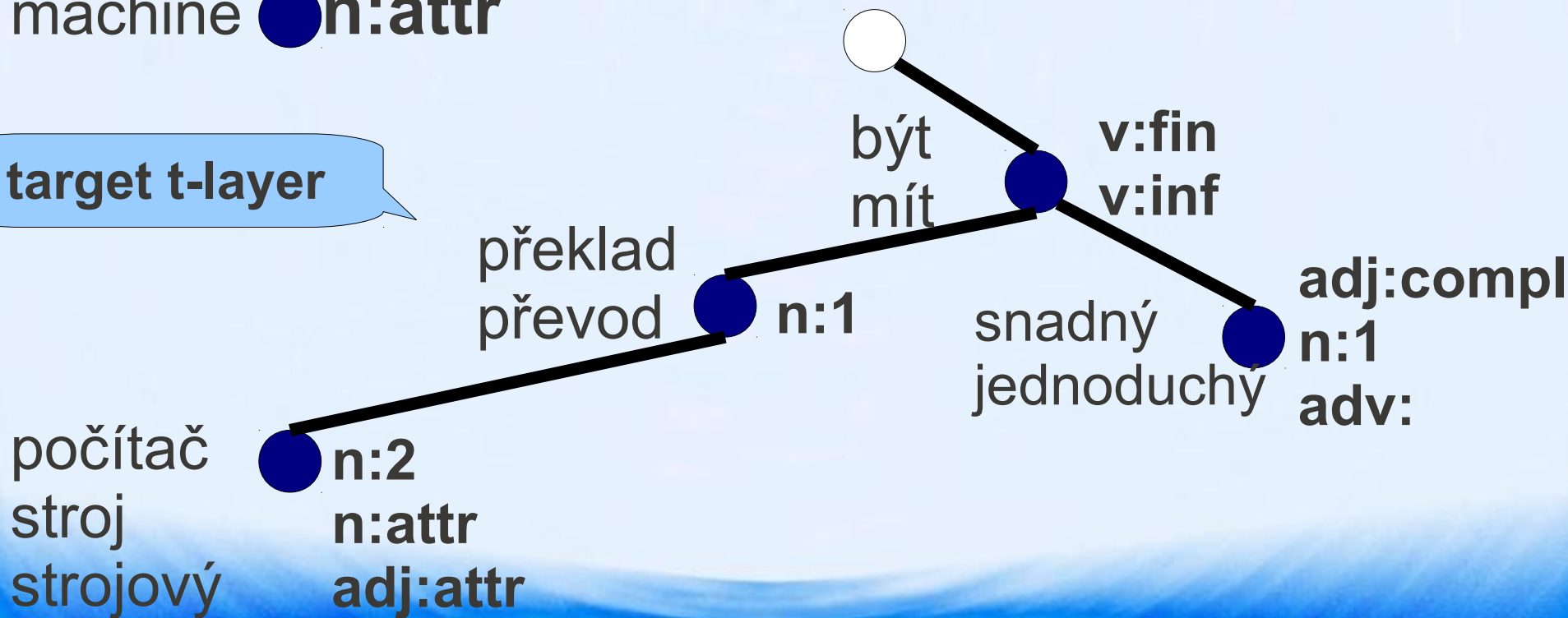
Demo Translation – Transfer

Get translation variants for lemmas and formems

source t-layer



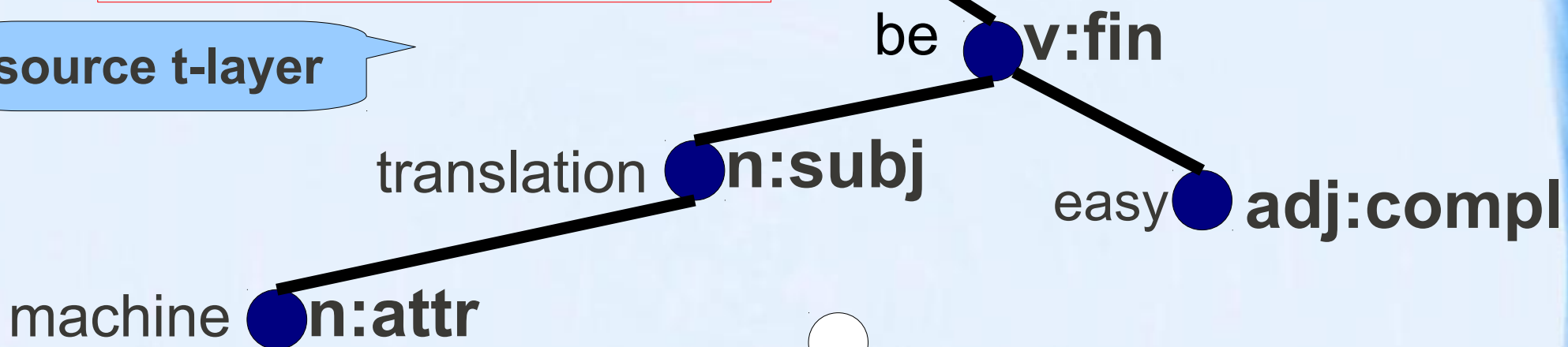
target t-layer



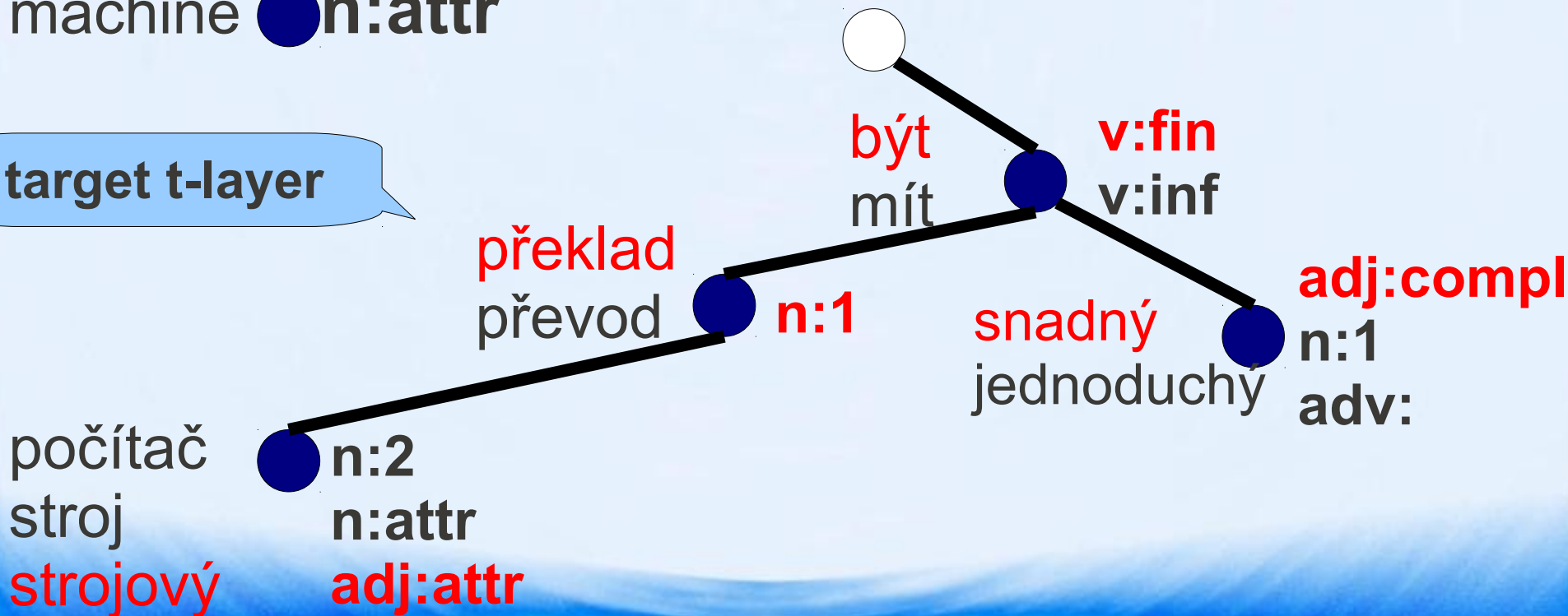
Demo Translation – Transfer

Select the best combination of lemmas and formems

source t-layer



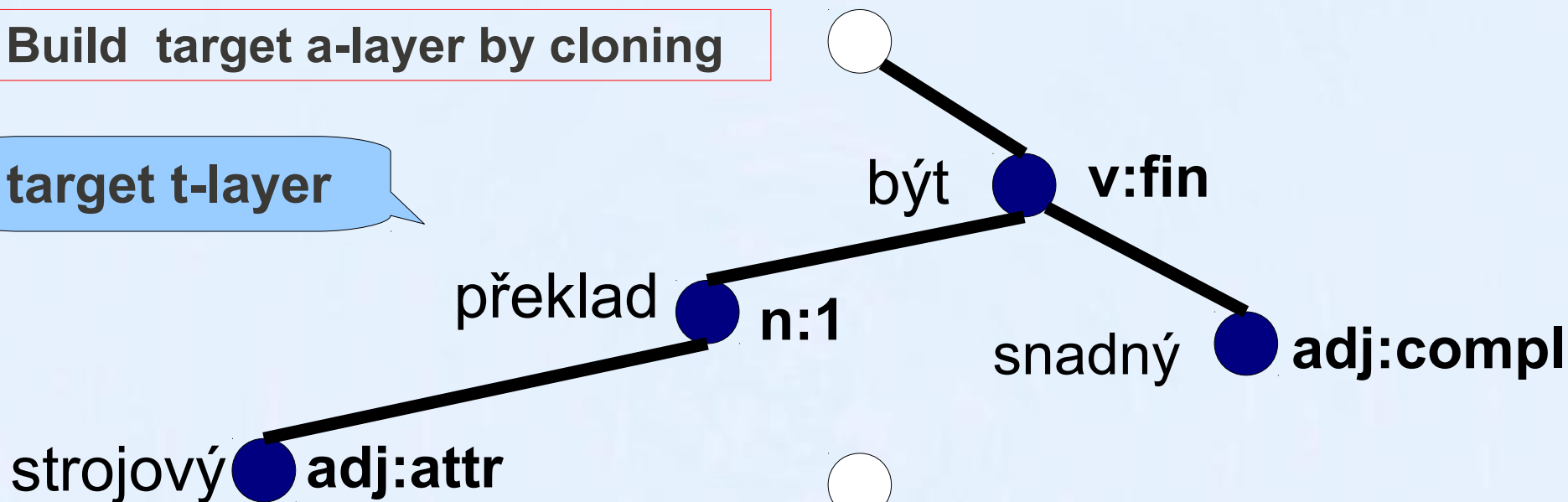
target t-layer



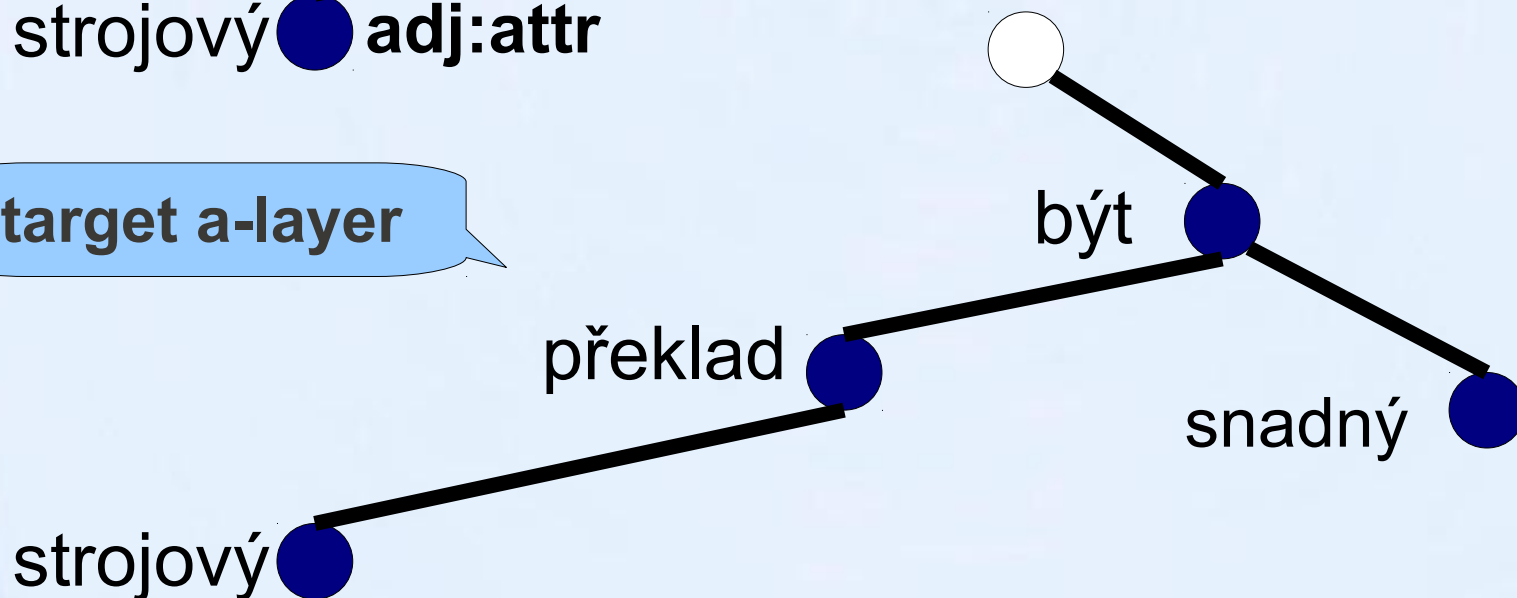
Demo Translation – Synthesis

Build target a-layer by cloning

target t-layer



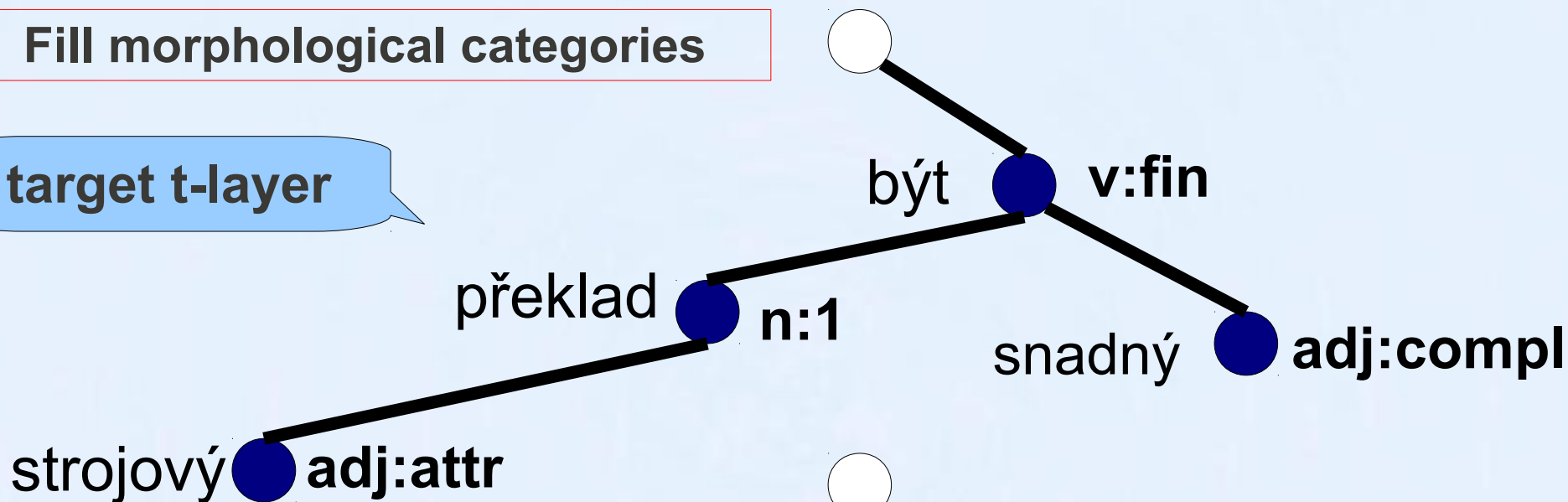
target a-layer



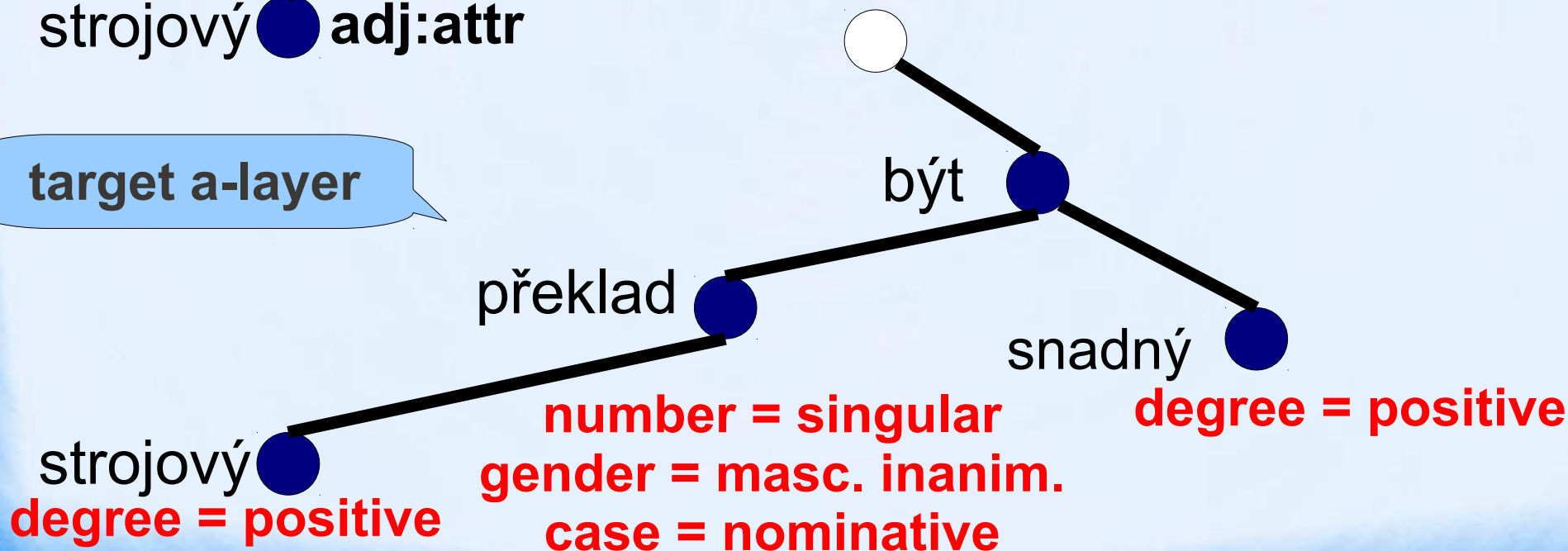
Demo Translation – Synthesis

Fill morphological categories

target t-layer



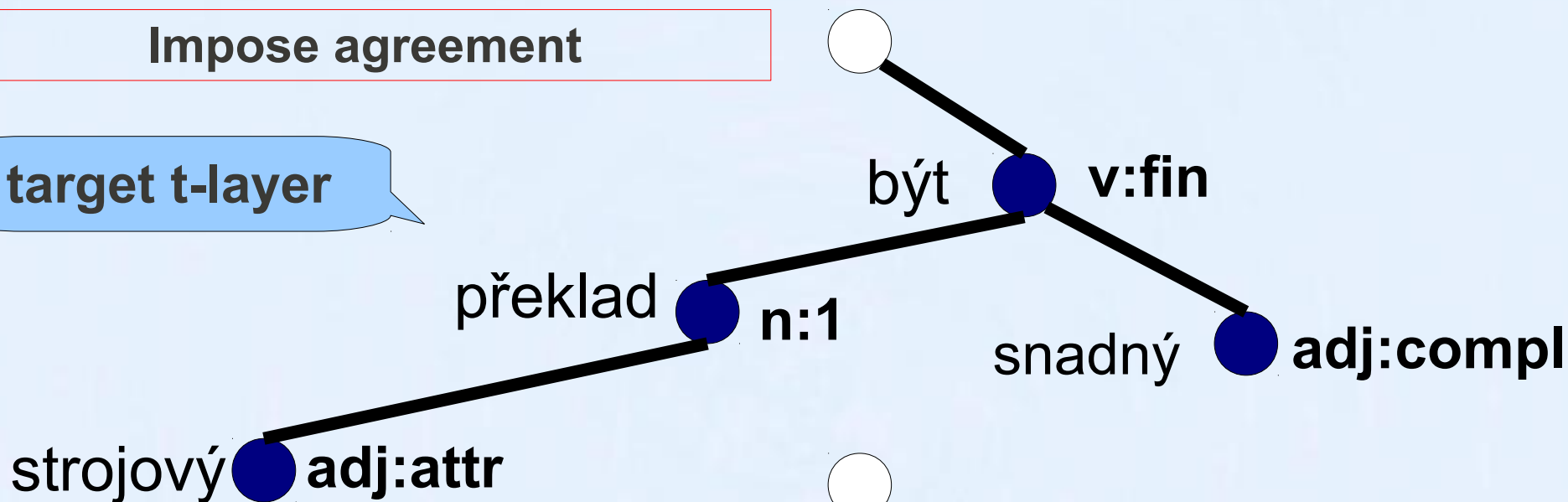
target a-layer



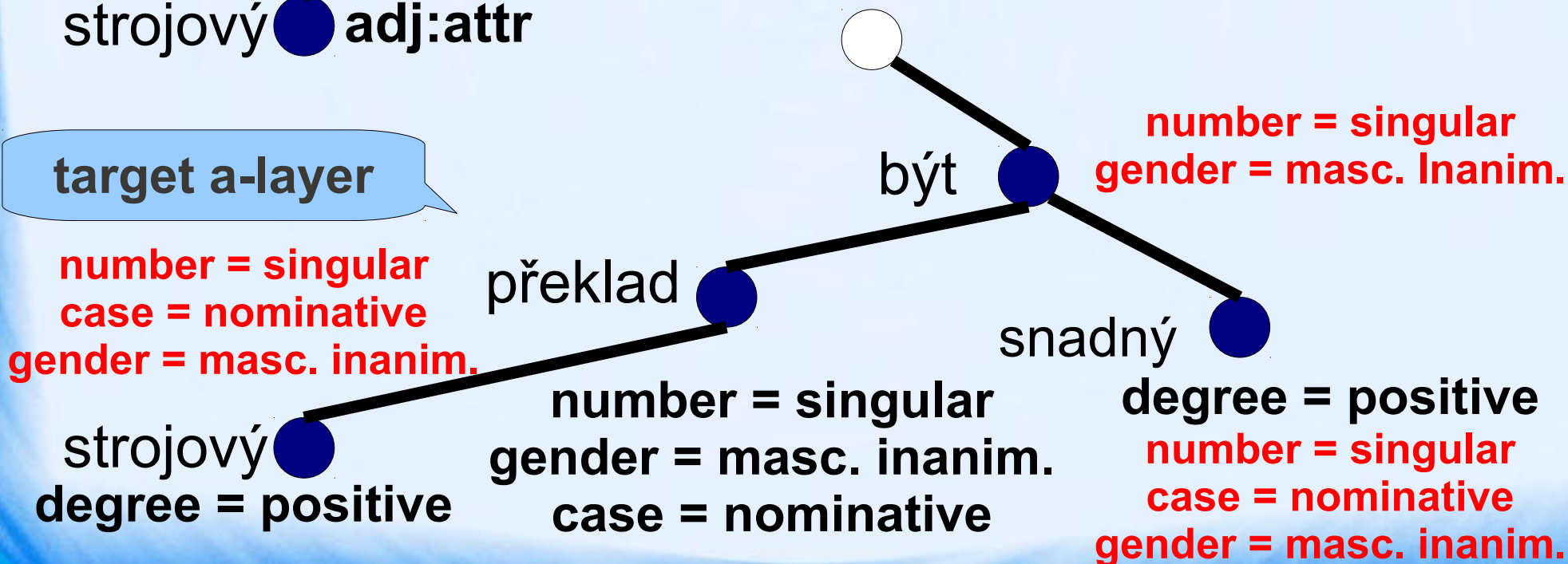
Demo Translation – Synthesis

Impose agreement

target t-layer



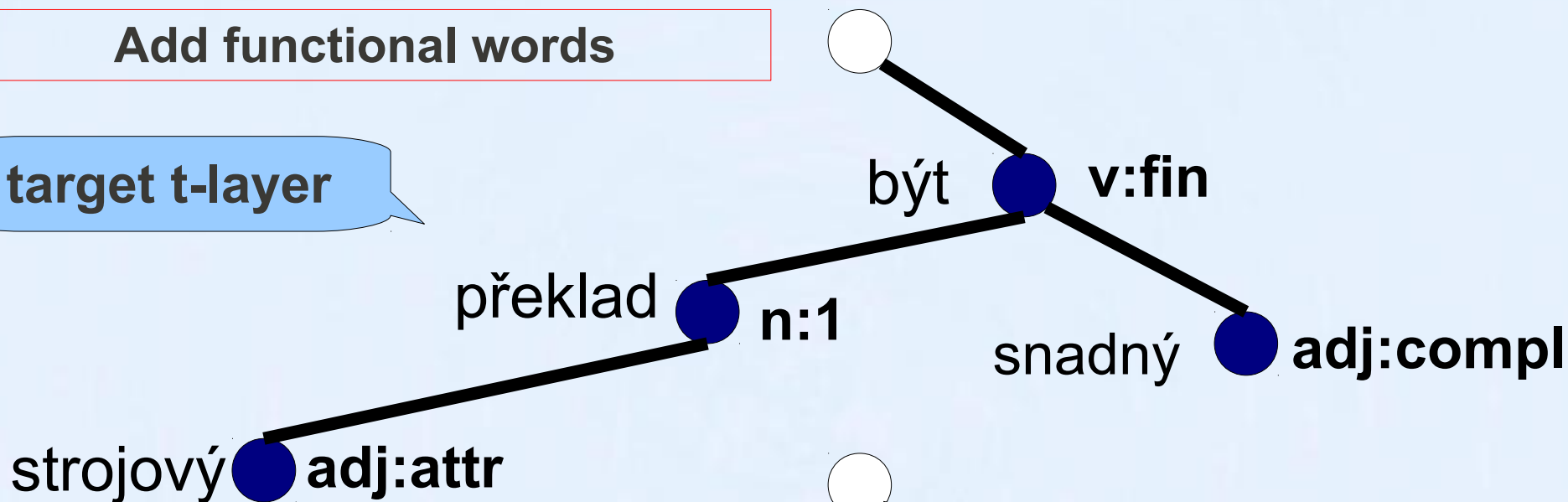
target a-layer



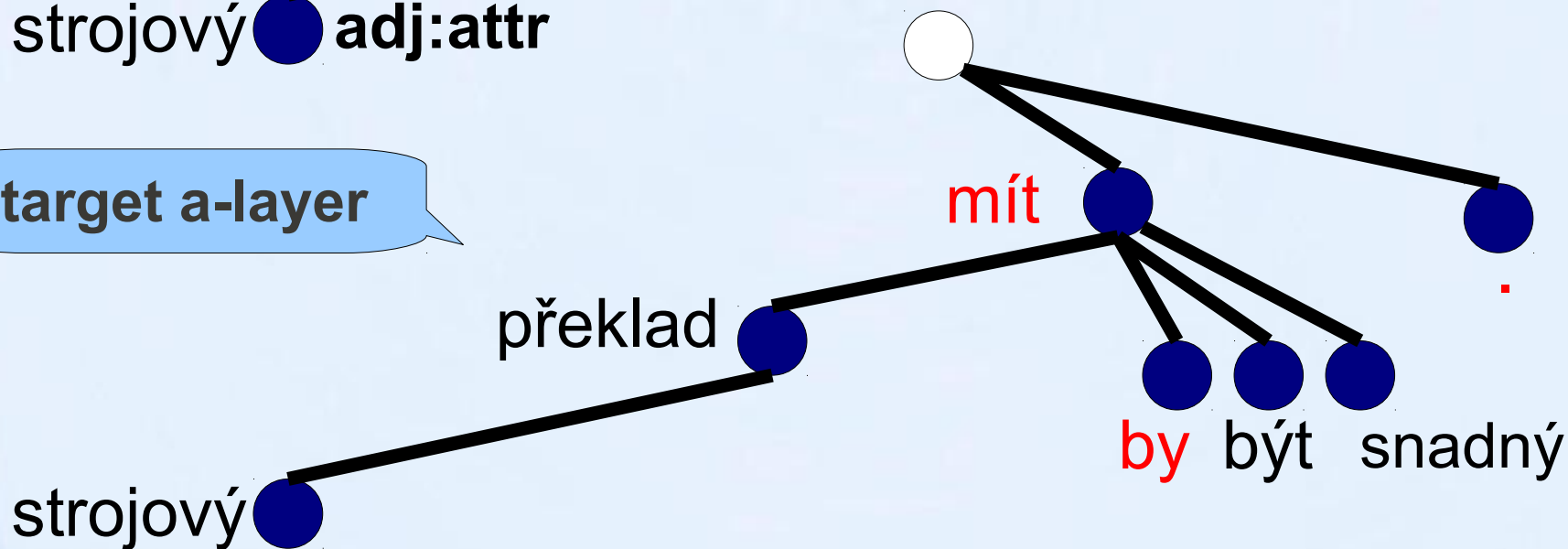
Demo Translation – Synthesis

Add functional words

target t-layer



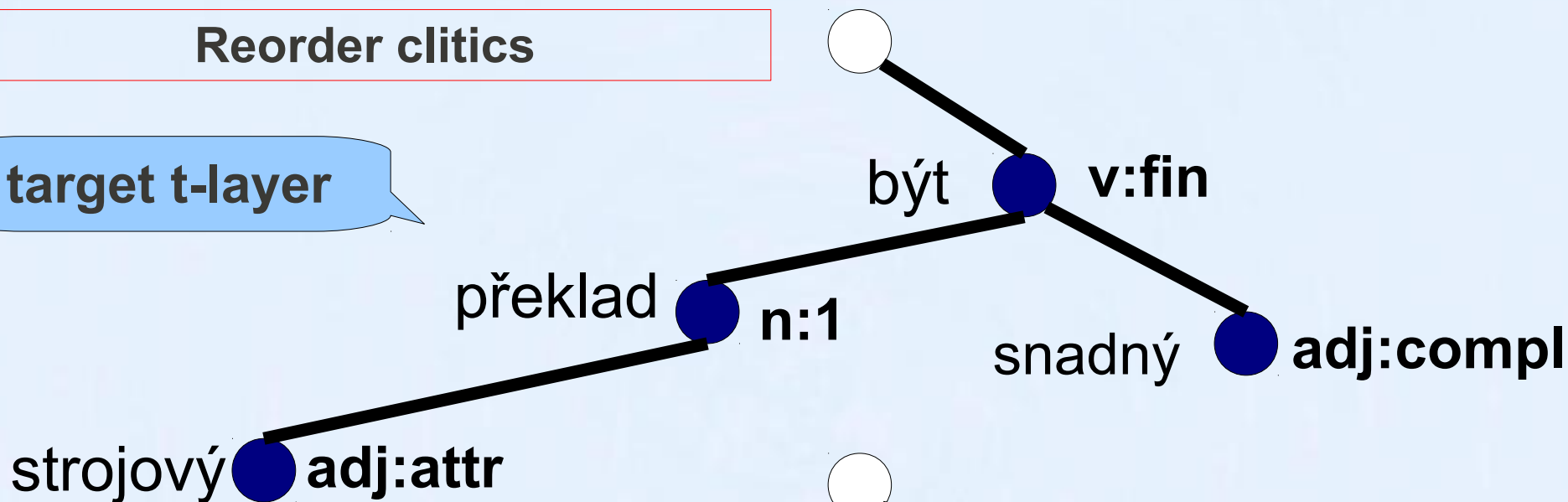
target a-layer



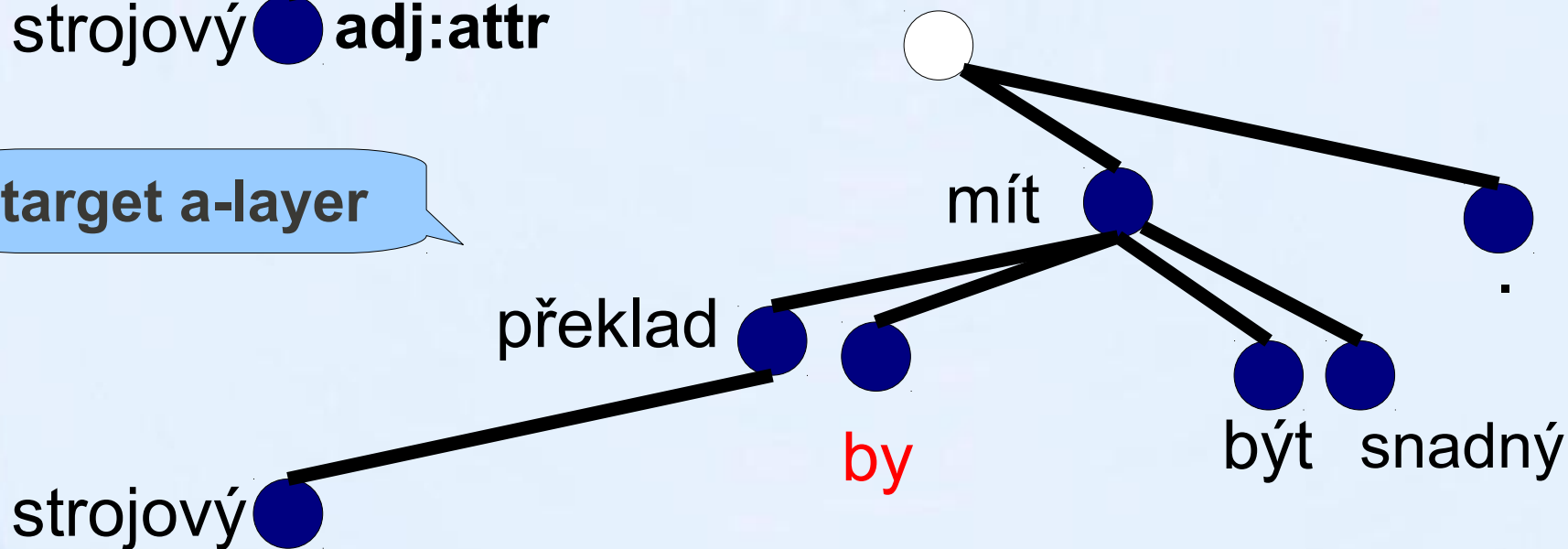
Demo Translation – Synthesis

Reorder clitics

target t-layer



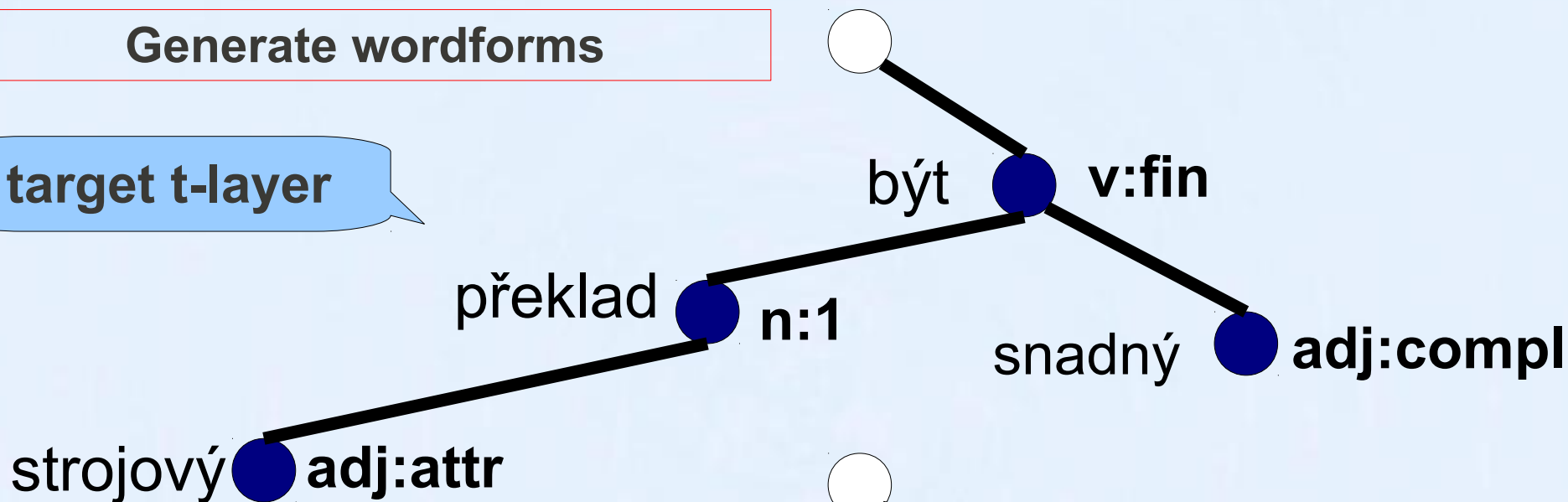
target a-layer



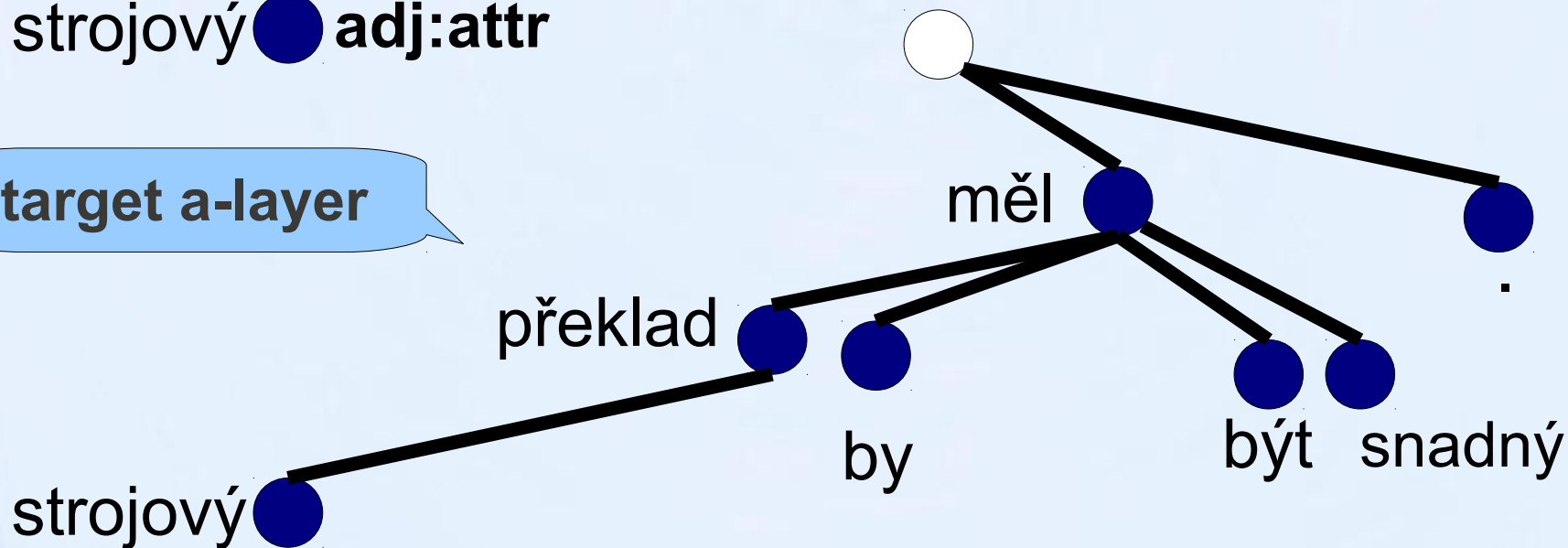
Demo Translation – Synthesis

Generate wordforms

target t-layer



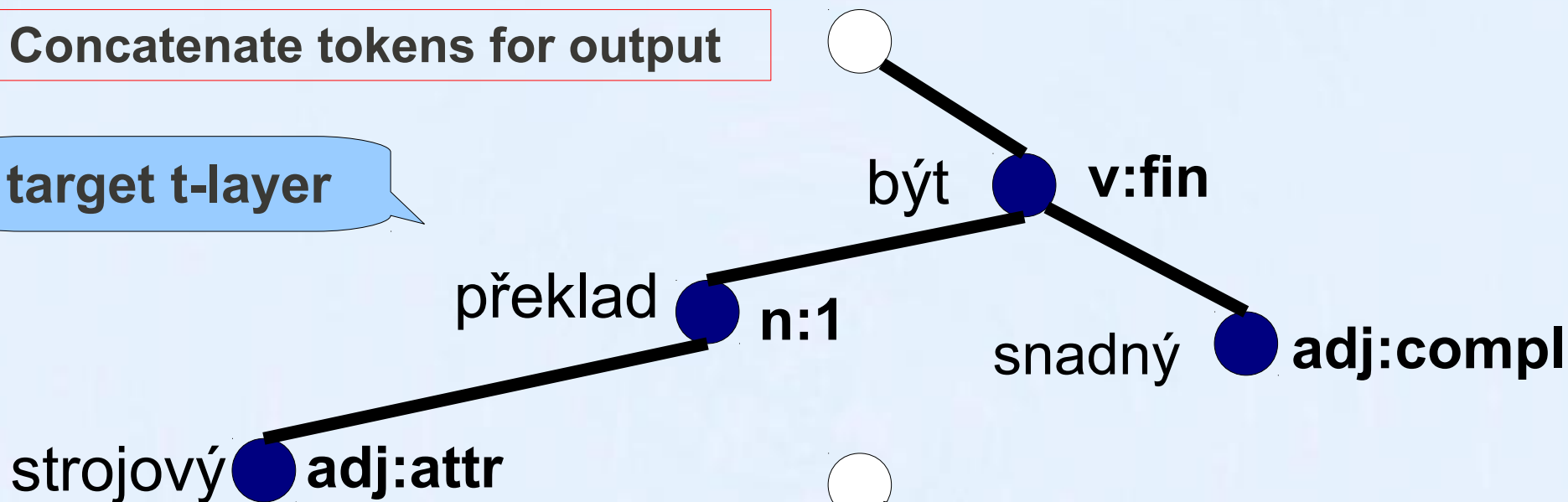
target a-layer



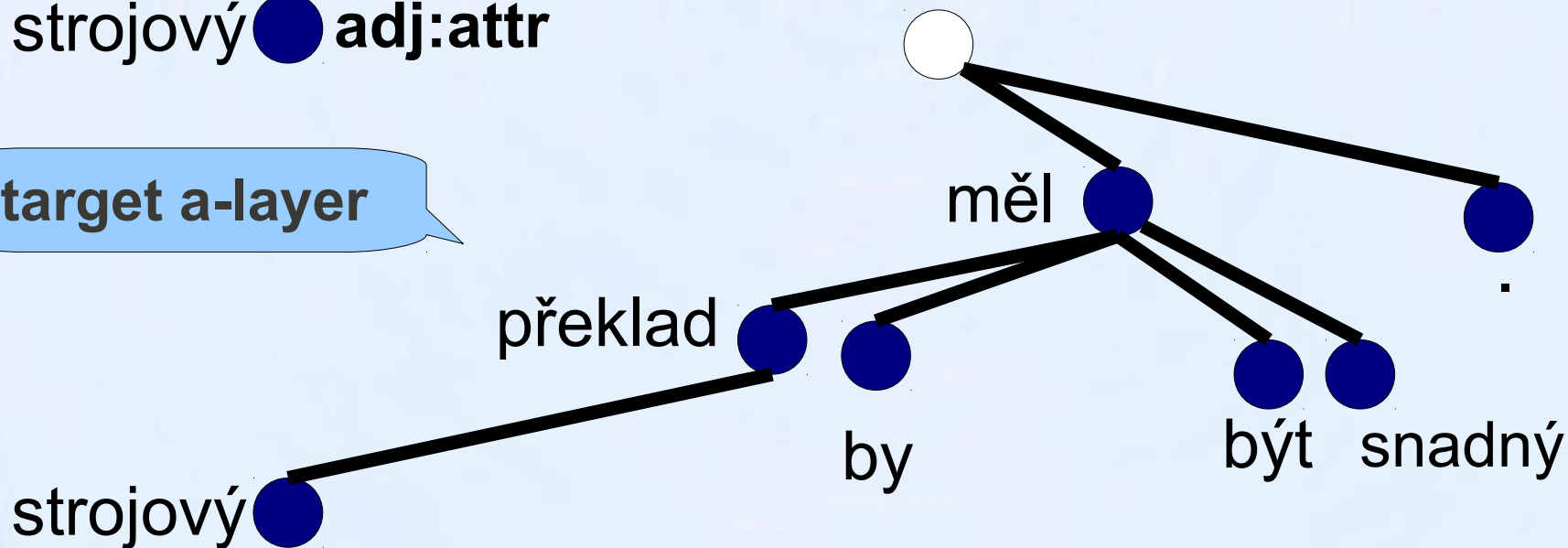
Demo Translation – Synthesis

Concatenate tokens for output

target t-layer



target a-layer



Strojový překlad by měl být snadný.

Demo Translation – Real 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_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

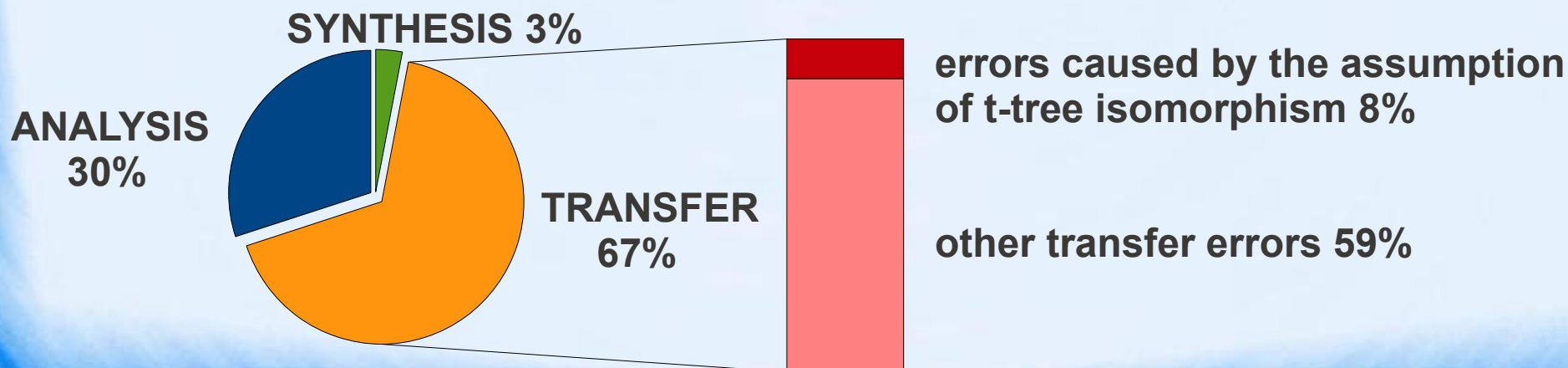
Ascii_quotes

Remove_repeated_tokens

Annotation of Translation Errors

sample of 250 sentences, 1463 errors in total

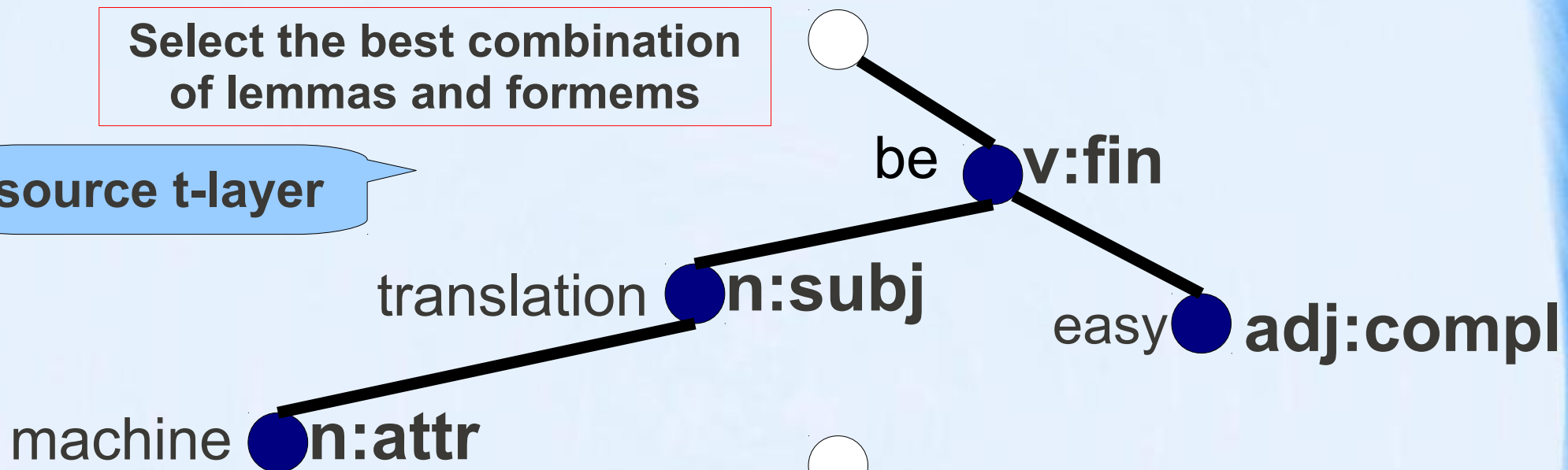
| | |
|----------------------|--|
| Type | lemma, formeme, gram., w. order,... |
| Subtype | gram: gender, person, tense,... |
| Seriousness | serious, minor |
| Circumstances | coordination, named entity, numbers |
| Source | tok, lem, tagger, parser, tecto, trans, x, syn, ? |



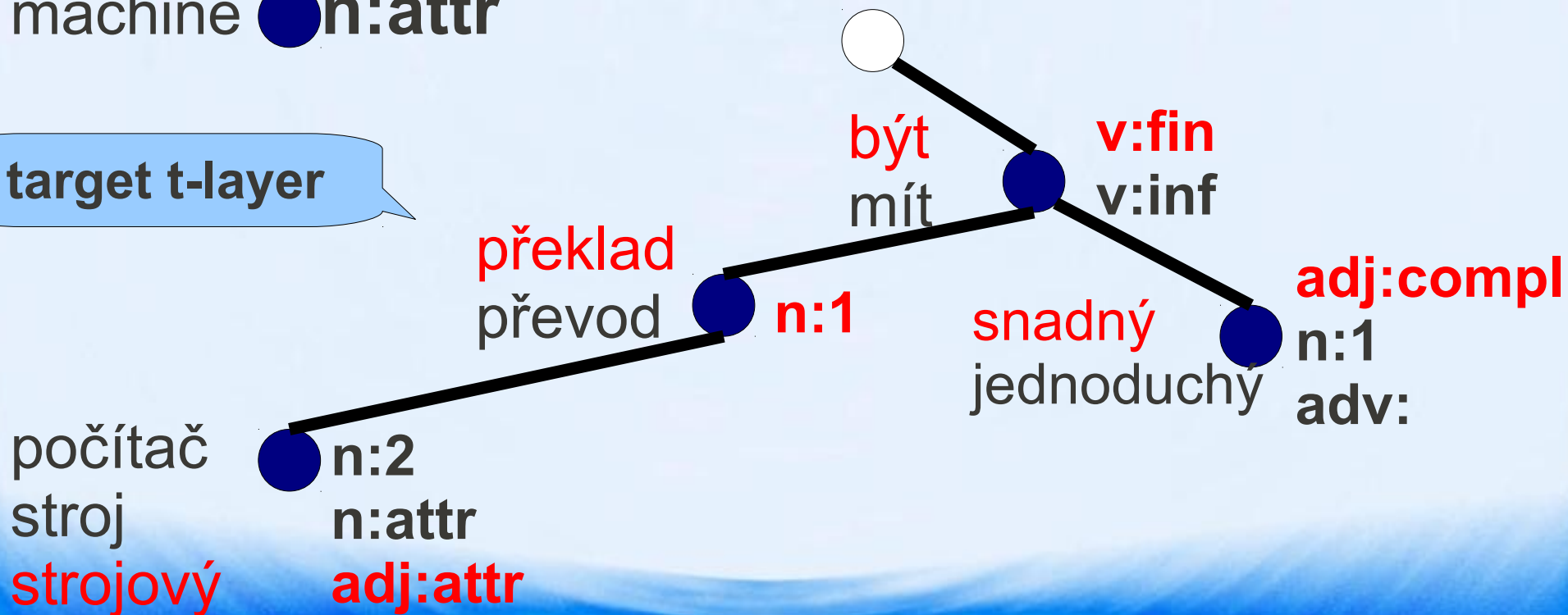
HMTM – Motivation

Select the best combination of lemmas and formems

source t-layer



target t-layer



HMTM – Motivation

Select the best label for each node

source t-layer

translation
n:subj

be v:fin

easy adj:compl

machine n:attr

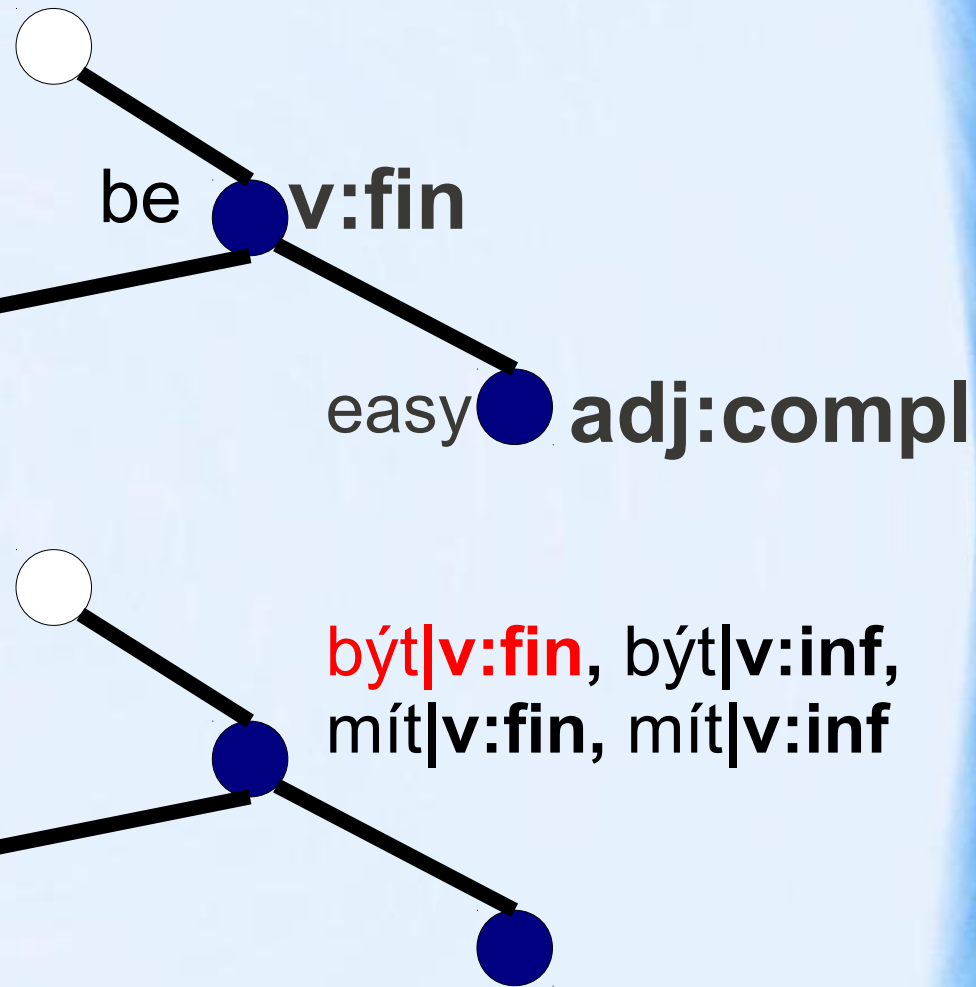
target t-layer

překlad|n:1,
převod|n:1

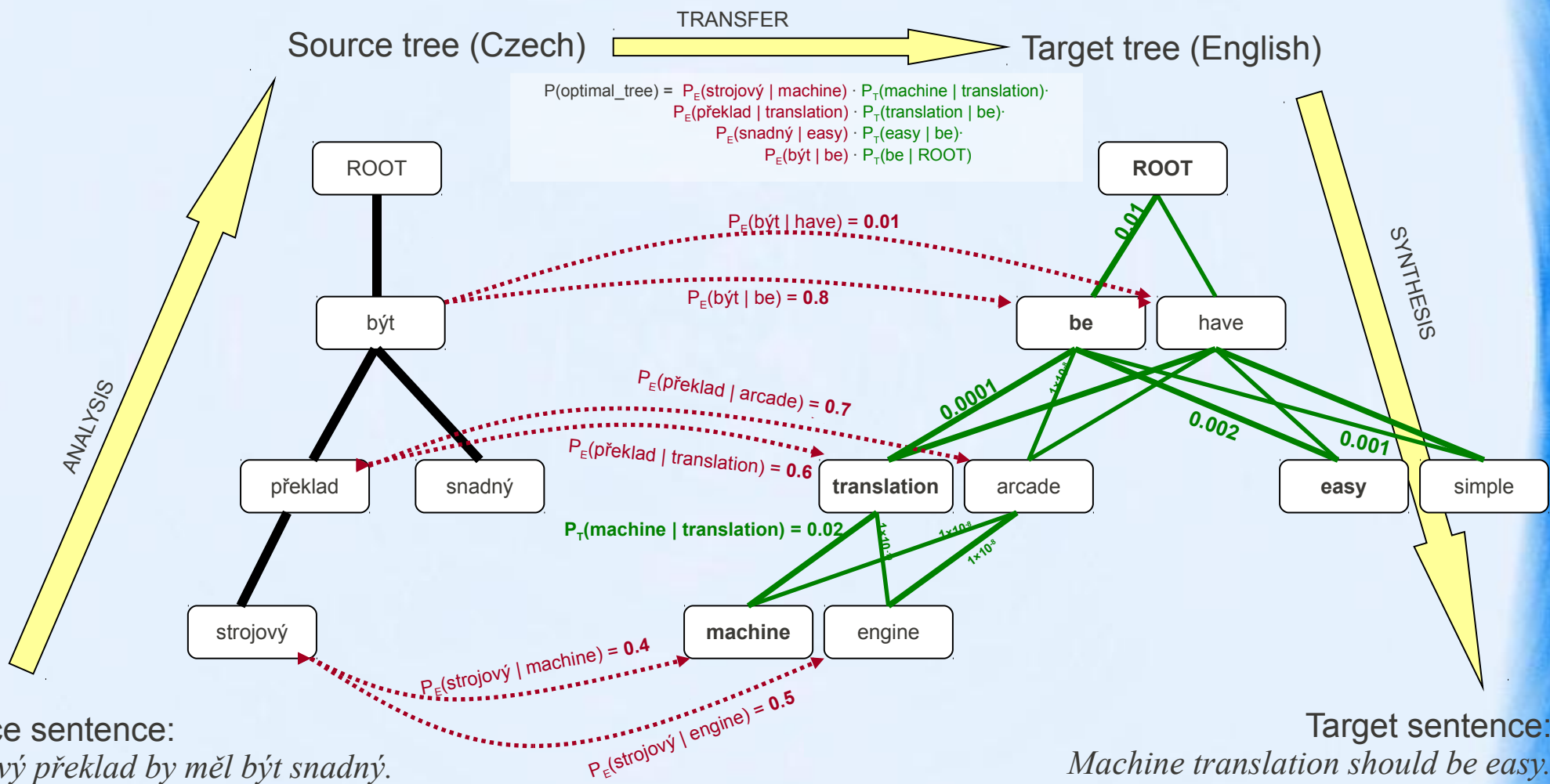
být|v:fin, být|v:inf,
mít|v:fin, mít|v:inf

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

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



HMTM in MT



$P_E(\text{source | target})$... emission probabilities ... **translation model**
 $P_T(\text{dependent | governing})$... transition probabilities ... **target-language tree model**

Combining Dictionaries

- new general interface (for lemmas and formems)
`$dict->get_translations($input_label, $features)`
returns a list of translation variants including probabilities
- OOP style, dictionary constructor can take another dictionary (or more) as a parameter → hierarchy

- Four basic types of dictionaries:

Static plain

loaded from a file „lemma → lemma“

Context

loaded from a file „lemma,features → lemma“

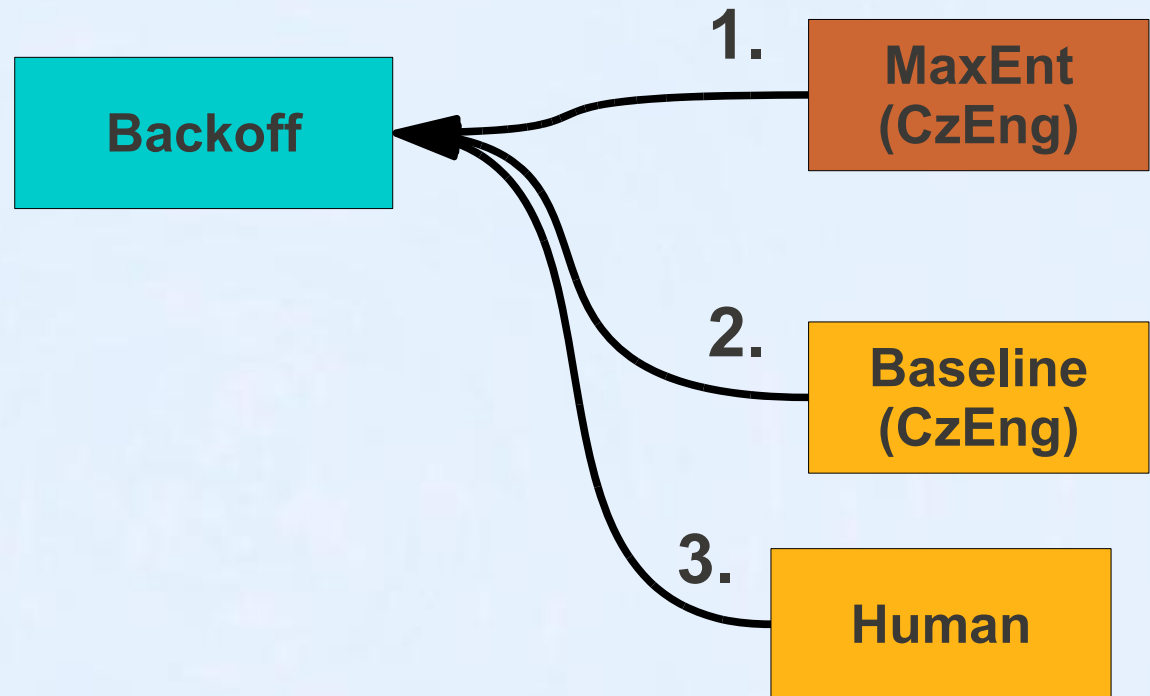
Derivational

translations derived dynamically, input dictionary

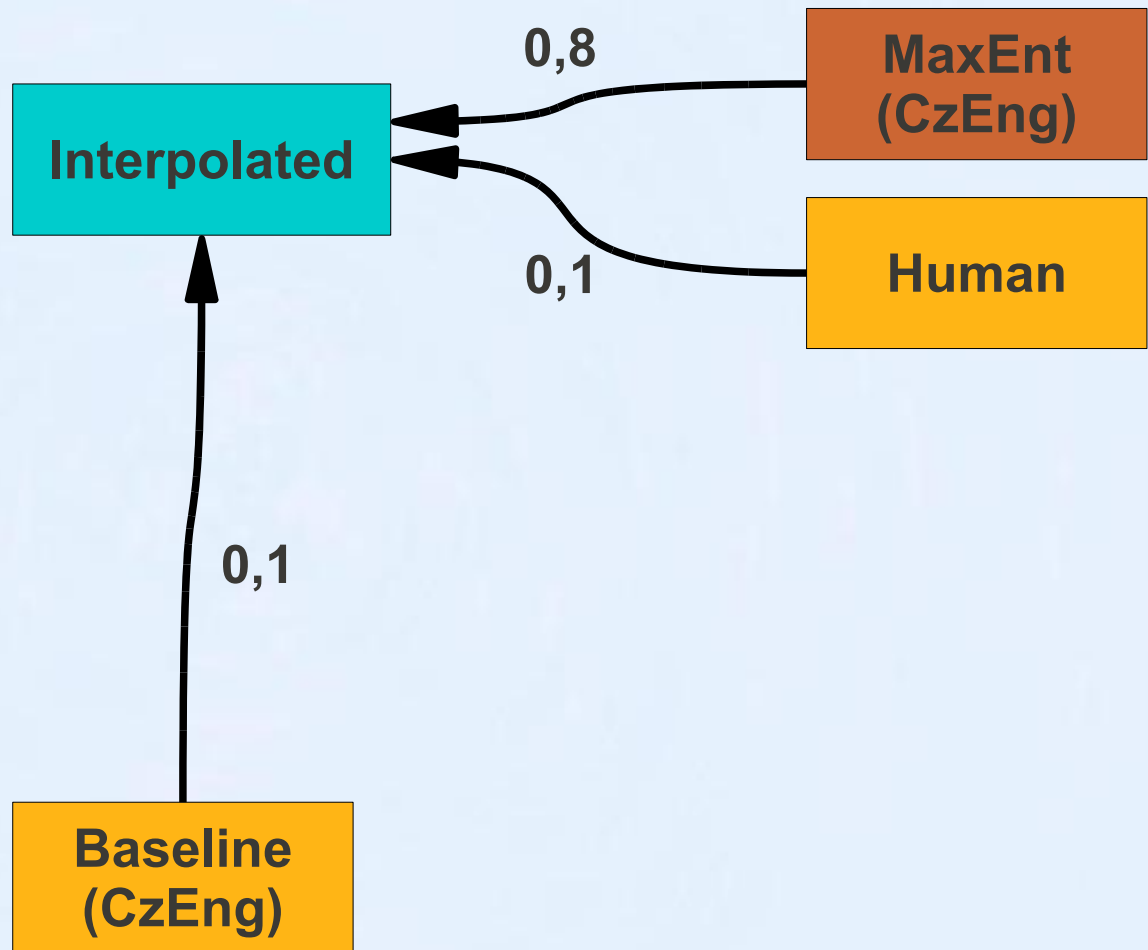
Combinational

combination of more input dictionaries

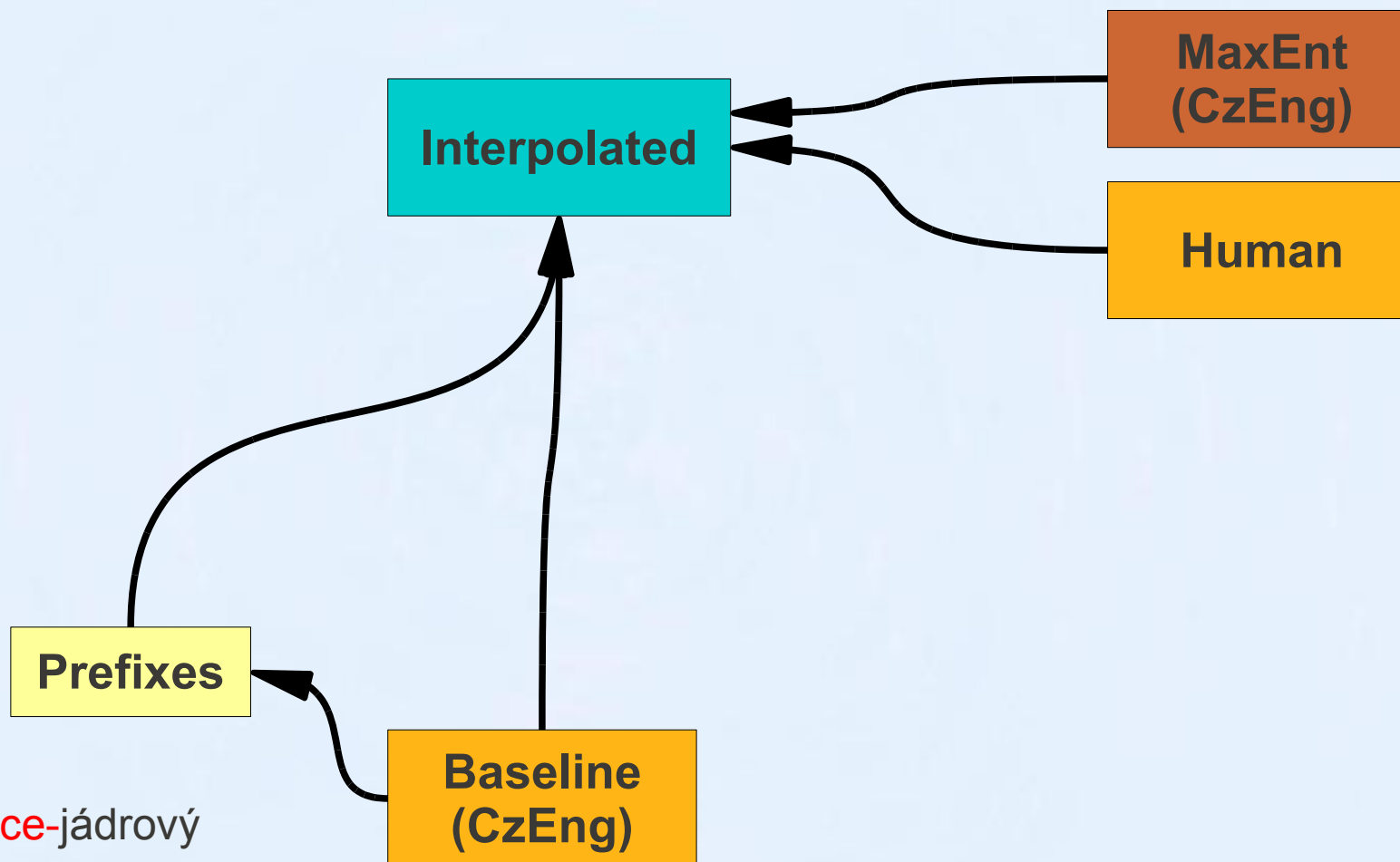
Hierarchy of lemma dictionaries



Hierarchy of lemma dictionaries

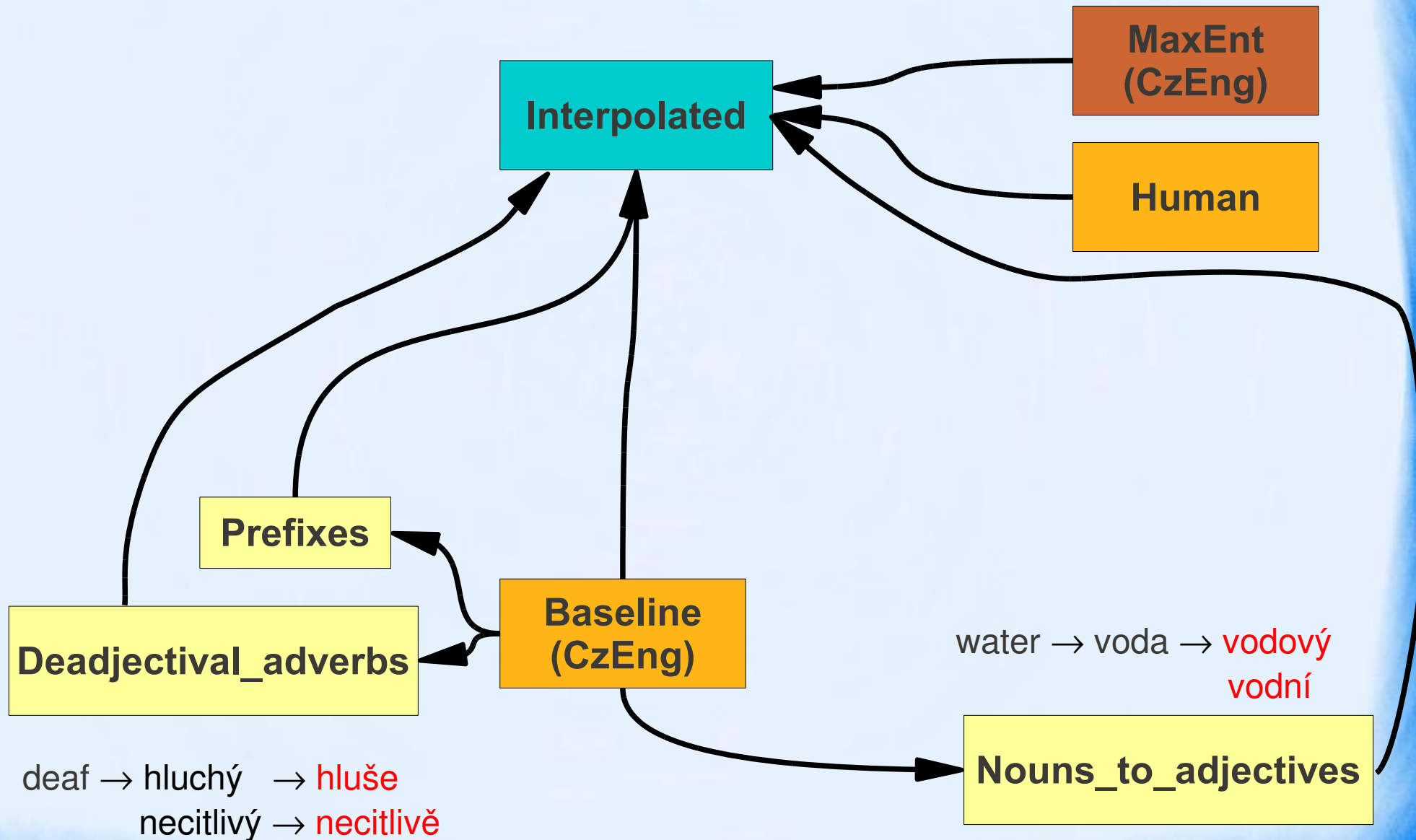


Hierarchy of lemma dictionaries

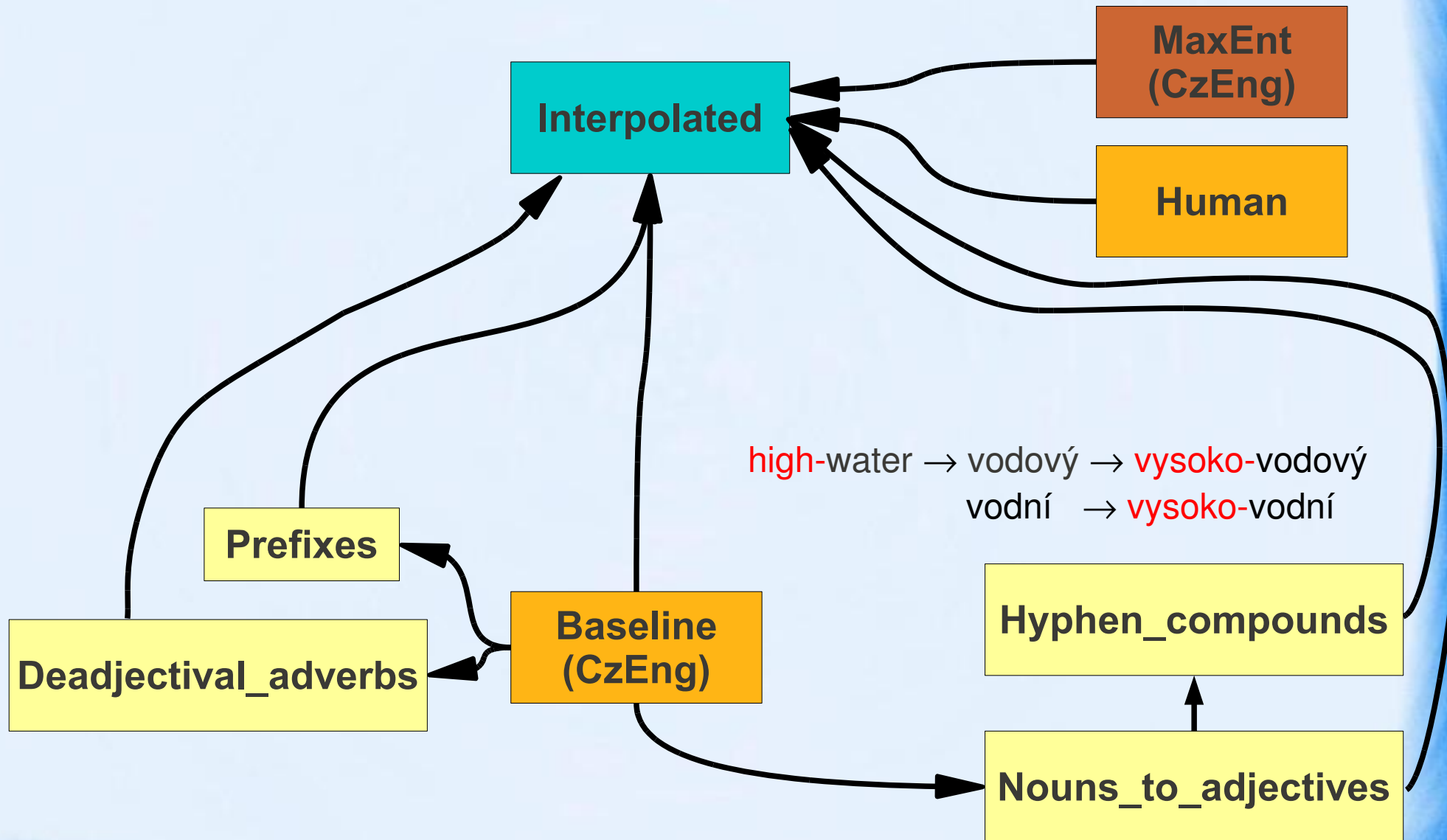


multi-core → více-jádrový
více-jádro
multi-jádrový
multi-jádro

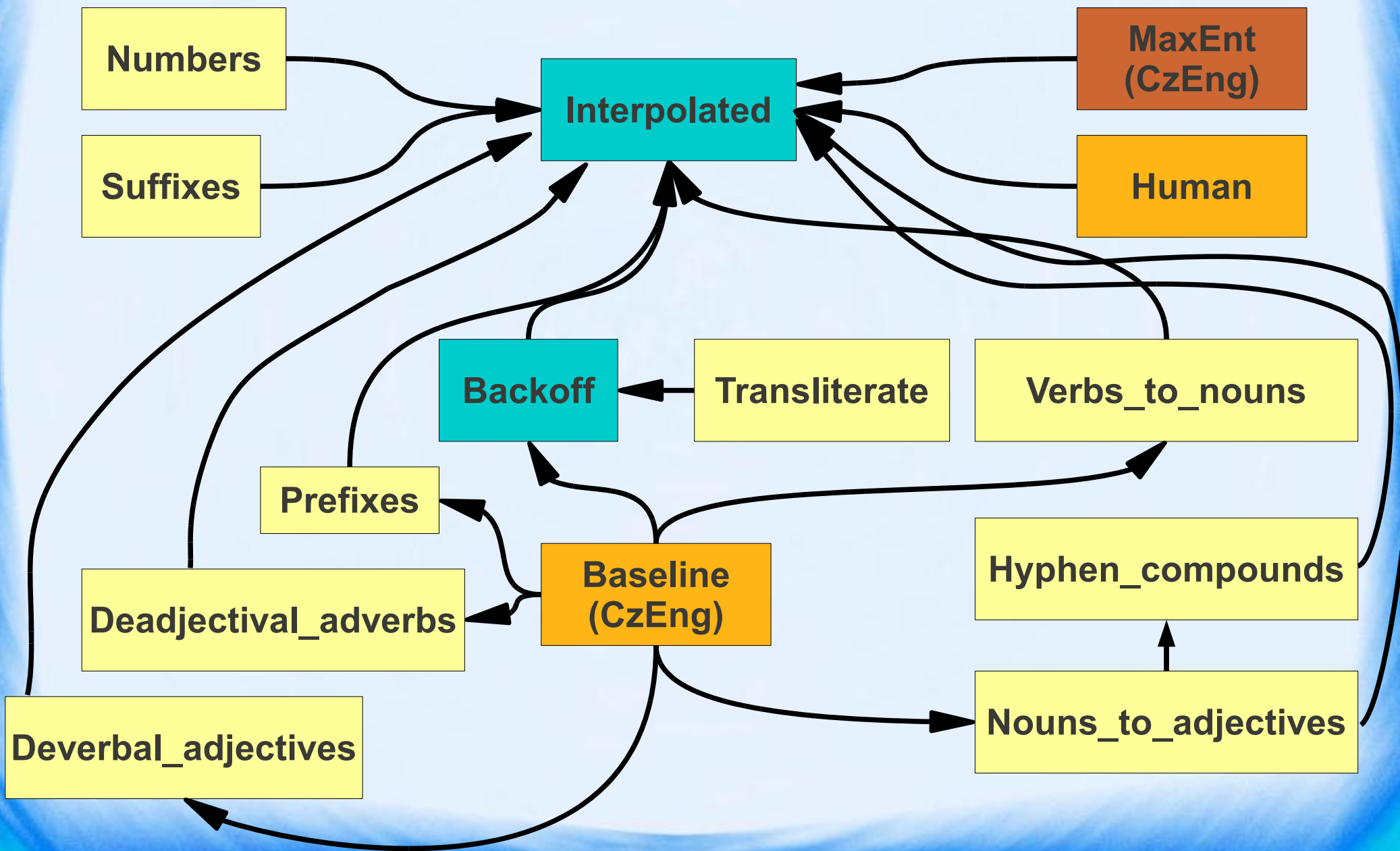
Hierarchy of lemma dictionaries



Hierarchy of lemma dictionaries



Hierarchy of lemma dictionaries



Maximum Entropy Dictionary

Baseline Dictionary

$$p(y|x) = \frac{\text{count}(x, y)}{\text{count}(x)}$$

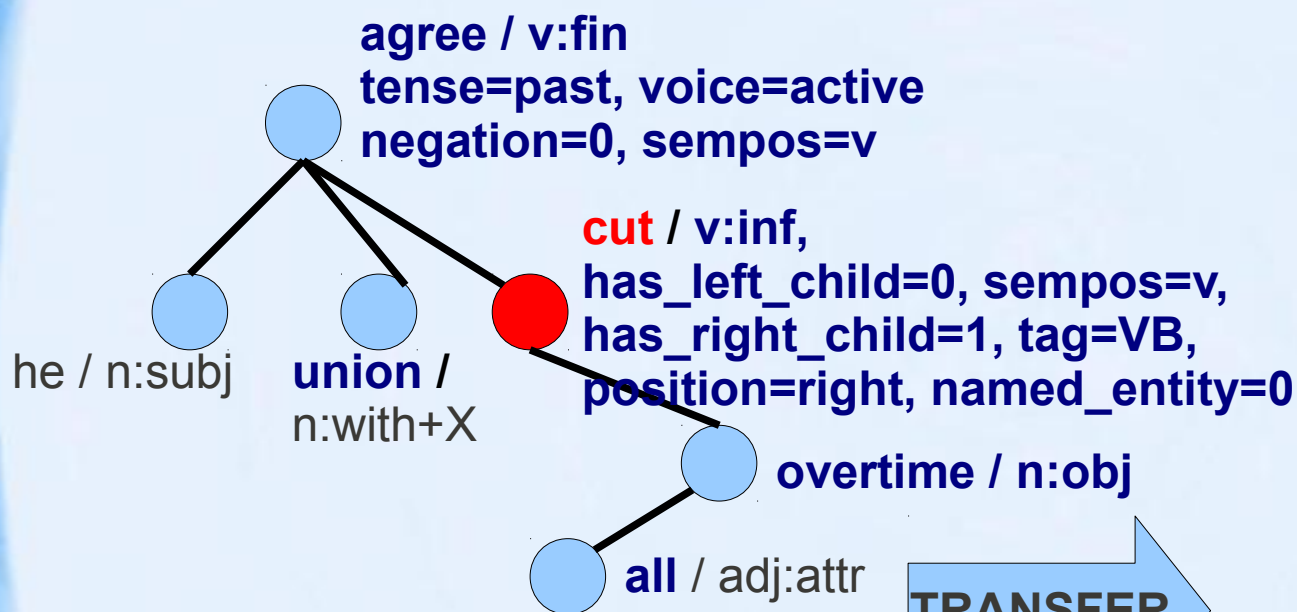
- Maximum likelihood estimates
(from the training sections of CzEng 0.9)
- Pruned by thresholds on $p(x|y)$ and $p(y|x)$
- No context used
x = source lemma
y = target lemma

MaxEnt Dictionary

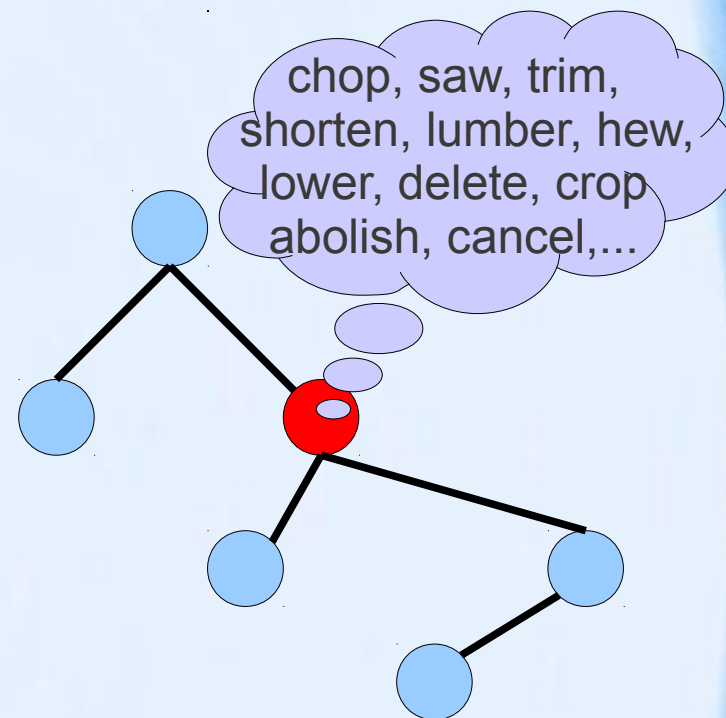
$$p(y|x) = \frac{1}{Z(x)} \exp \sum_i \lambda_i f_i(x, y)$$

- One MaxEnt model for each source lemma
(same training data as for the Baseline Dict.)
- Interpolated with Baseline Dict. (due to pruning)
- Context features used (x = source context)
 - local tree context
 - local linear context
 - morphological & syntactic categories
 - ...

Maximum Entropy Dictionary



TRANSFER



ANALYSIS

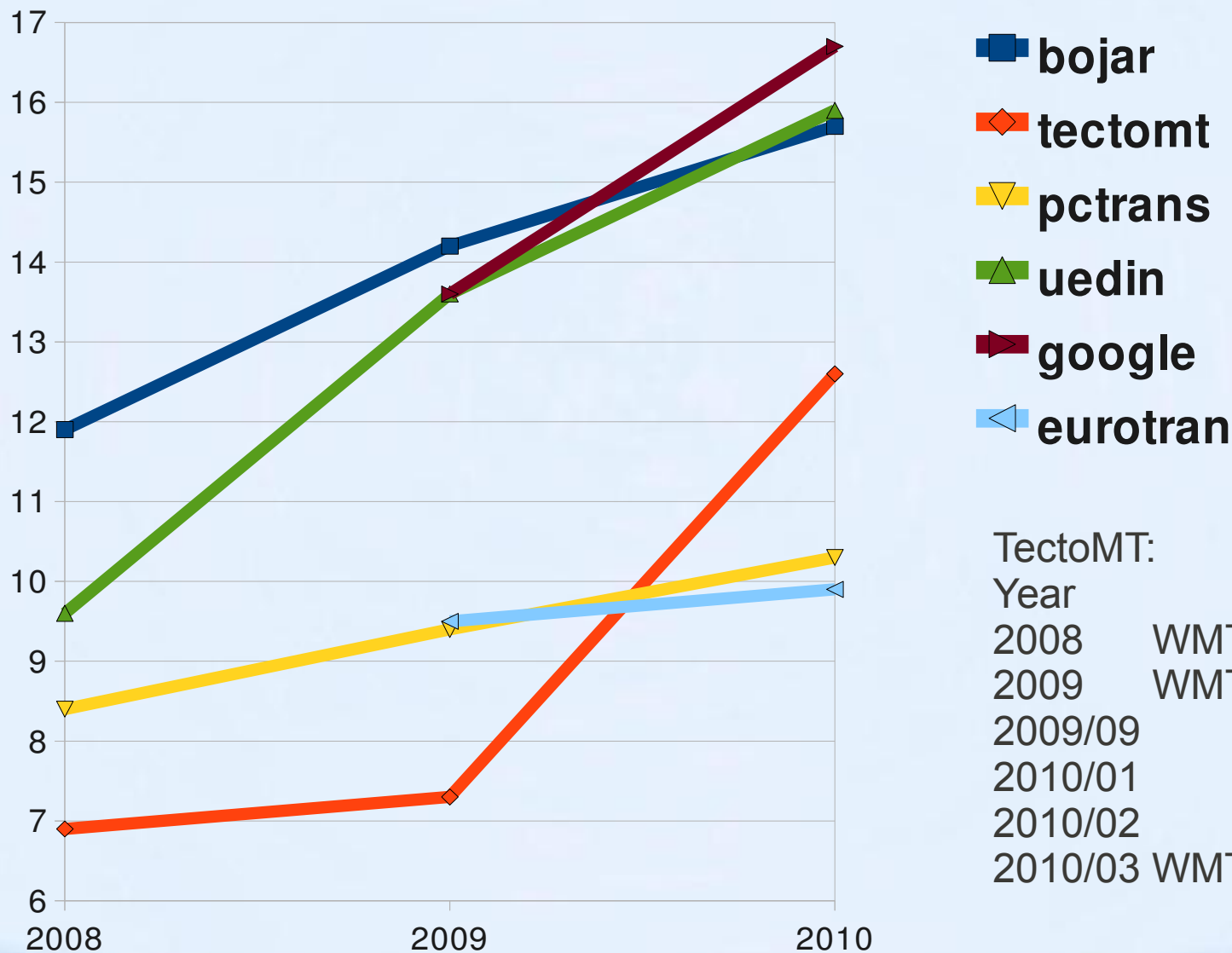
SYNTHESIS

He agreed with the unions to cut all overtime.

Dohodl se s odbory na zrušení všech přesčasů.

Results – BLEU

WMT = Workshop on Statistical Machine Translation

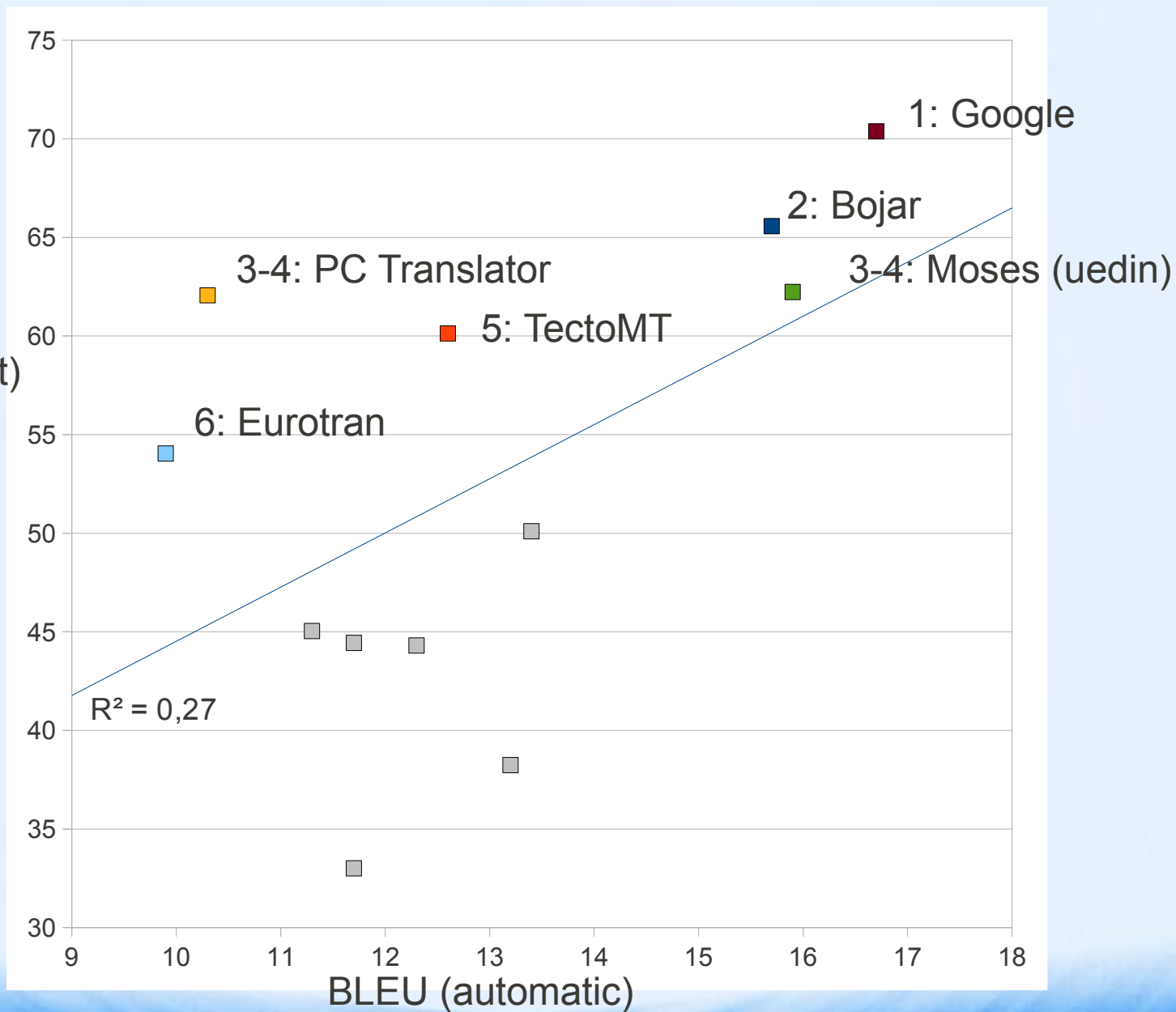


TectoMT:

| Year | WMT | 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

Rank
(human judgement)



Examples of Translation (2009)

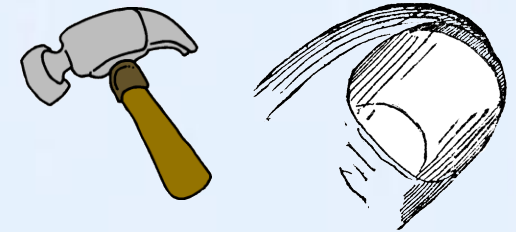
A miss by an inch
is a miss by a mile.

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.

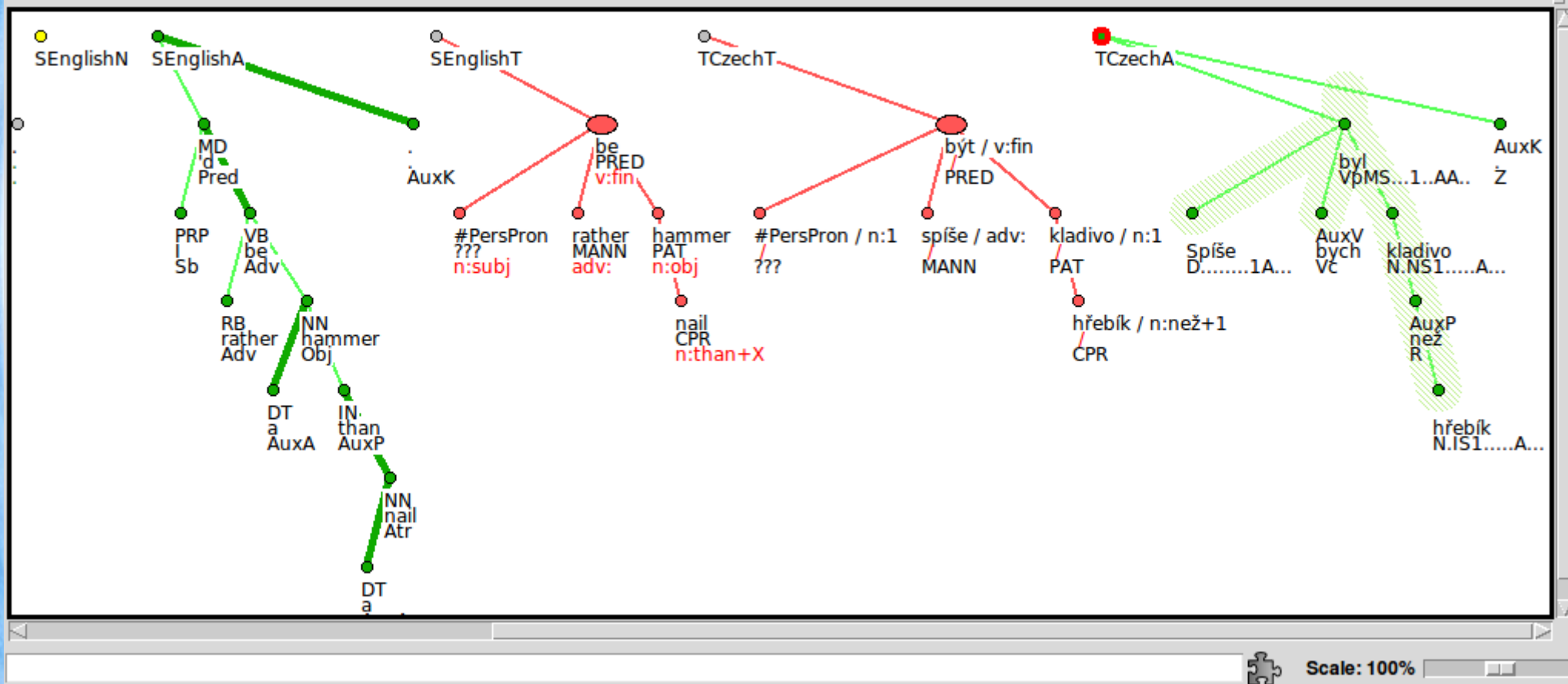
Pták v ruce je cenný
dvakrát v Bushovi.



Example of Translation (2010)

I'd rather be a hammer than a nail.

Spíše bych byl kladivo než hřebík.



Scale: 100%

Sample of MaxEnt Features

input_label=nail

output_label=hřebík#N (metal nail)

| | |
|--------------------------------|--------------------------|
| child_formeme_n:in+X=1 | 1.64483855116042 |
| is_member=1 | 1.30042900630692 |
| child_formeme_v:fin=1 | 1.04422203176176 |
| next_node_tlemma=down | 0.838961007712912 |
| is_capitalized=1 | 0.792130821958927 |
| position=right | 0.747785245407306 |
| tense_g=post | 0.744919903760696 |
| voice_g=active | 0.659489975893991 |
| prev_node_tlemma=drive | 0.655357850937254 |
| parent_capitalized=1 | 0.622953832124697 |
| formeme=n:from+X | 0.599348506643414 |
| prev_node_tlemma=hammer | 0.592276691427986 |
| child_tlemma_few=1 | 0.553464629114697 |
| child_tlemma_remove=1 | 0.546698831608057 |
| sempos=n.denot | 0.504719359514573 |
| next_node_tlemma=and | 0.502529618088752 |
| formeme_g=v:until+fin | 0.491064112122981 |
| child_tlemma_rusty=1 | 0.428884558837039 |
| tag_g=VBP | 0.422967377093101 |
| next_node_tlemma=screw | 0.344701934524519 |
| ... | |

output_label=nehet#N (fingernail or toenail)

| | |
|------------------------|--------------------------|
| child_formeme_n:poss=1 | 1.32717038827268 |
| child_tlemma_finger=1 | 1.07509772743853 |
| child_formeme_n:of+X=1 | 0.982021327950337 |
| position=left | 0.886912864256063 |
| prev_node_tlemma=black | 0.770671304450658 |
| child_tlemma_broken=1 | 0.761077744287099 |
| child_formeme_v:attr=1 | 0.700099311992958 |
| formeme=n:at+X | 0.674547829214778 |
| formeme_g=n:attr | 0.673367412957367 |
| child_tlemma_long=1 | 0.673158400394094 |
| next_node_tlemma=file | 0.600496248030202 |
| child_tlemma_false=1 | 0.584236638145312 |
| prev_node_tlemma=false | 0.584236638145312 |
| number=sg | 0.563056142428995 |
| formeme=n:obj | 0.533943098032196 |
| formeme=n:by+X | 0.528852315800188 |
| ... | |