

TectoMT: Machine Translation System

Martin Popel

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



FEAST (Forum Entwicklung und Anwendung von Sprach-Technologien)
November 16, 2010, Saarbrücken, Germany

Outline

- Treex (framework) vs. TectoMT (MT system)
- Demo translation step by step
- Annotation of translation errors
- Improvements
 - Hidden Markov Tree Models (HMTM)
 - Combining dictionaries
 - Maximum Entropy dictionary
- Results
- Examples of translation

Treex vs. TectoMT



2005

...

2011

NLP framework
TectoMT

multi-purpose
NLP framework
Treex

MT system
TectoMT

lemmatization

tagging

parsing

MT system
TectoMT

lemmatization

tagging

parsing

Main author:
Zdeněk Žabokrtský

coreference

PEDT preprocessing

CzEng analysis

treebank conversions

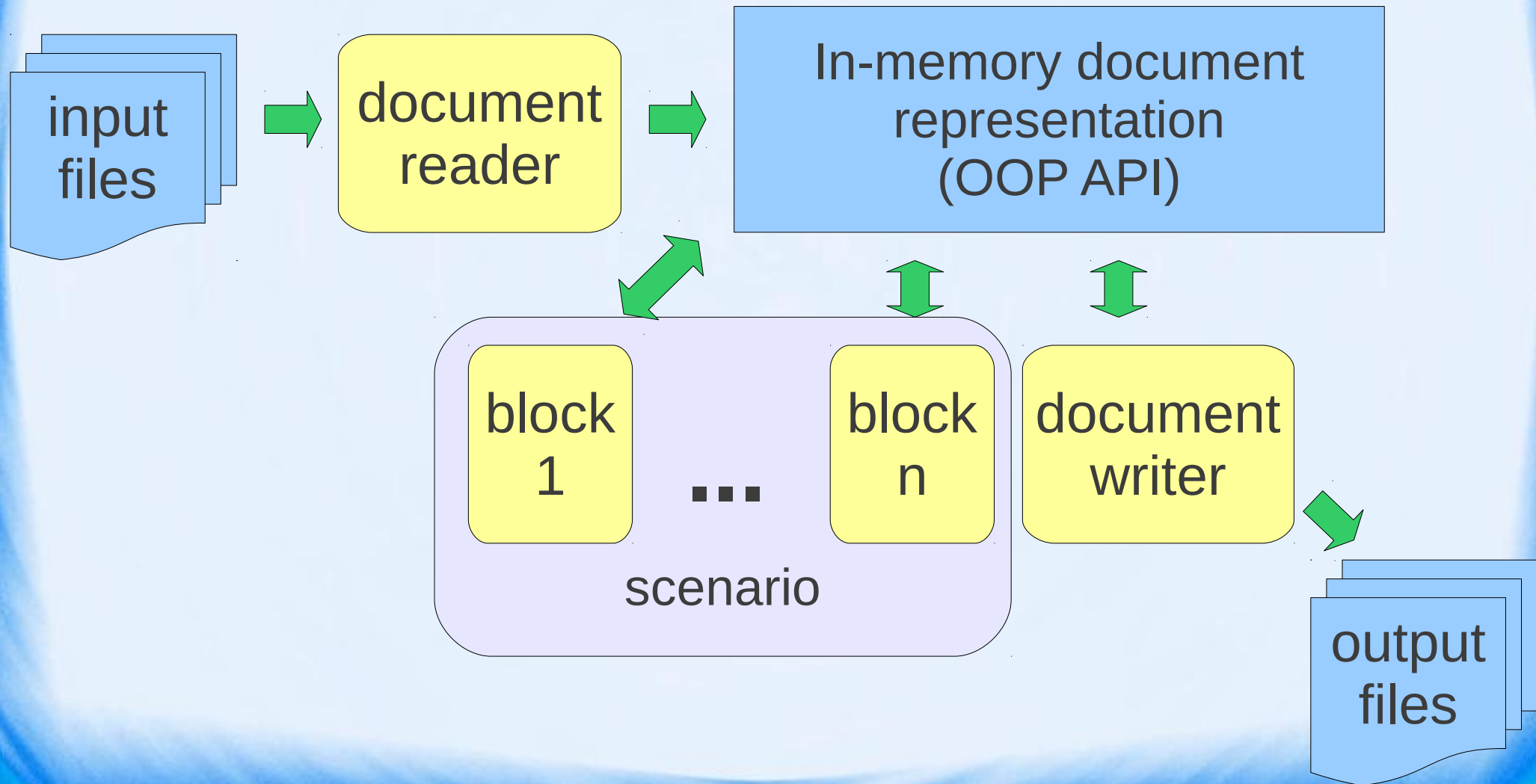
named entity r.

alignment (word,tree)

SMT preproc.

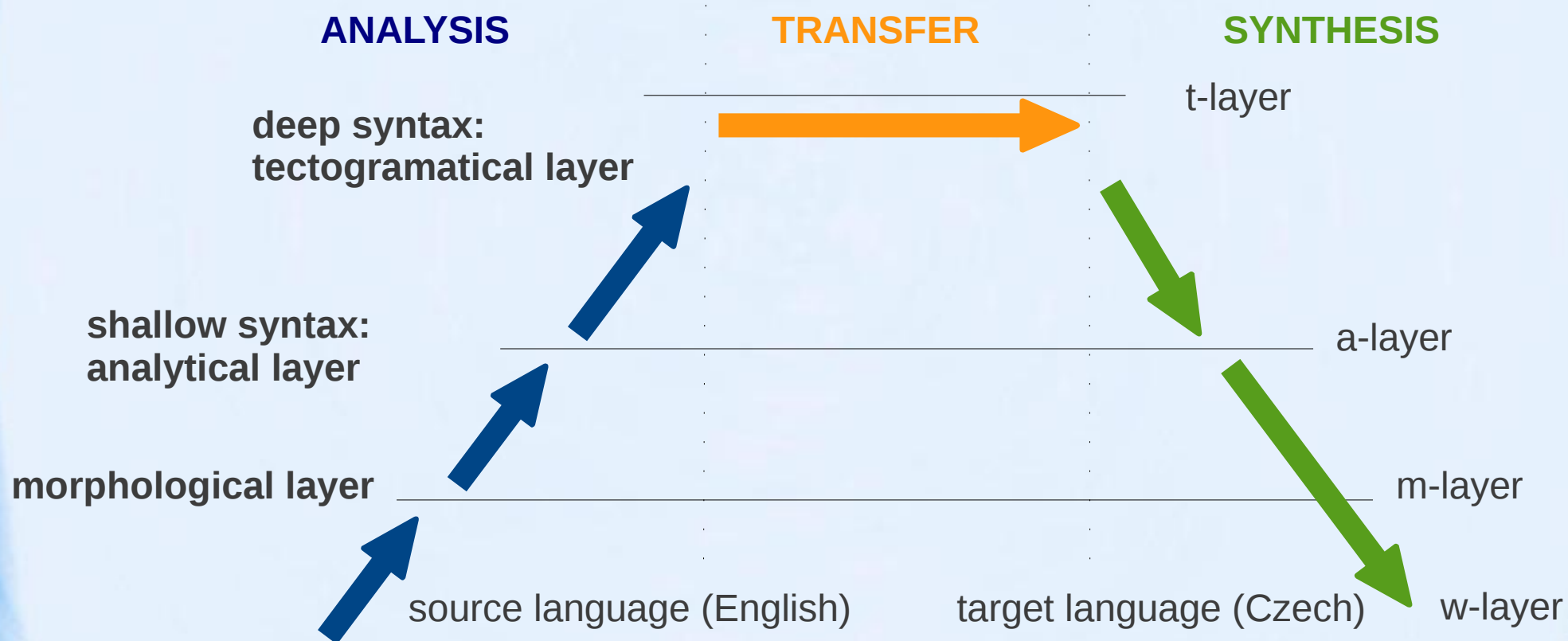
etc.

Treex architecture



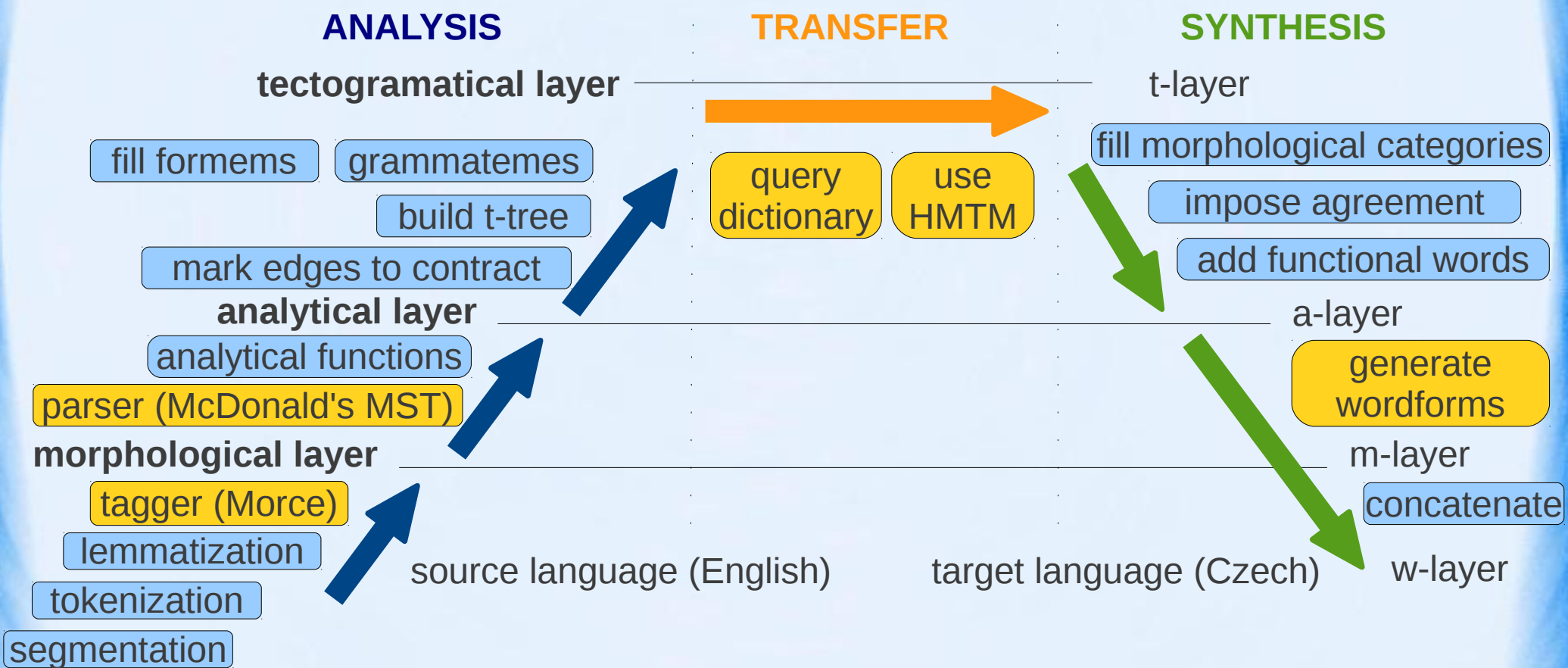
Translation scheme

transfer over the tectogrammatical layer



Translation scheme

rule based & statistical blocks



Demo Translation – Analysis

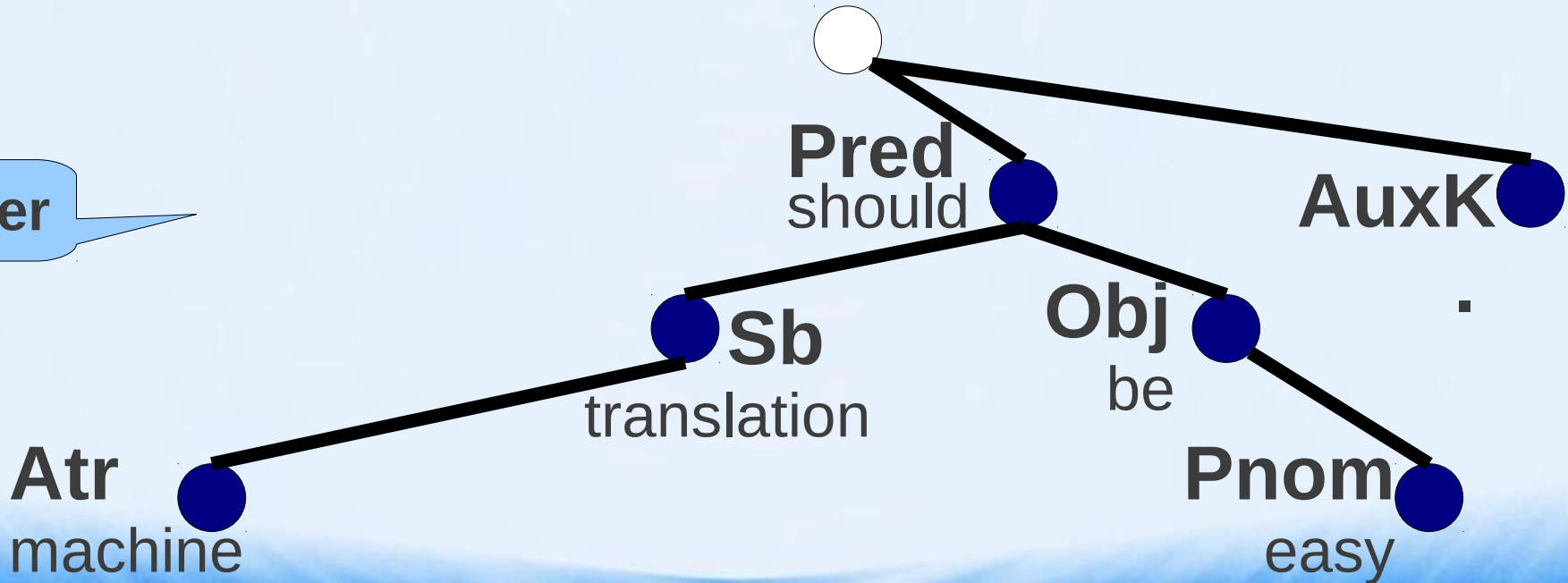
raw text

Machine translation should be easy.

m-layer

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

a-layer



Demo Translation – Analysis

raw text

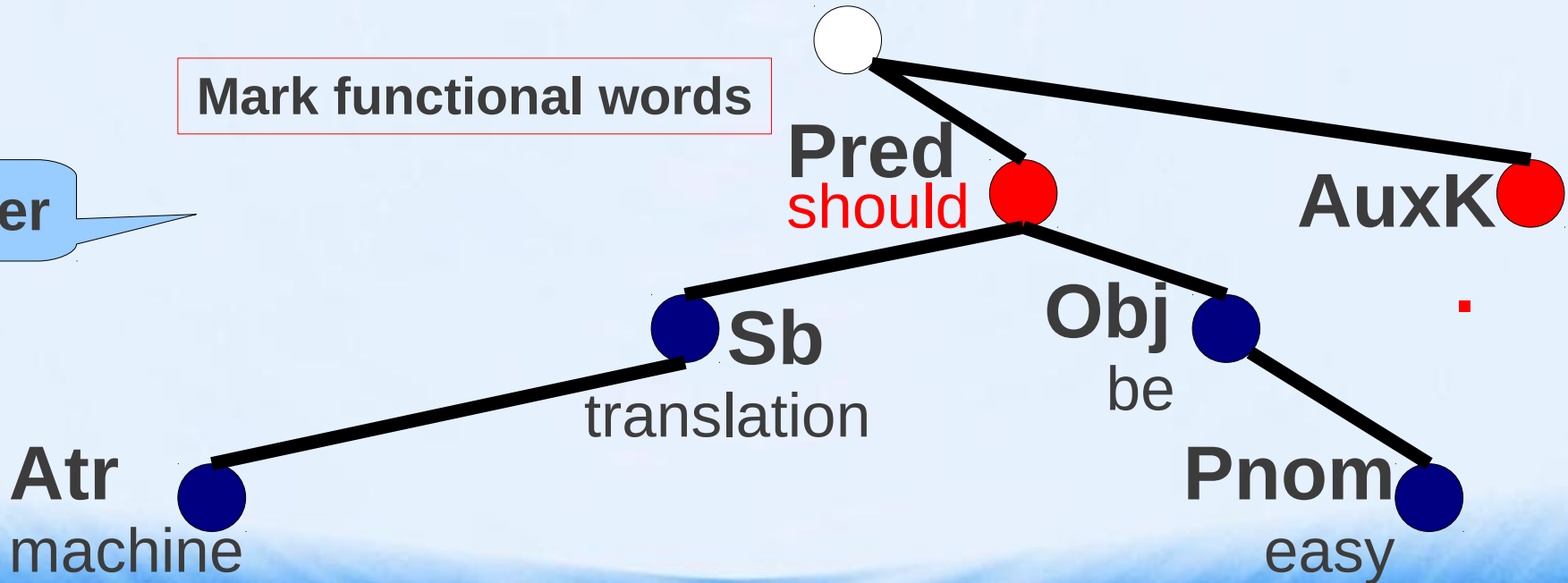
Machine translation should be easy.

m-layer

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

Mark functional words

a-layer



Demo Translation – Analysis

raw text

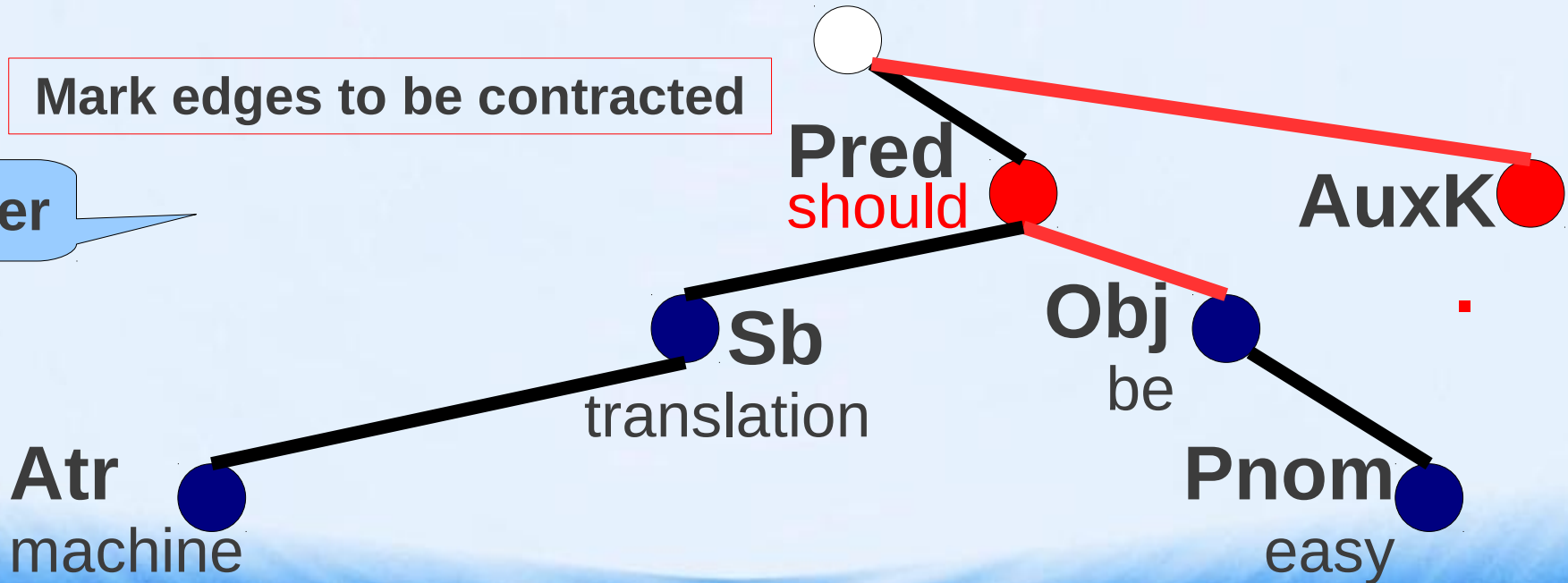
Machine translation should be easy.

m-layer

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

Mark edges to be contracted

a-layer



Demo Translation – Analysis

raw text

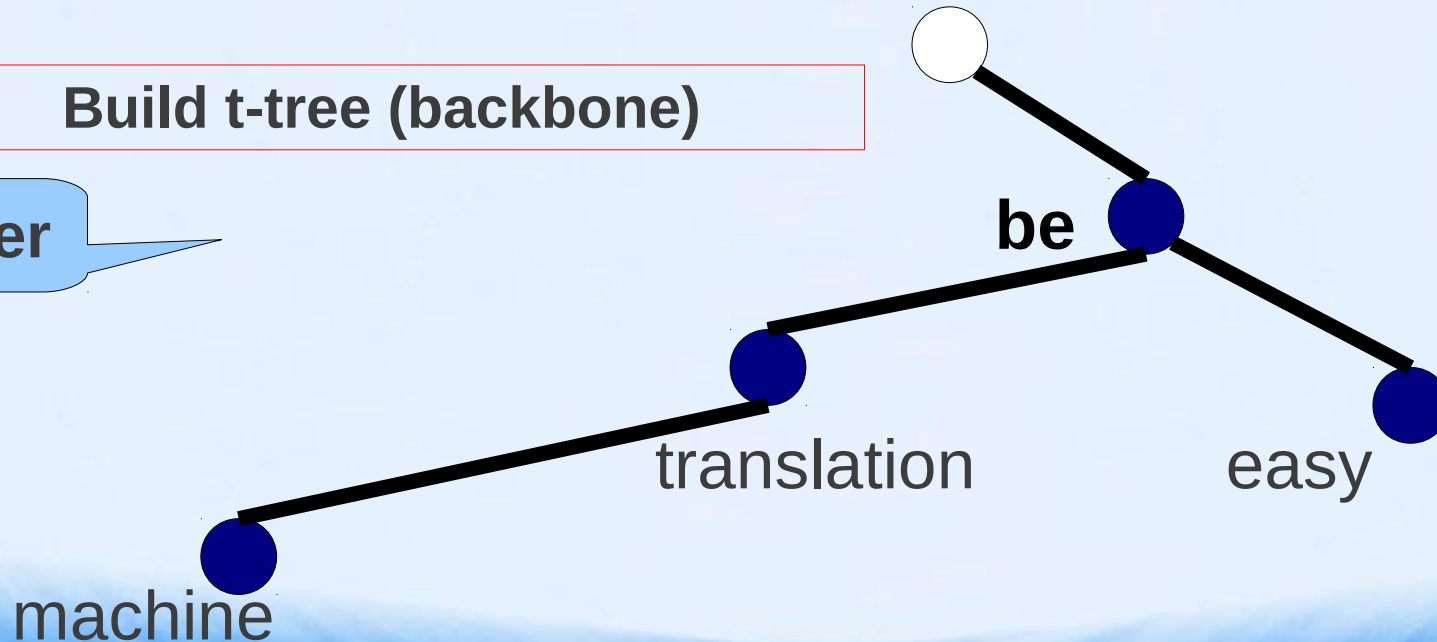
Machine translation should be easy.

m-layer

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

Build t-tree (backbone)

t-layer



Demo Translation – Analysis

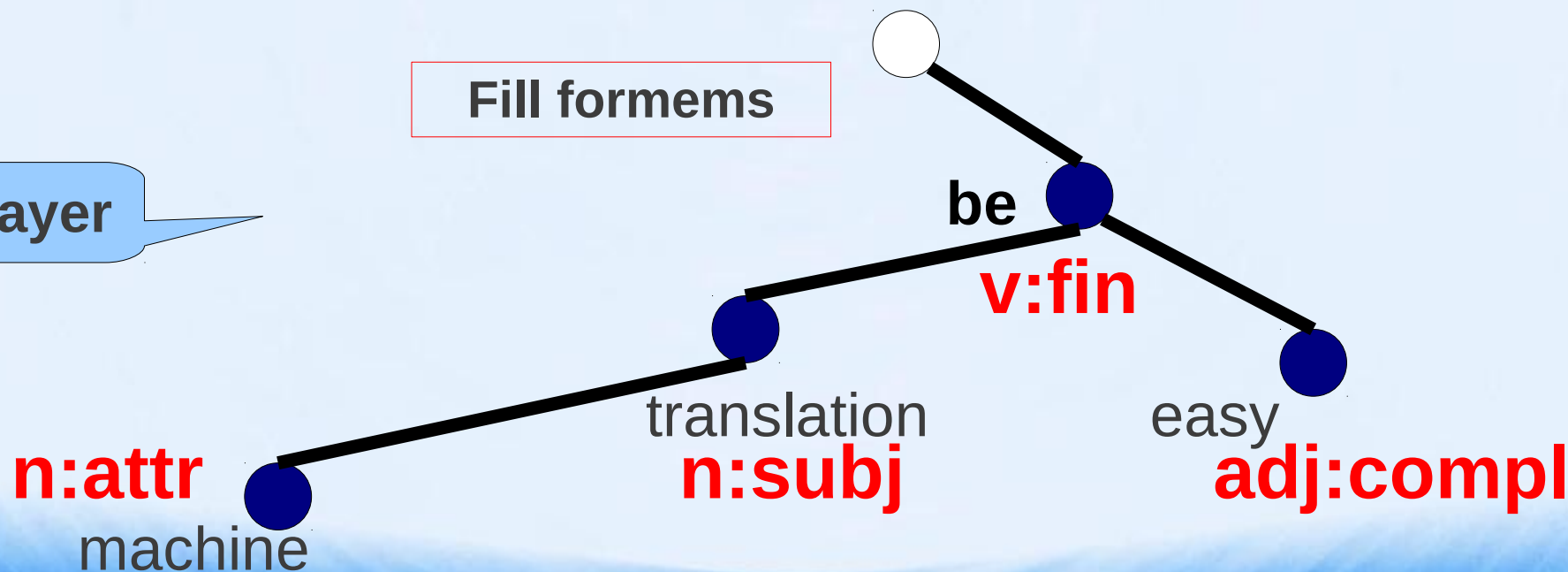
raw text

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

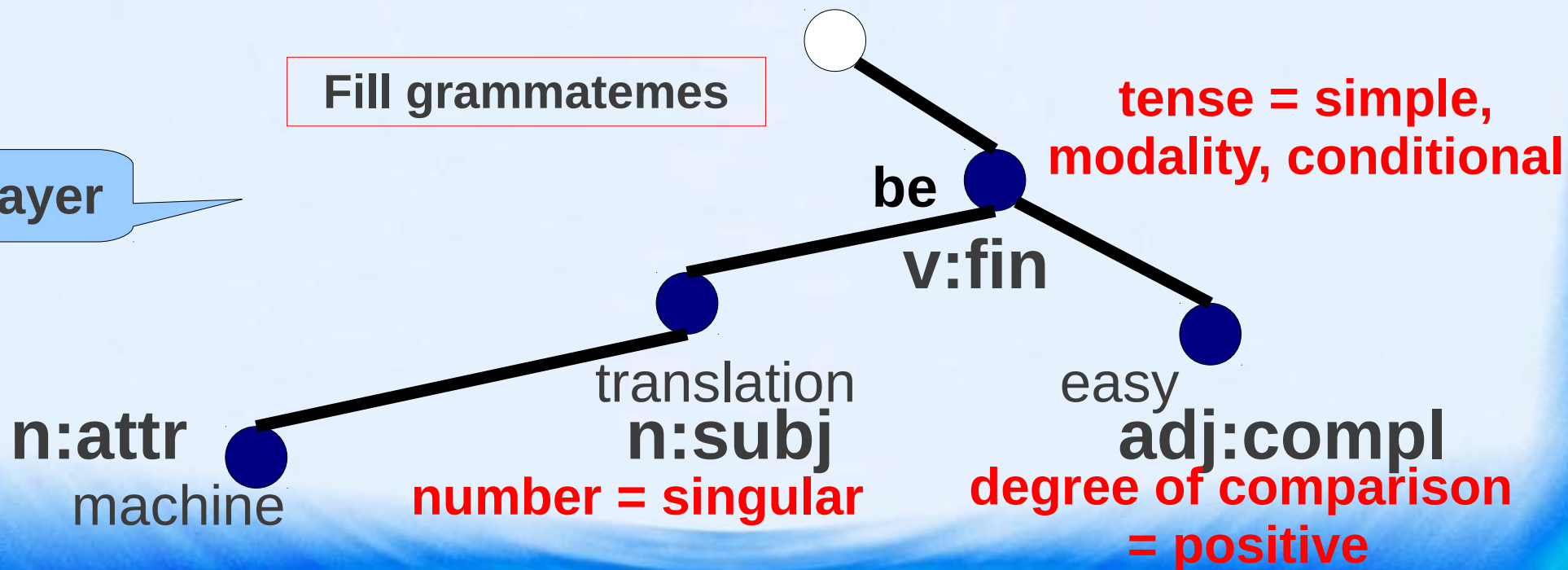
● ● ● ● ● ●

machine translation should be easy .

 NN NN MD VB JJ .

Fill grammatememes

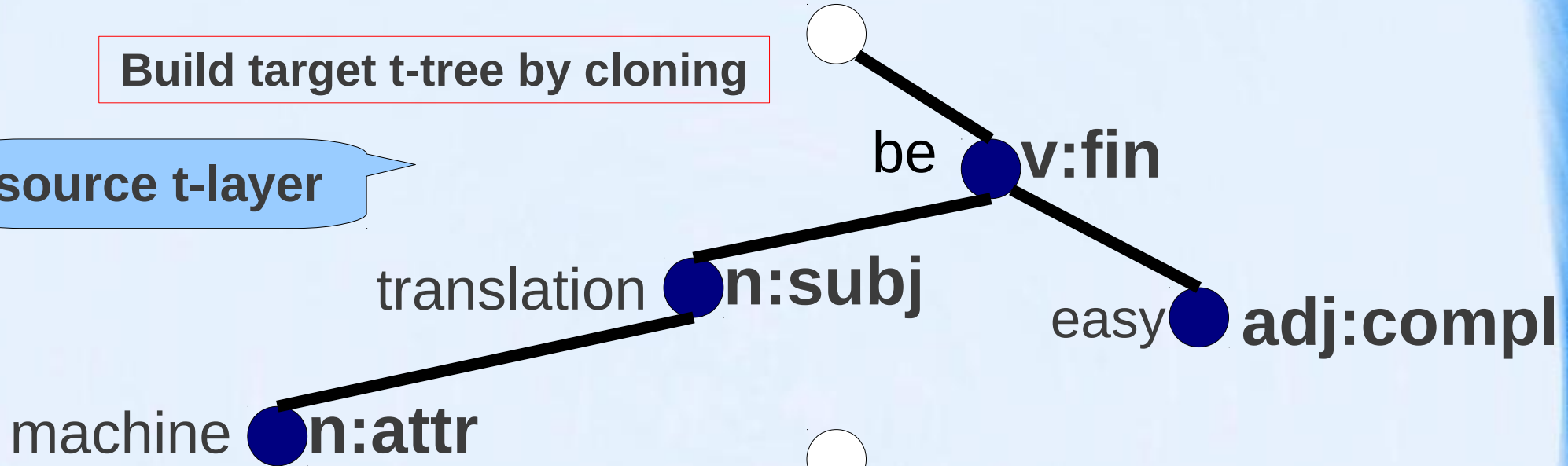
t-layer



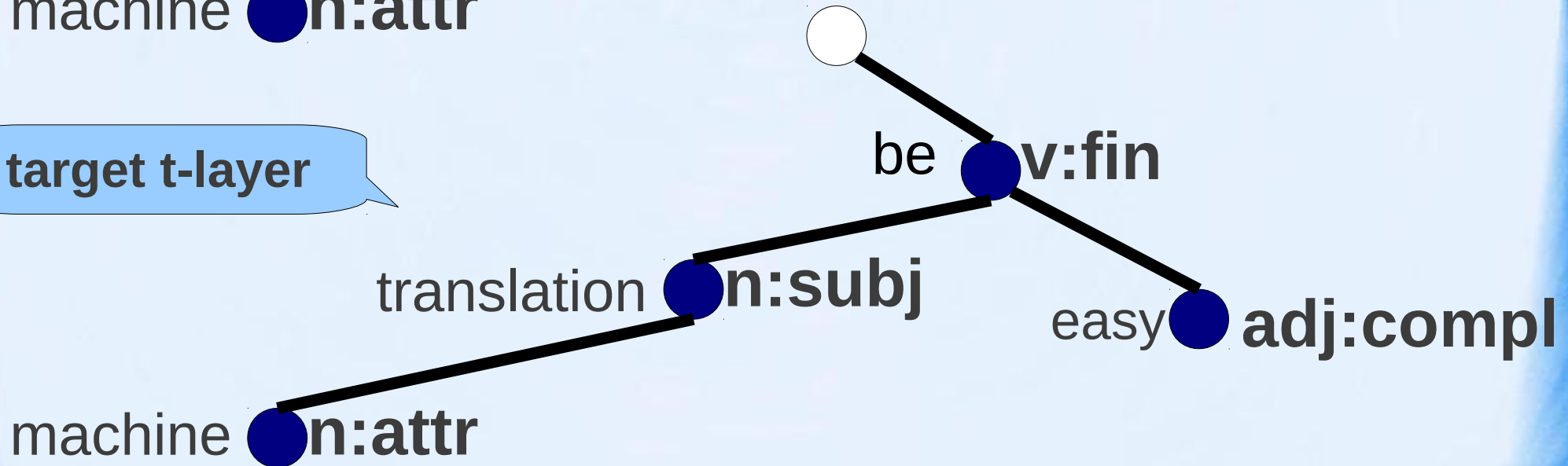
Demo Translation – Transfer

Build target t-tree by cloning

source t-layer



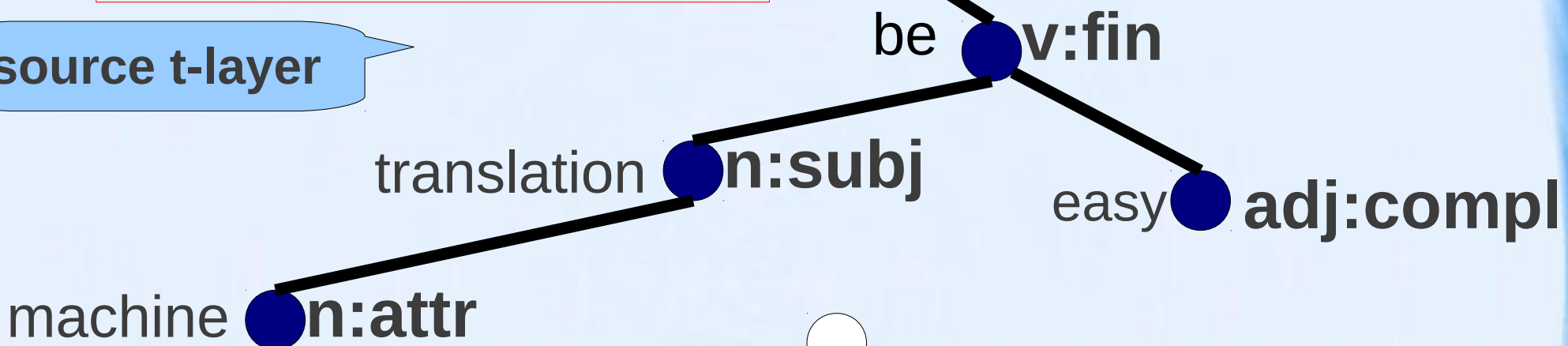
target t-layer



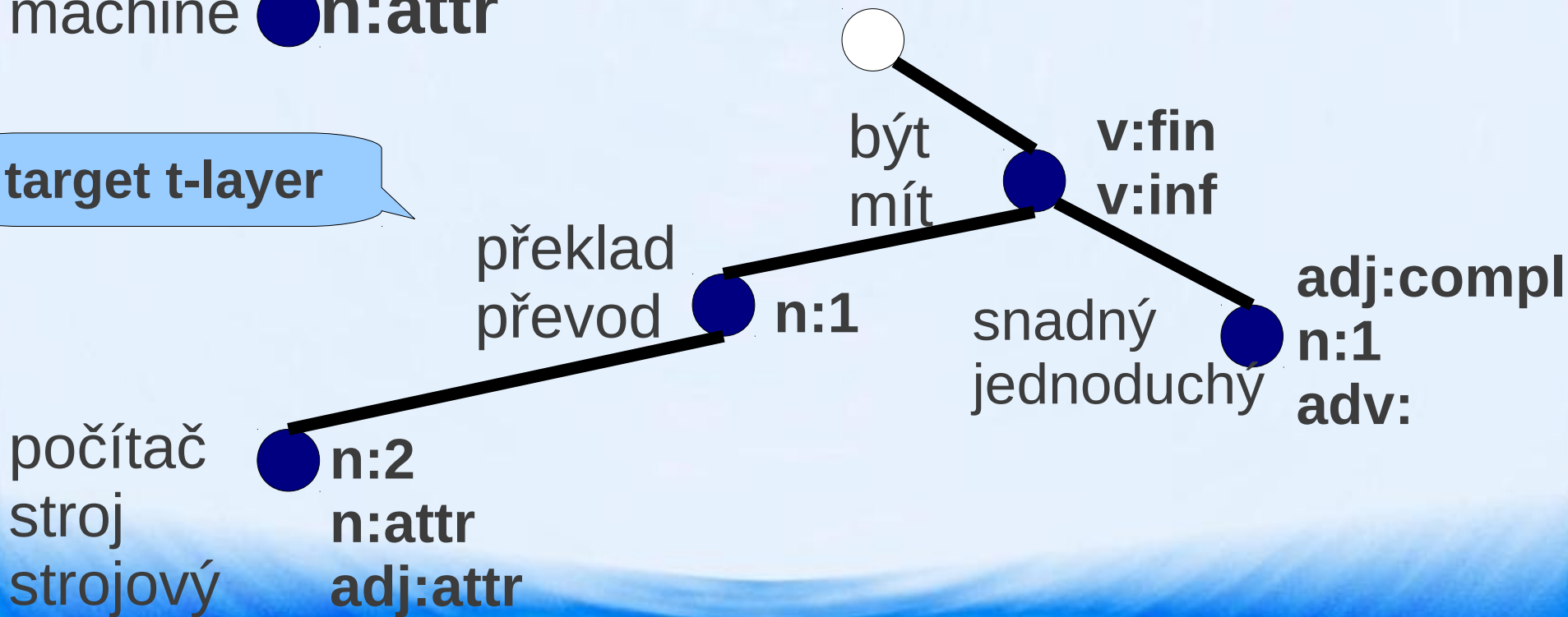
Demo Translation – Transfer

Get translation variants for lemmas and formems

source t-layer



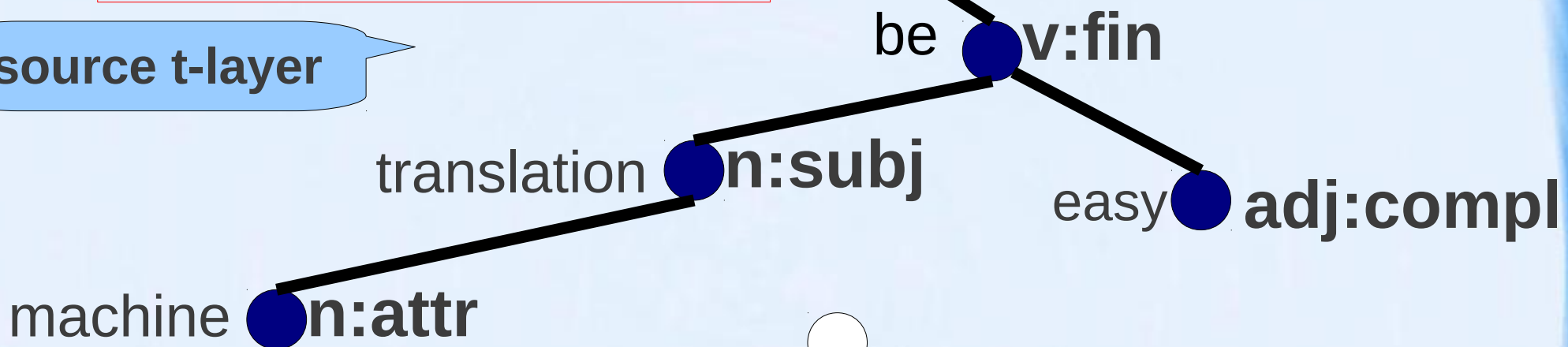
target t-layer



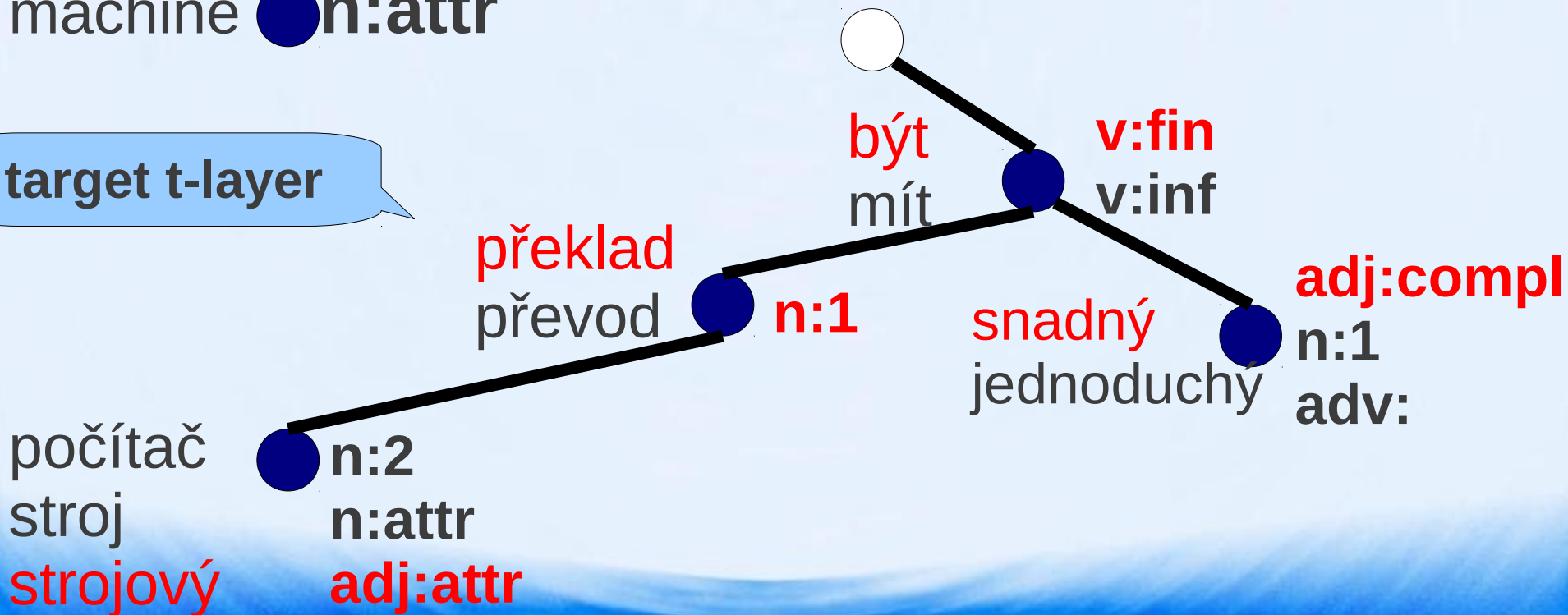
Demo Translation – Transfer

Select the best combination of lemmas and formems

source t-layer



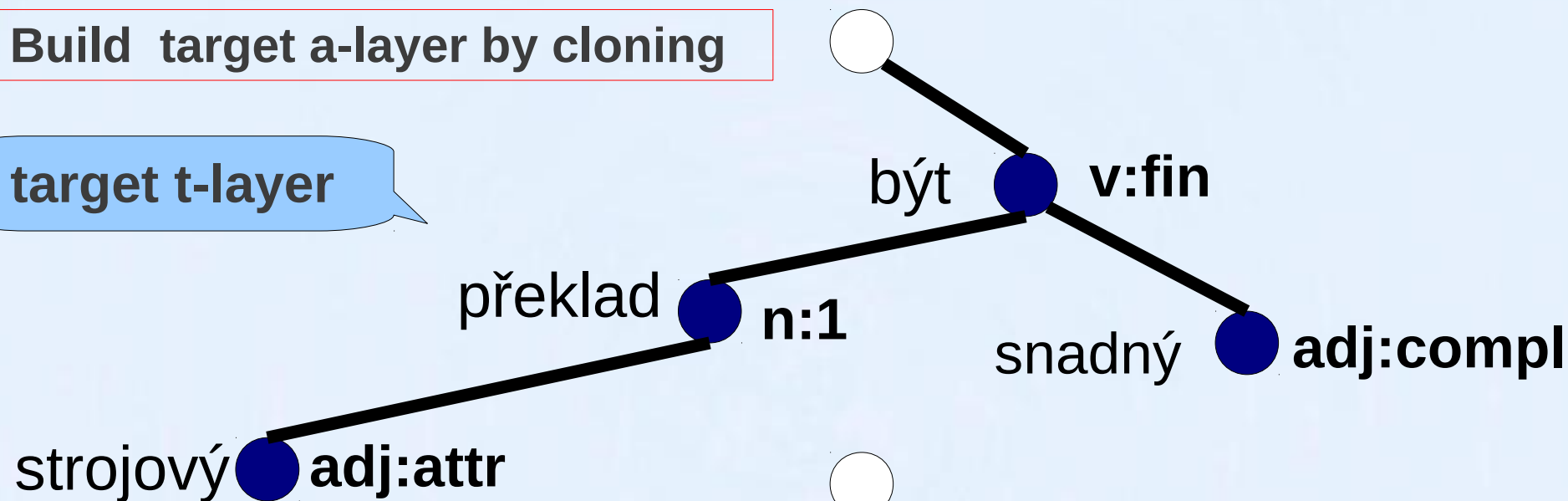
target t-layer



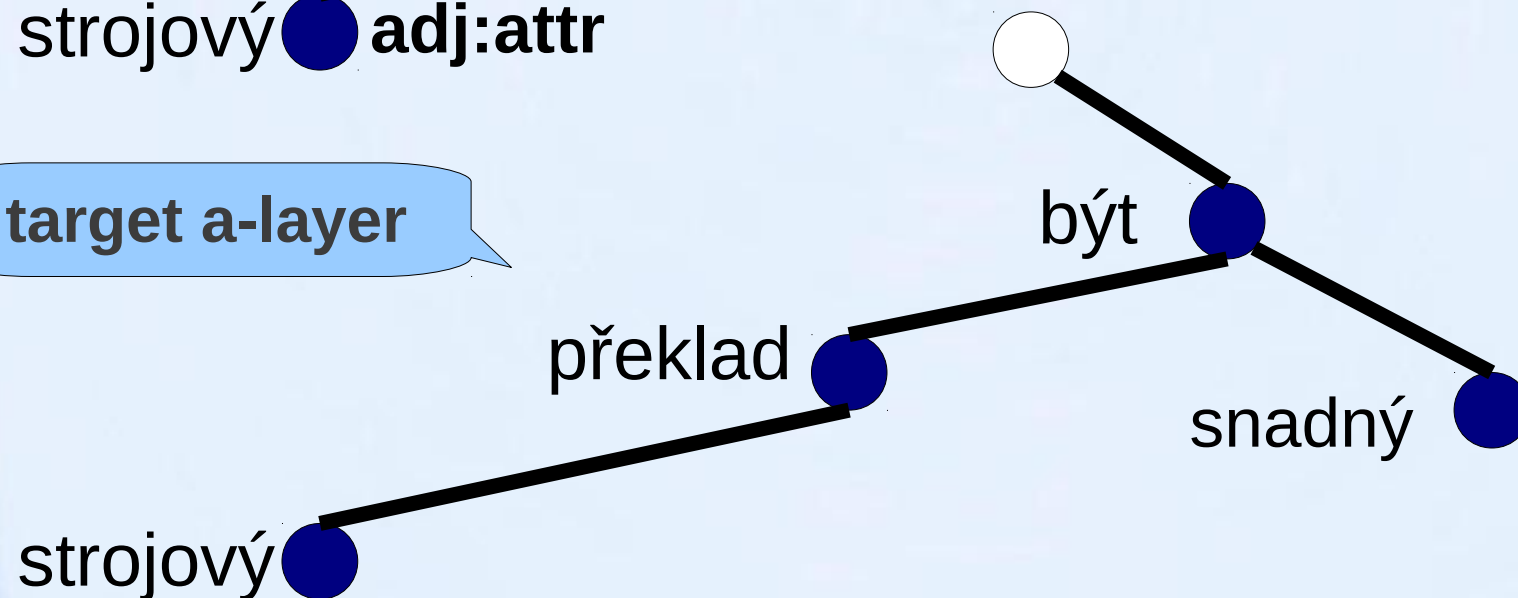
Demo Translation – Synthesis

Build target a-layer by cloning

target t-layer



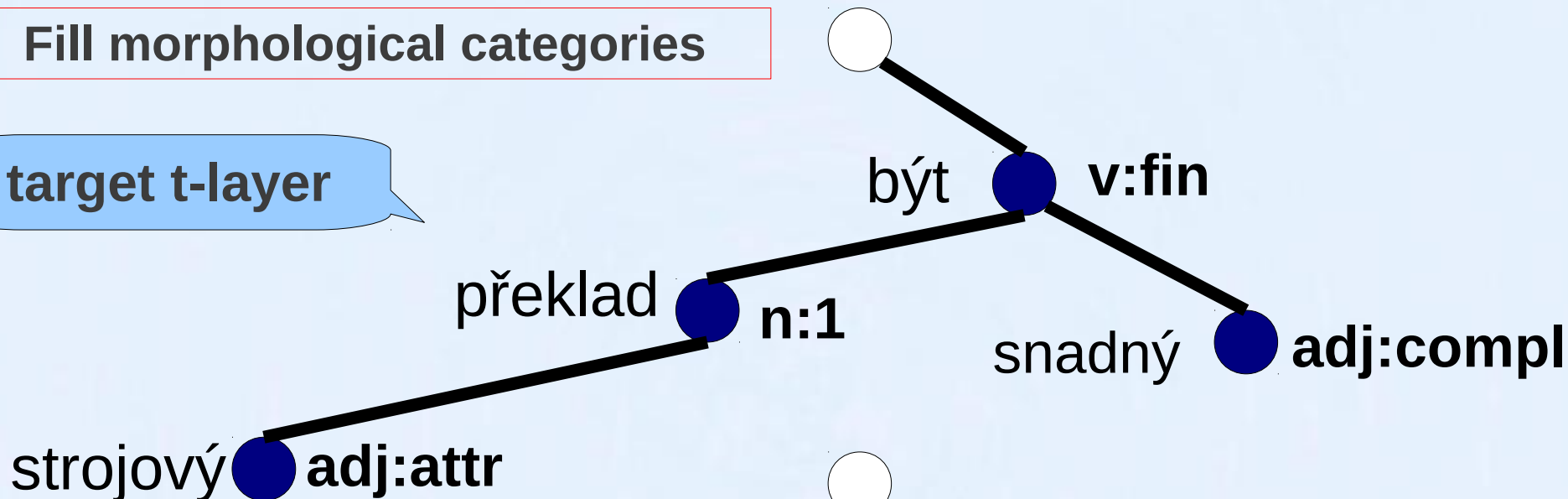
target a-layer



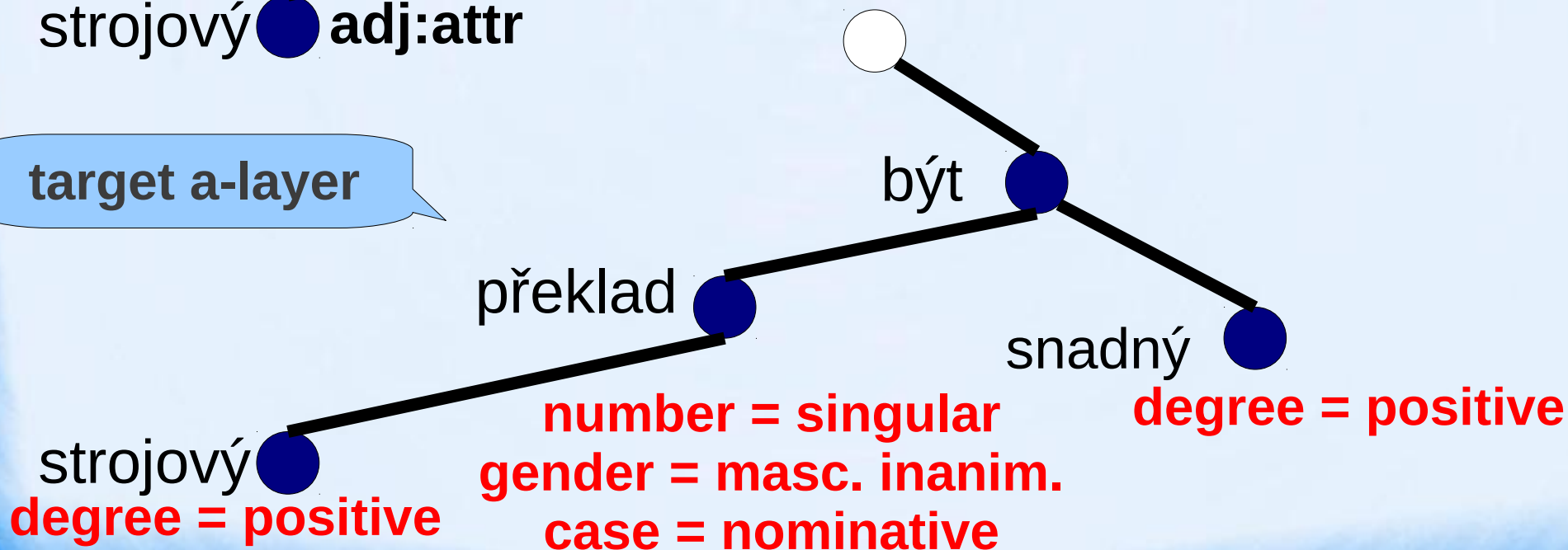
Demo Translation – Synthesis

Fill morphological categories

target t-layer



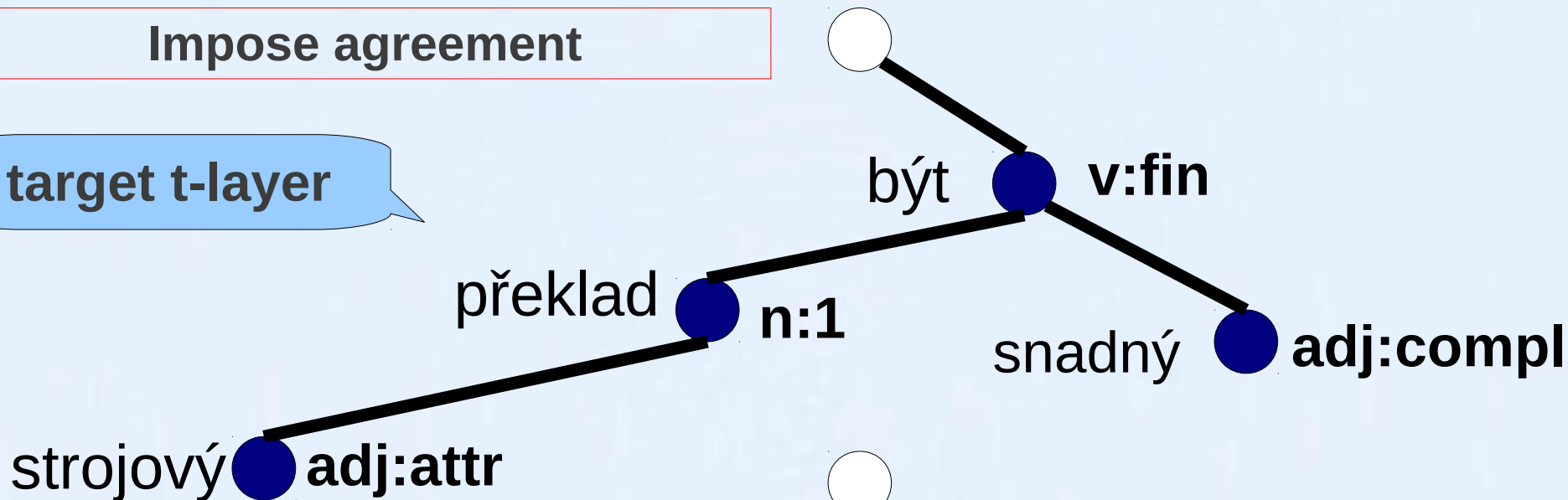
target a-layer



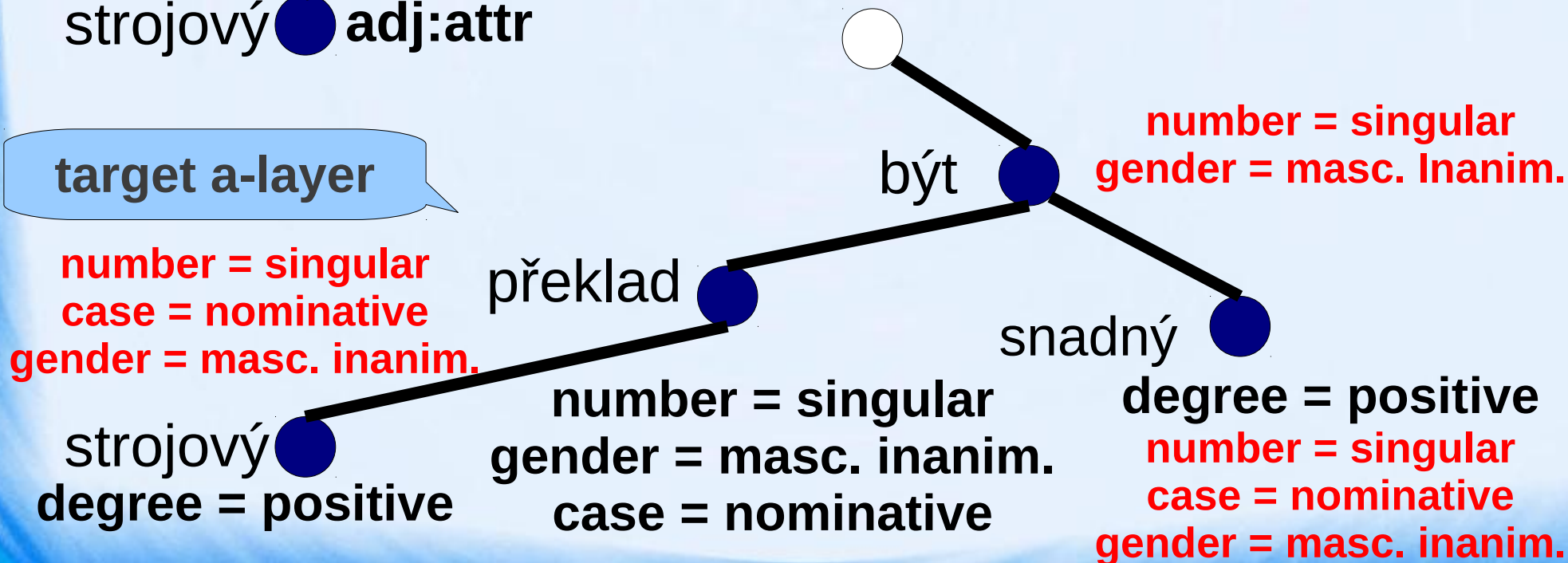
Demo Translation – Synthesis

Impose agreement

target t-layer



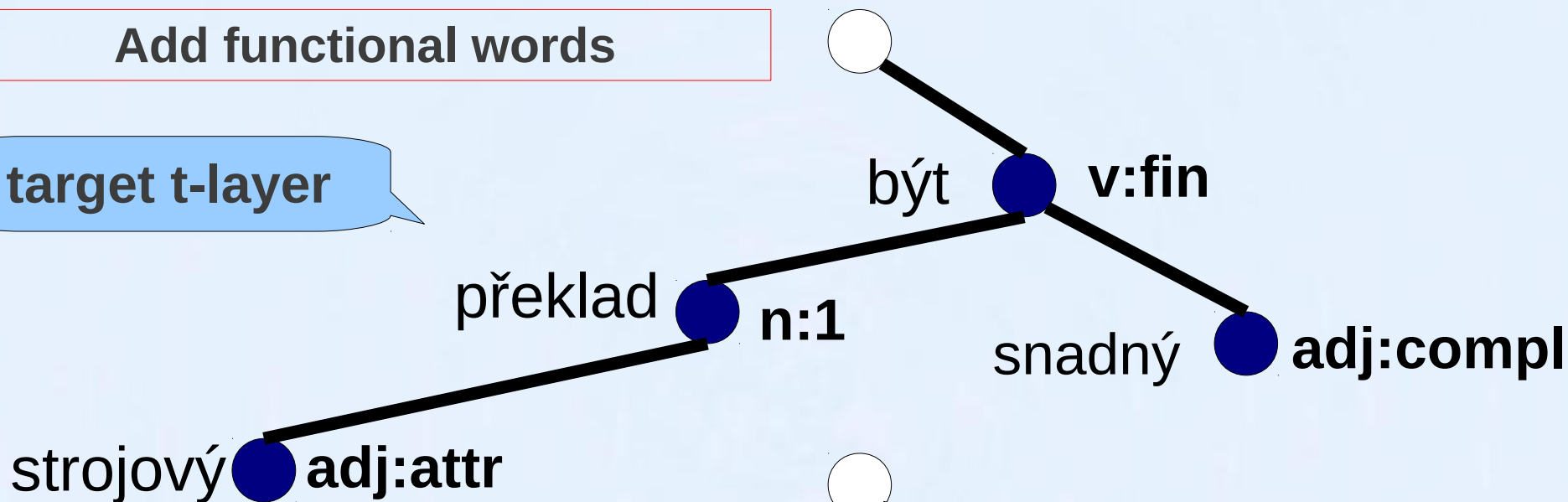
target a-layer



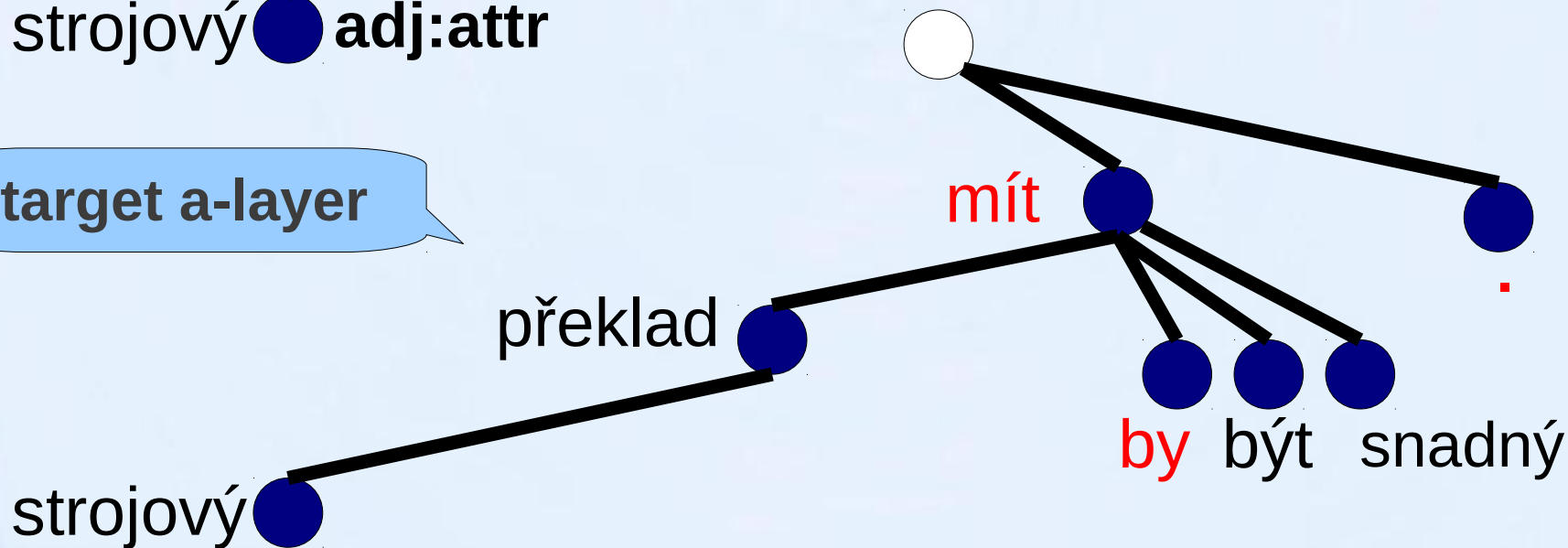
Demo Translation – Synthesis

Add functional words

target t-layer



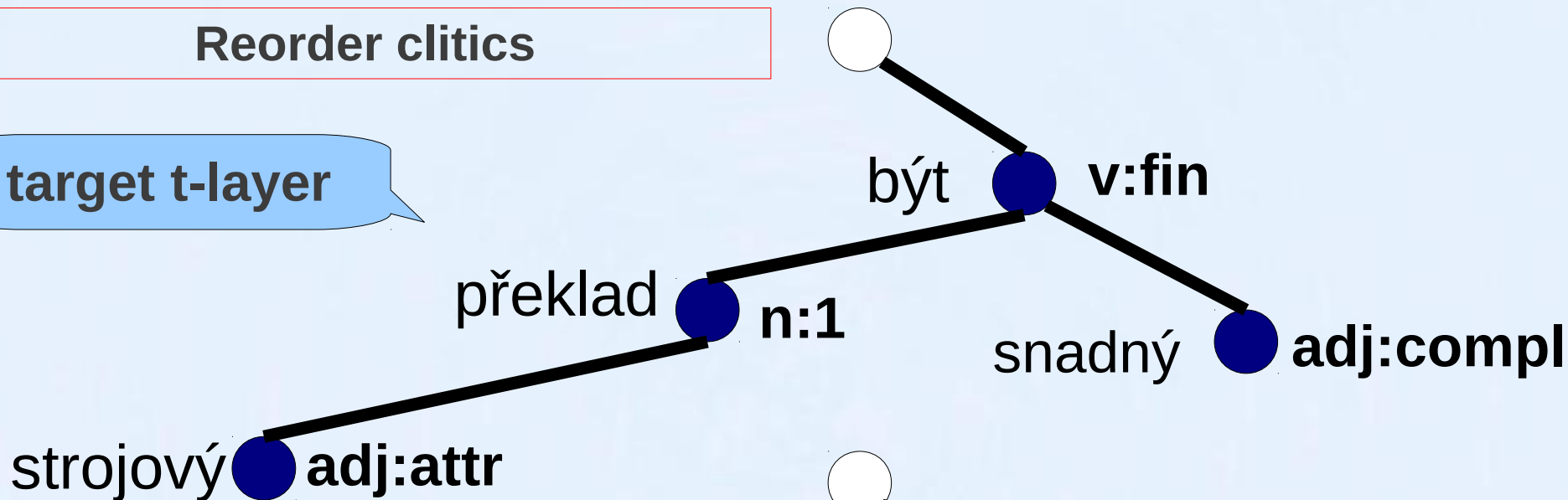
target a-layer



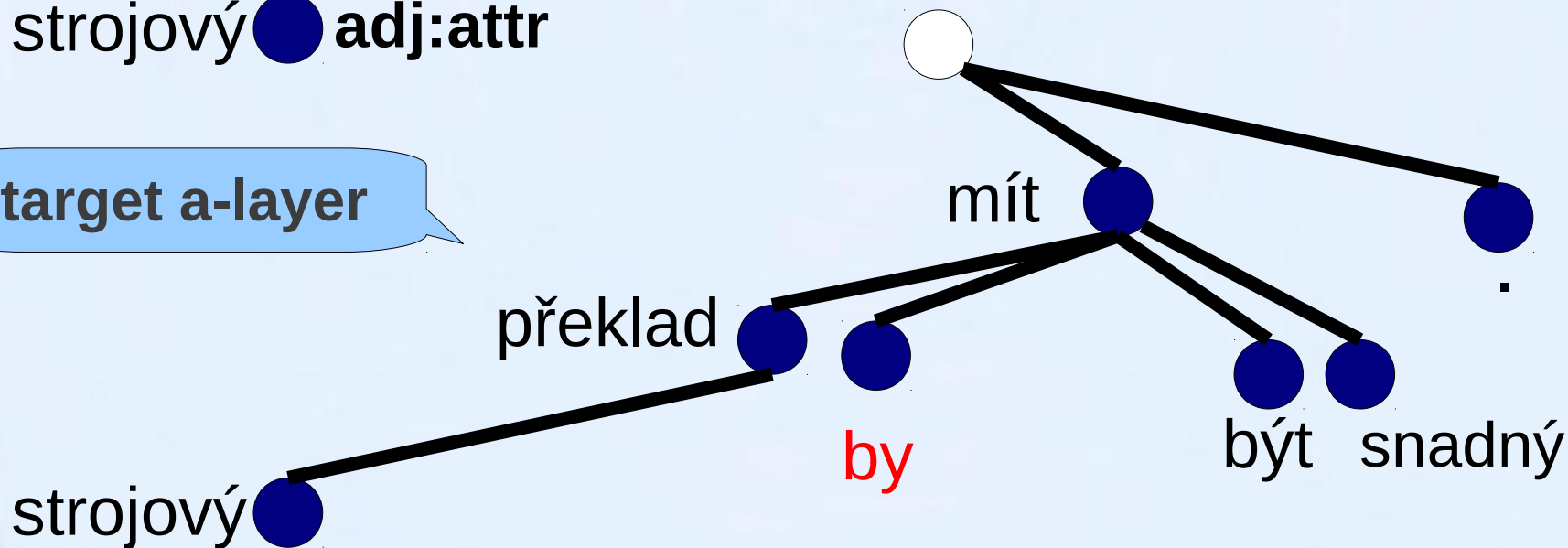
Demo Translation – Synthesis

Reorder clitics

target t-layer



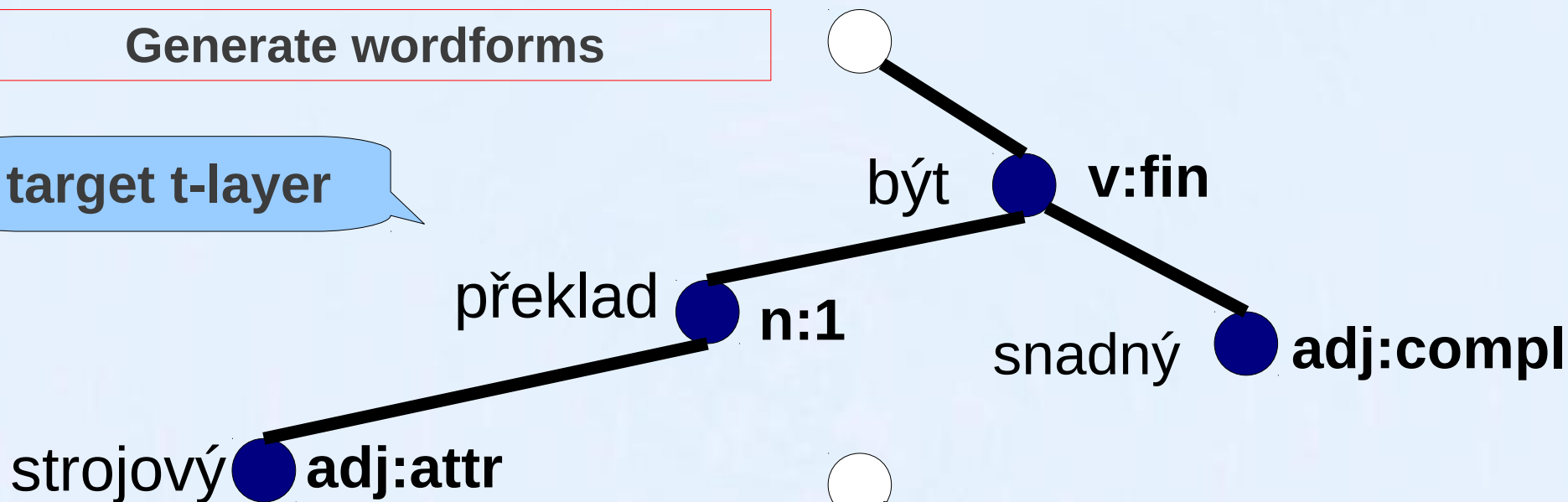
target a-layer



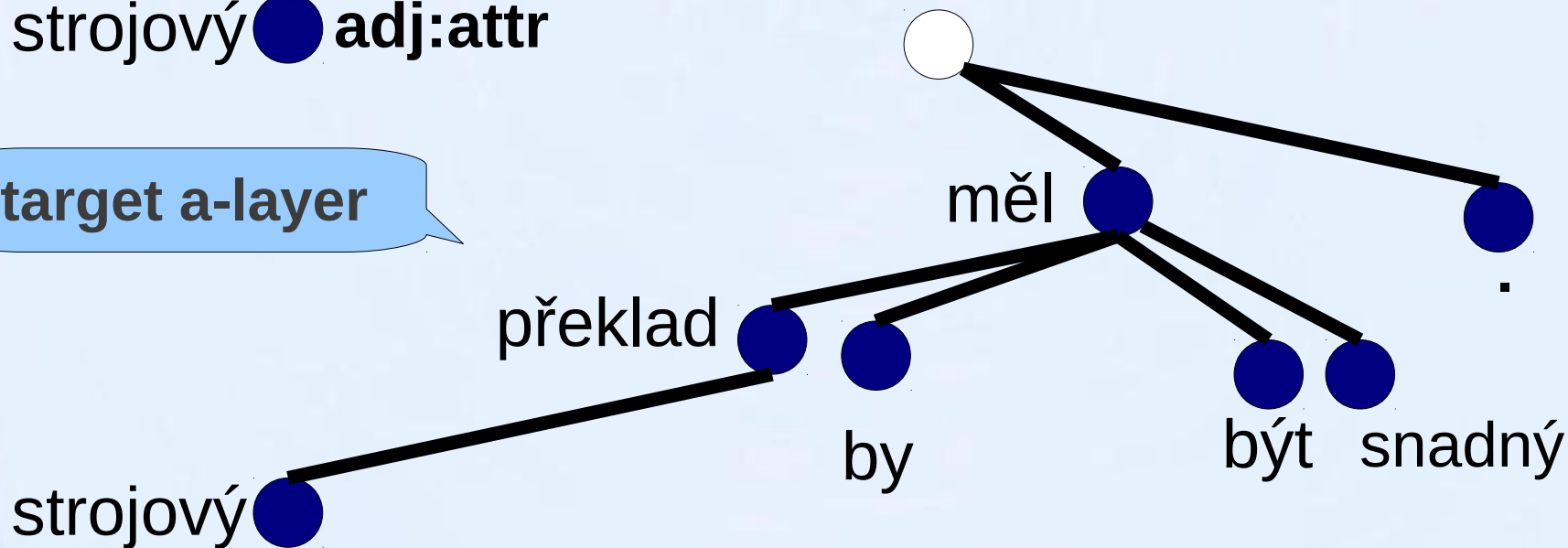
Demo Translation – Synthesis

Generate wordforms

target t-layer



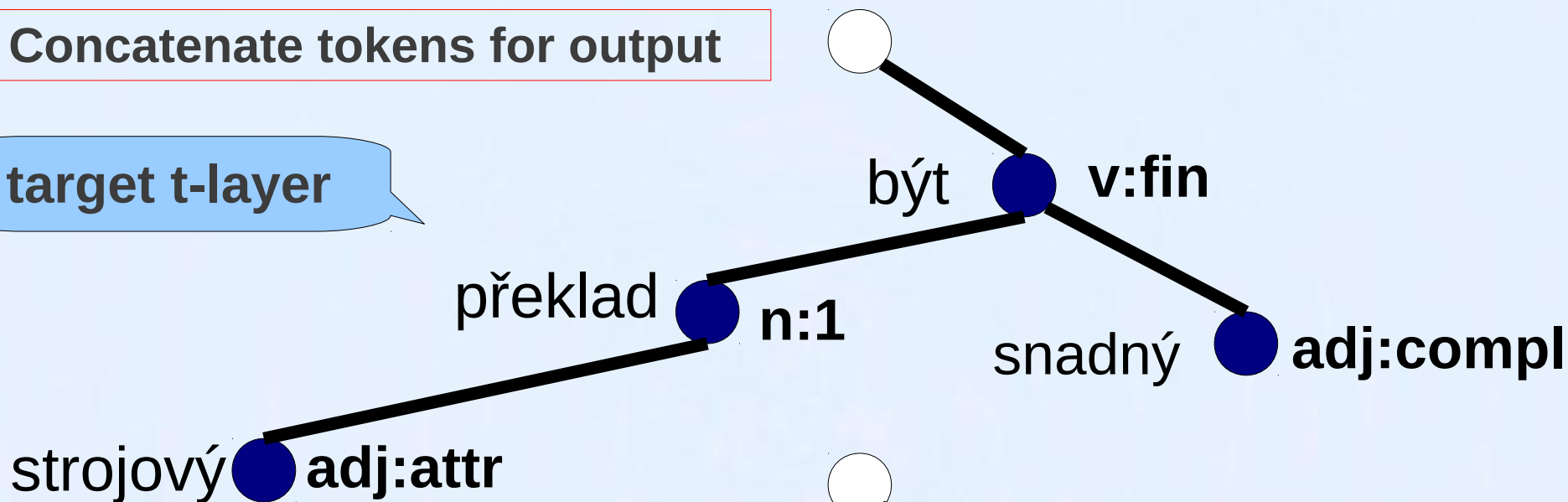
target a-layer



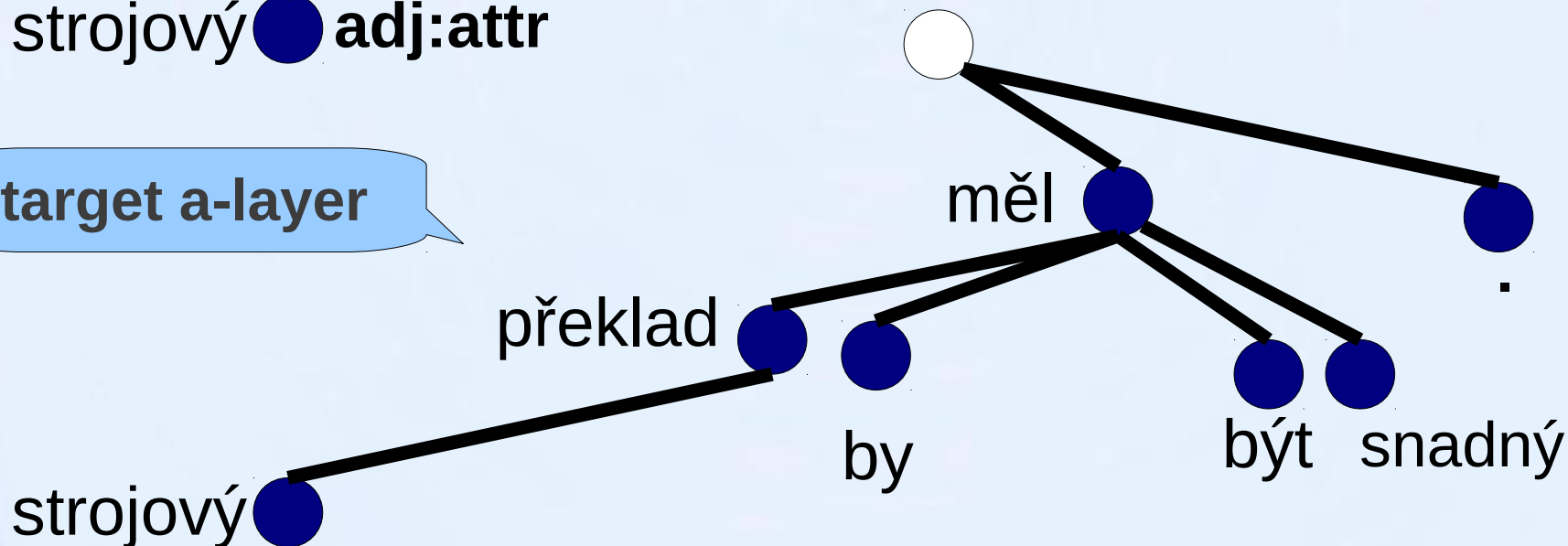
Demo Translation – Synthesis

Concatenate tokens for output

target t-layer



target a-layer



Strojový překlad by měl být snadný.

Demo Translation – Real Scenario



MORPHOLOGY:

ResegmentSentences

Tokenize

NormalizeForms

FixTokenization

TagMorce

FixTags

Lemmatize

NAMED ENTITIES:

StanfordNamedEntities

DistinguishPersonalNames

A-LAYER:

MarkChunks

ParseMST

SetIsMemberFromDeprel

RehangConllToPdtStyle

FixNominalGroups

FixIsMember

FixAtree

FixMultiwordPrepAndConj

FixDicendiVerbs

SetAfunAuxCPCoord

SetAfun

T-LAYER:

MarkEdgesToCollapse

MarkEdgesToCollapseNeg

BuildTree

SetIsMember

MoveAuxFromCoordToMembers

FixTlemmas

SetCoapFunctors

FixEitherOr

FixIsMember

MarkClauseHeads

MarkPassives

SetFunctors

MarkInfin

MarkRelClauseHeads

MarkRelClauseCoref

MarkDspRoot

MarkParentheses

SetNodetype

SetGrammatemes

SetFormeme

RehangSharedAttr

SetVoice

FixImperatives

SetIsNameOfPerson

SetGenderOfPerson

AddCorAct

FindTextCoref

TRANSFER:

CopyTree

TrLFPPhrases

TrLFJointStatic

DeleteSuperfluousTnodes

TrFTryRules

TrFAddVariants

TrFRerank

TrLTryRules

TrLAddVariants

TrLFNumeralsByRules

TrLFilterAspect

TransformPassiveConstructions

PrunePersonalNameVariants

RemoveUnpassivableVariants

TrLFCompounds

CutVariants

RehangToEffParents

TrLFTreeViterbi

RehangToOrigParents

CutVariants

FixTransferChoices

ReplaceVerbWithAdj

DeletePossPronBeforeVlastni

TrLFFemaleSurnames

AddNounGender

MarkNewRelClauses

AddRelpronBelowRc

ChangeCorToPersPron

AddPersPronBelowVfin

AddVerbAspect

FixDateTime

FixGrammatemesAfterTransfer

FixNegation

MoveAdjsBeforeNouns

MoveGenitivesRight

MoveRelClauseRight

MoveDicendiCloserToDsp

MovePersPronNextToVerb

MoveEnoughBeforeAdj

MoveJesteBeforeVerb

FixMoney

OverridePpWithPhraseTr

FindGramCorefForRefIIPron

NeutPersPronGenderFromAntec

ValencyRelatedRules

SetClauseNumber

TurnTextCorefToGramCoref

SYNTHESIS TO A-LAYER:

CopyTree

DistinguishHomonymous.

ReverseNumberNounDep.

InitMorphcat

FixPossessiveAdjs

MarkSubject

ImposePronZAgr

ImposeRelPronAgr

ImposeSubjpredAgr

ImposeAttrAgr

ImposeComplAgr

DropSubjPersProns

AddPrepos

AddSubconj

AddReflexParticles

AddAuxVerbCompoundPassive

AddAuxVerbModal

AddAuxVerbCompoundFuture

AddAuxVerbConditional

AddAuxVerbCompoundPast

AddClausalExpletivePronouns

ResolveVerbs

ProjectClauseNumber

AddParentheses

AddSentFinalPunct

AddSubordClausePunct

AddCoordPunct

AddAppositionPunct

ChooseMlemmaForPersPron

GenerateWordforms

MoveCliticsToWackernagel

DeleteSuperfluousPrepos

DeleteEmptyNouns

VocalizePrepos

CapitalizeSentStart

CapitalizeNamedEntities.

FillTagFromMorphcat

SYNTHESIS TO TEXT:

ConcatenateTokens

ApplySubstitutions

DetokenizeUsingRules

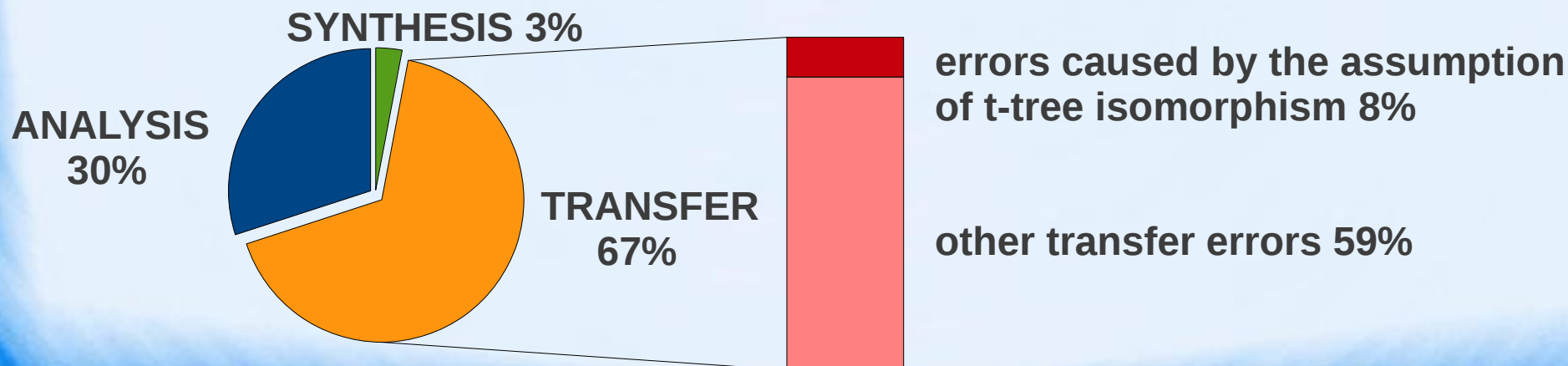
RemoveRepeatedTokens

NormalizePunctuationForWMT

Annotation of Translation Errors

sample of 250 sentences, 1463 errors in total

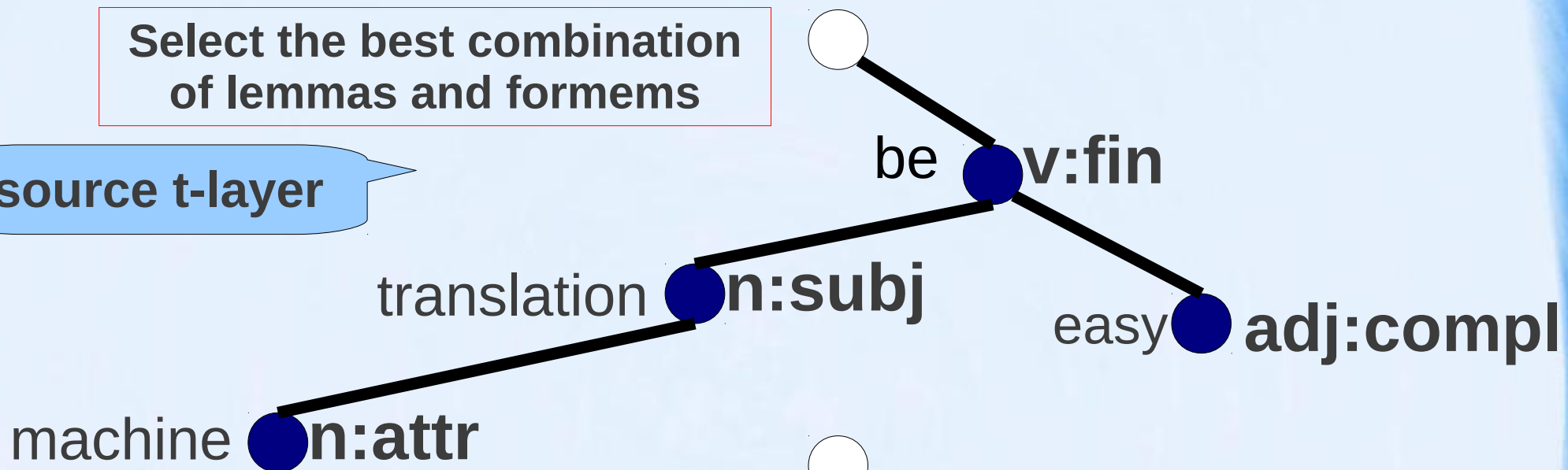
Type	lemma, formeme, gram., w. order,...
Subtype	gram: gender, person, tense,...
Seriousness	serious, minor
Circumstances	coordination, named entity, numbers
Source	tok, lem, tagger, parser, tecto, trans, x, syn, ?



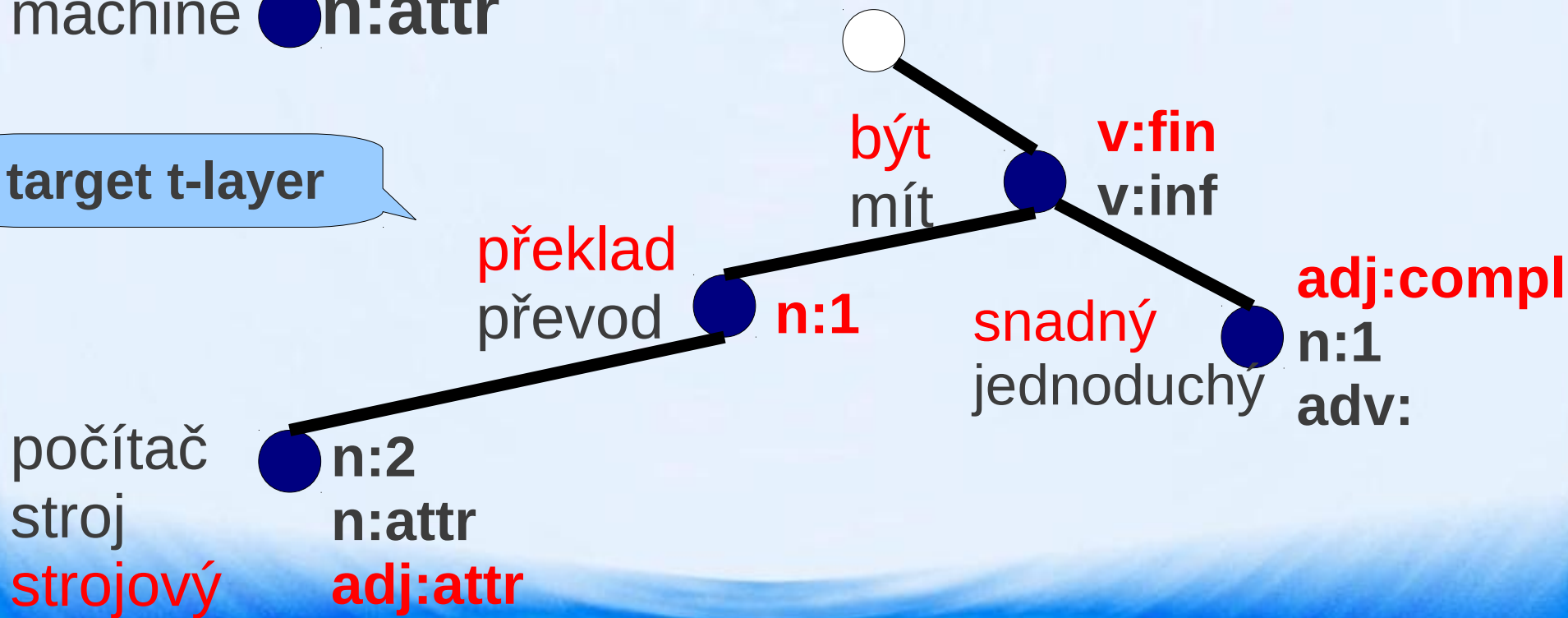
HMTM – Motivation

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

machine **n:attr**

translation
n:subj

be **v:fin**

easy **adj:compl**

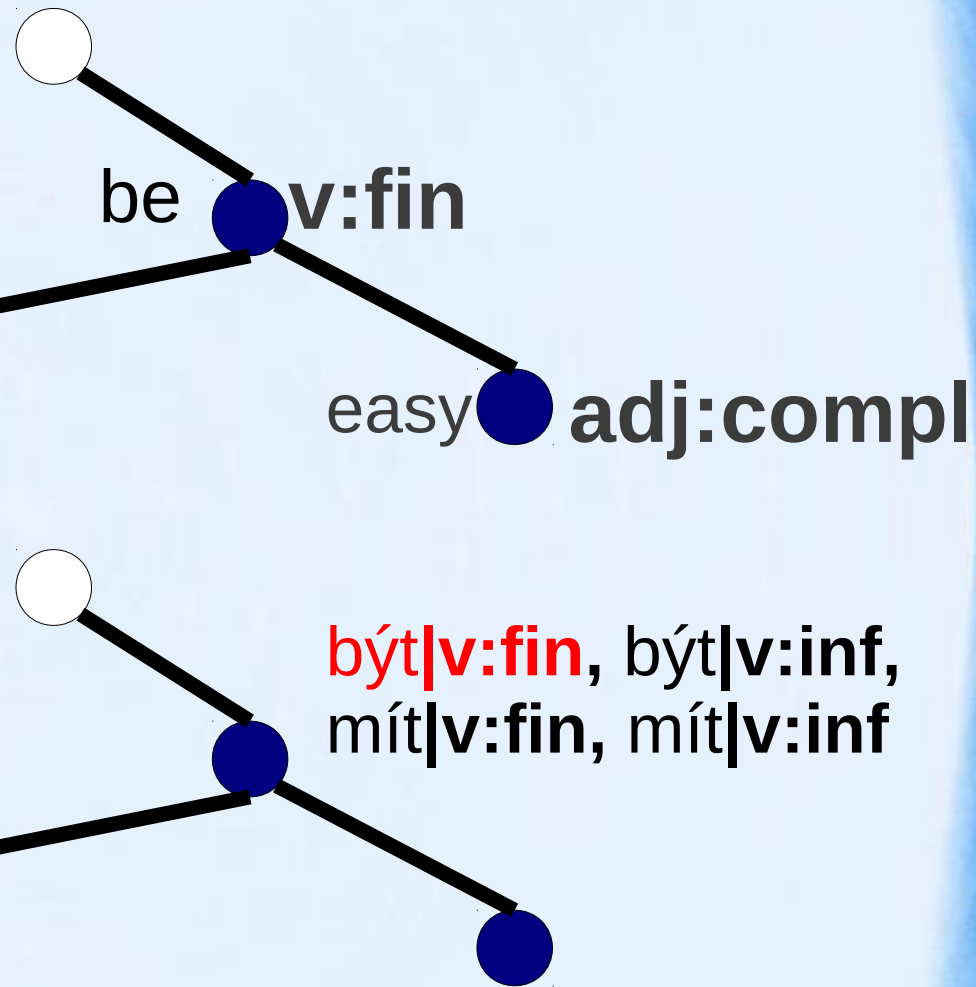
target t-layer

překlad|n:1,
převod|n:1

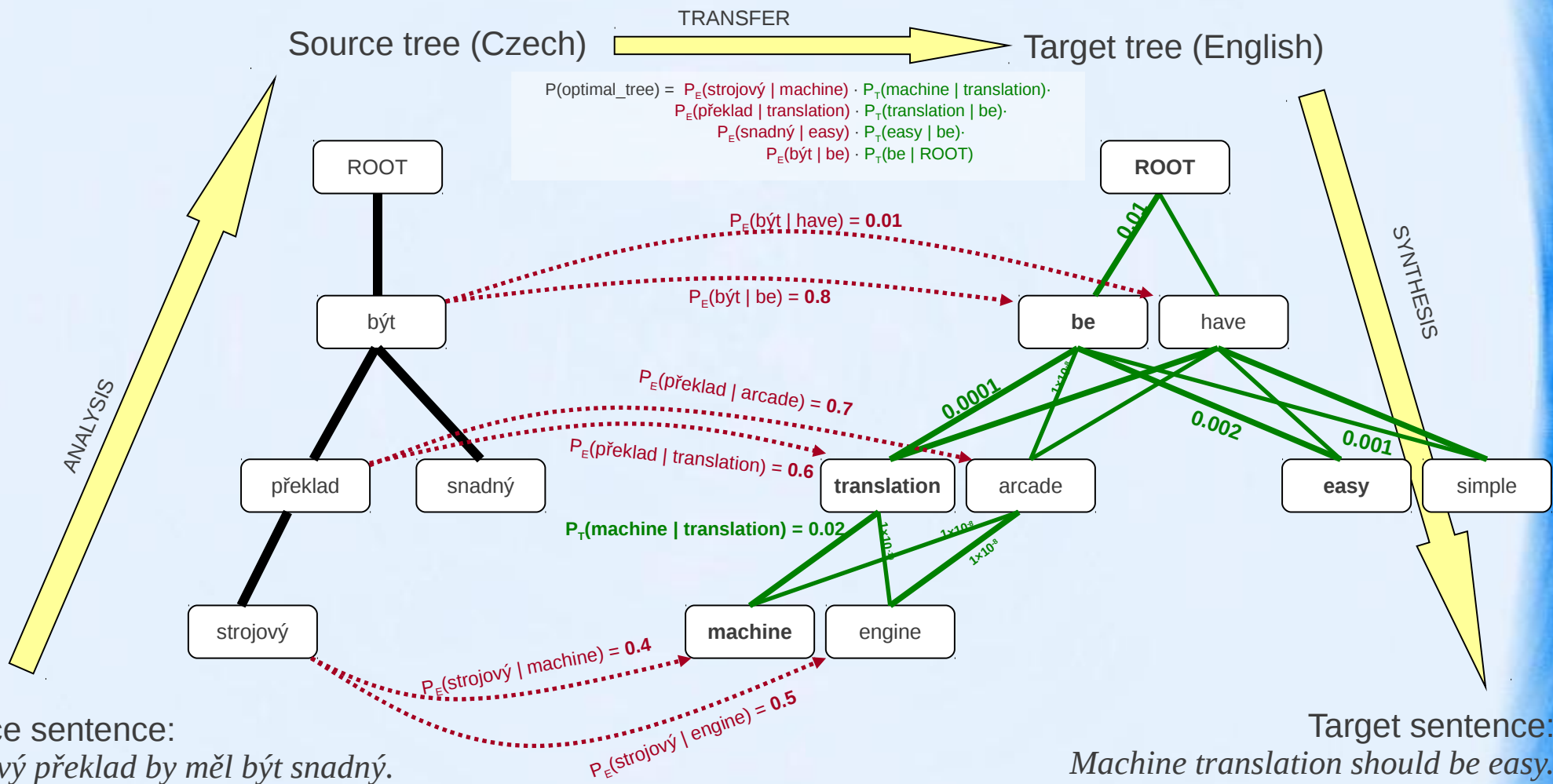
být|v:fin, být|v:inf,
mít|v:fin, mít|v:inf

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 → hierachy

- Four basic types of dictionaries:

Static plain

loaded from a file „lemma → lemma“

Context

loaded from a file „lemma,features → lemma“

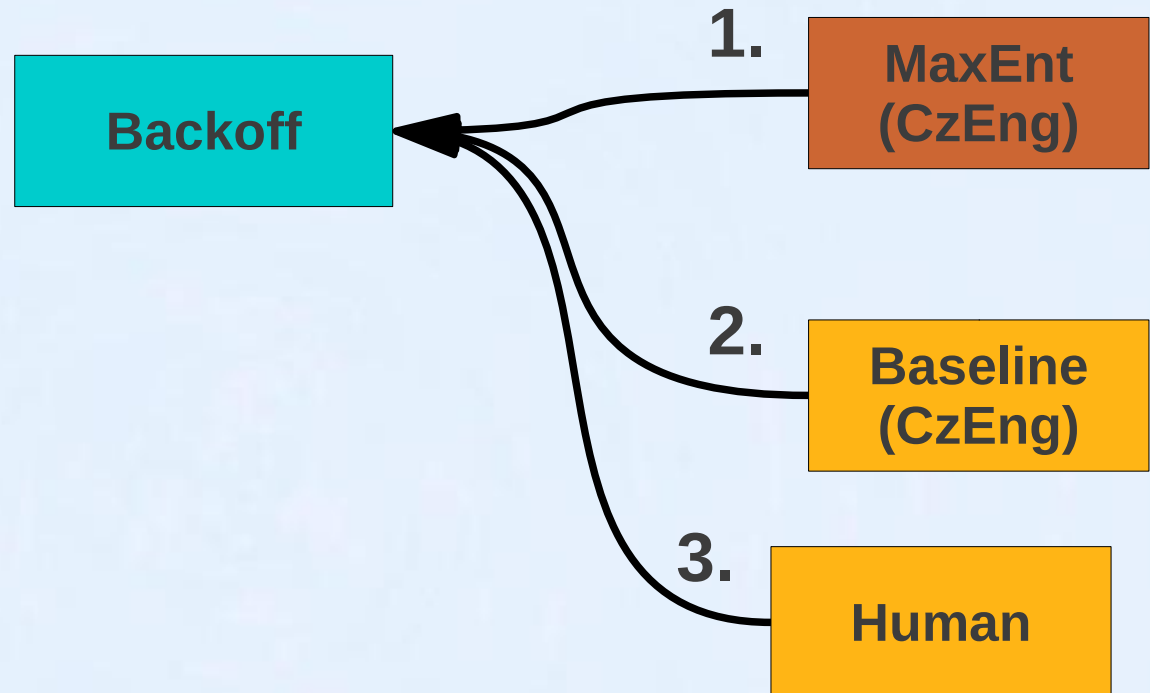
Derivational

translations derived dynamicaly, input dictionary

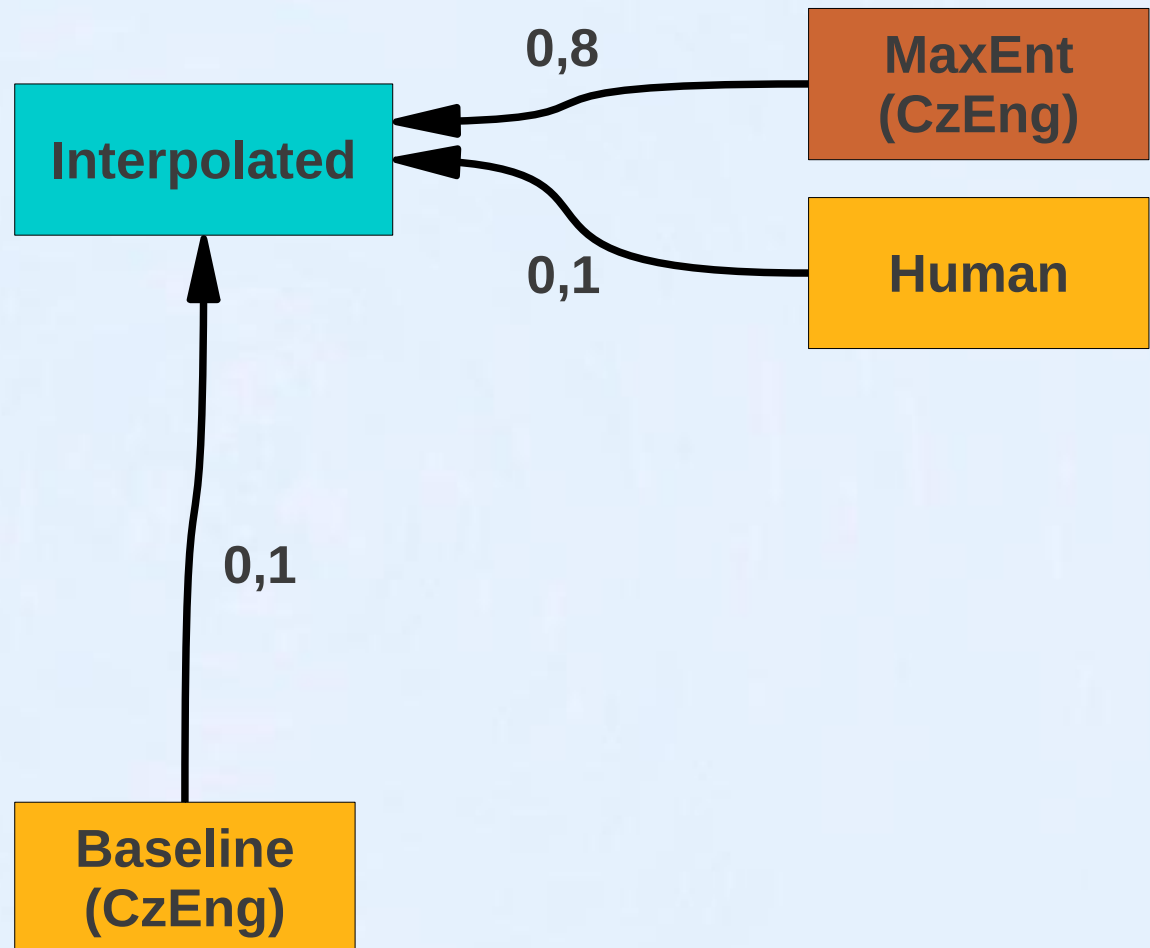
Combinaional

combination of more input dictionaries

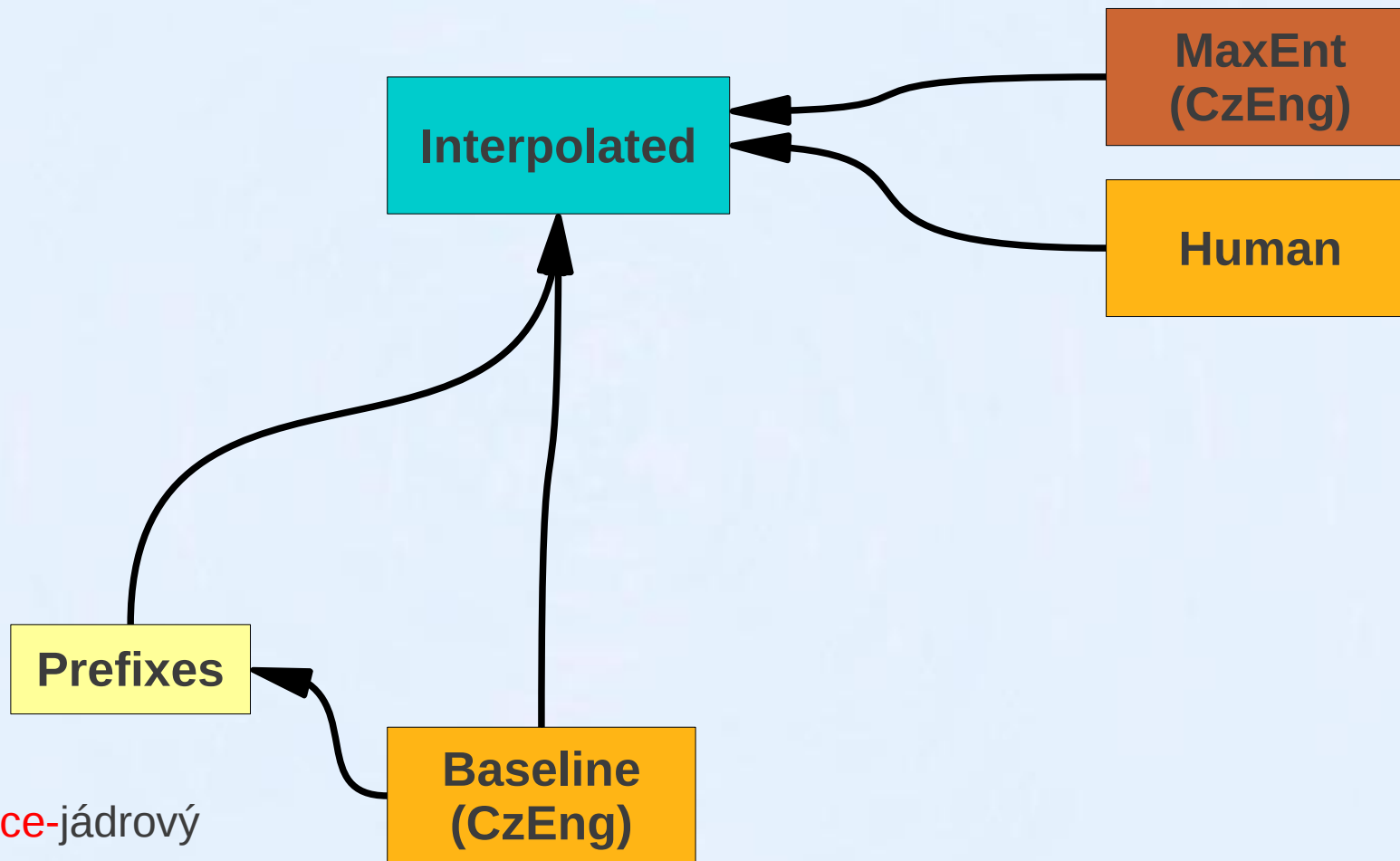
Hierarchy of lemma dictionaries



Hierarchy of lemma dictionaries

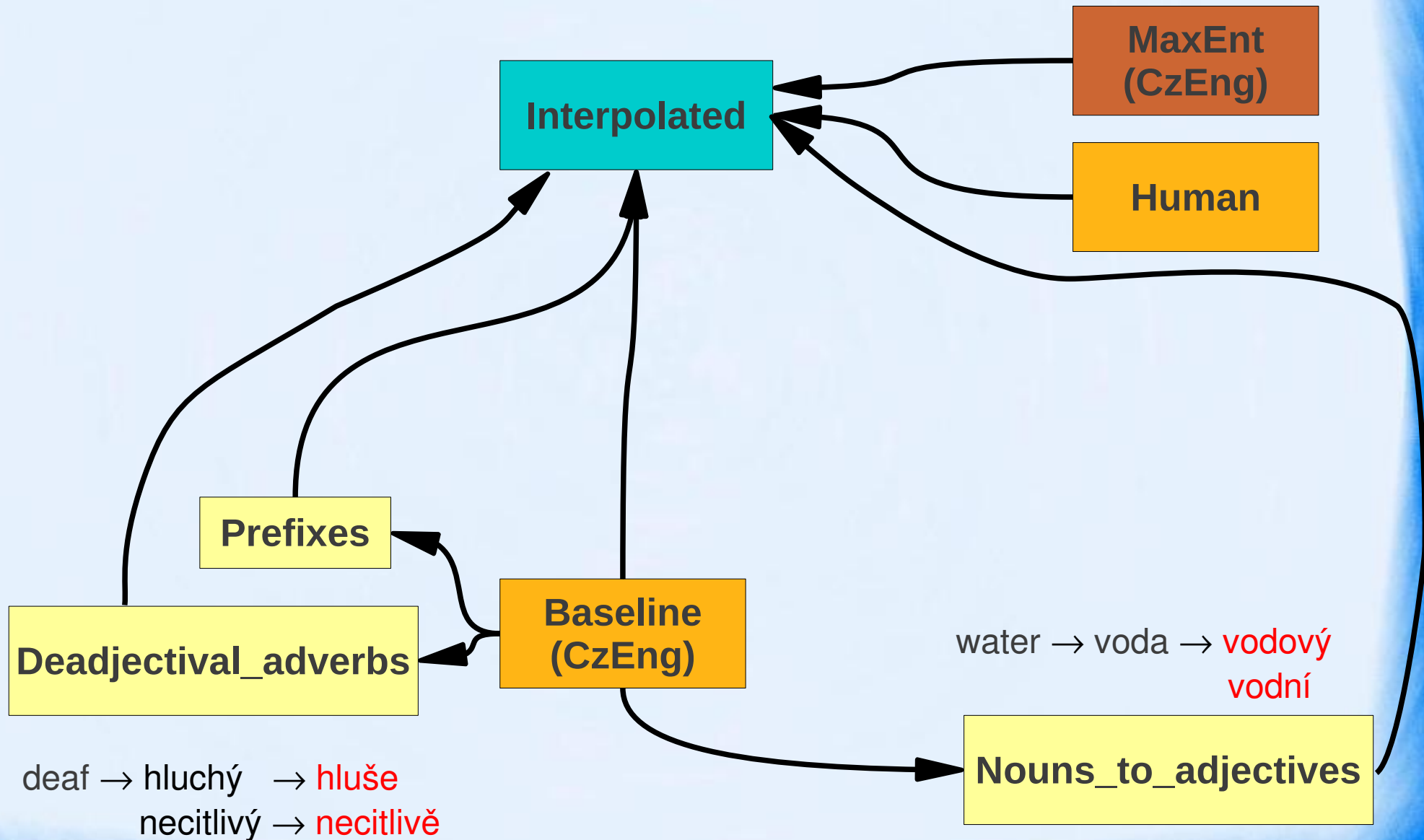


Hierarchy of lemma dictionaries

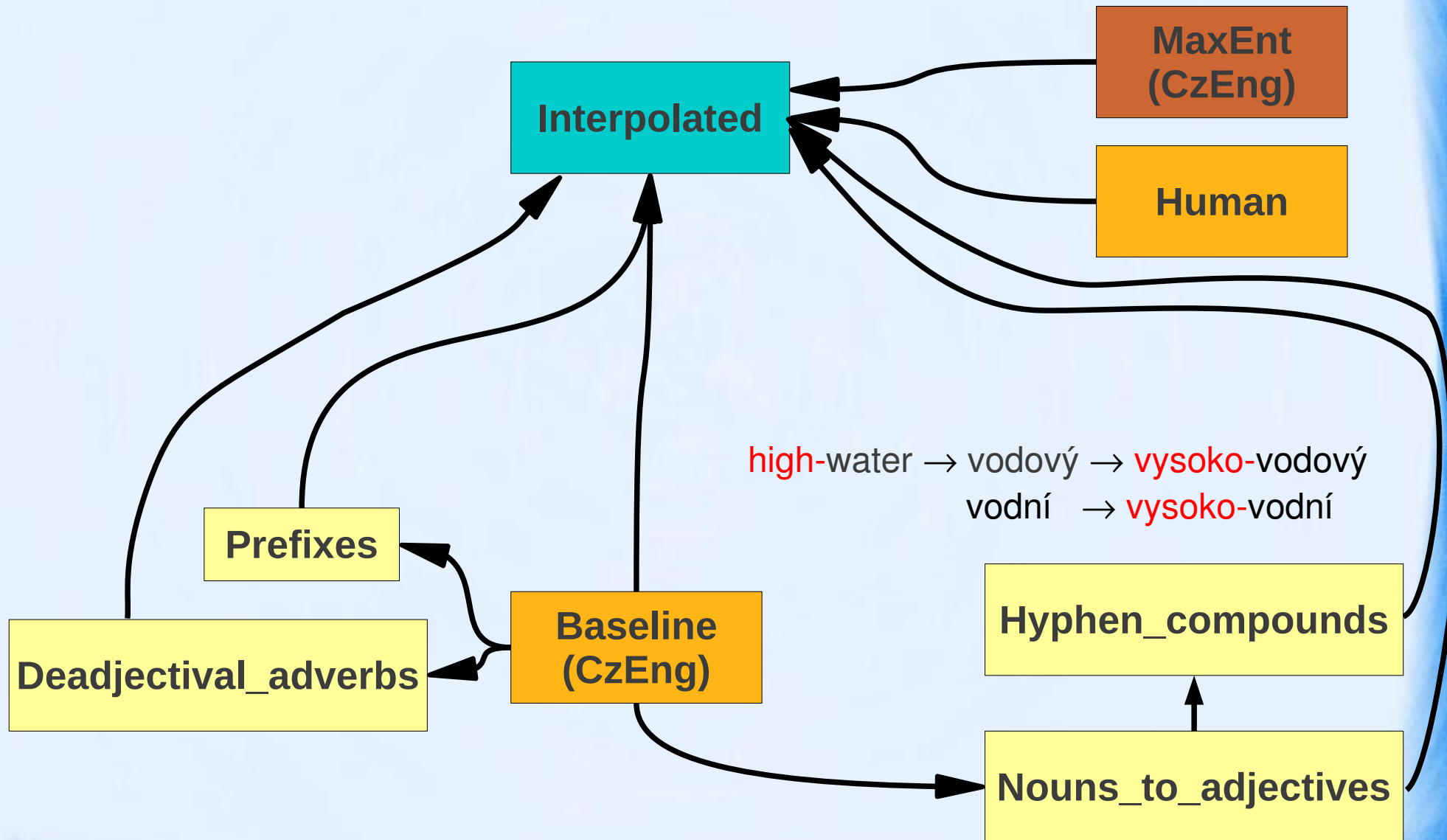


multi-core → více-jádrový
více-jádro
multi-jádrový
multi-jádro

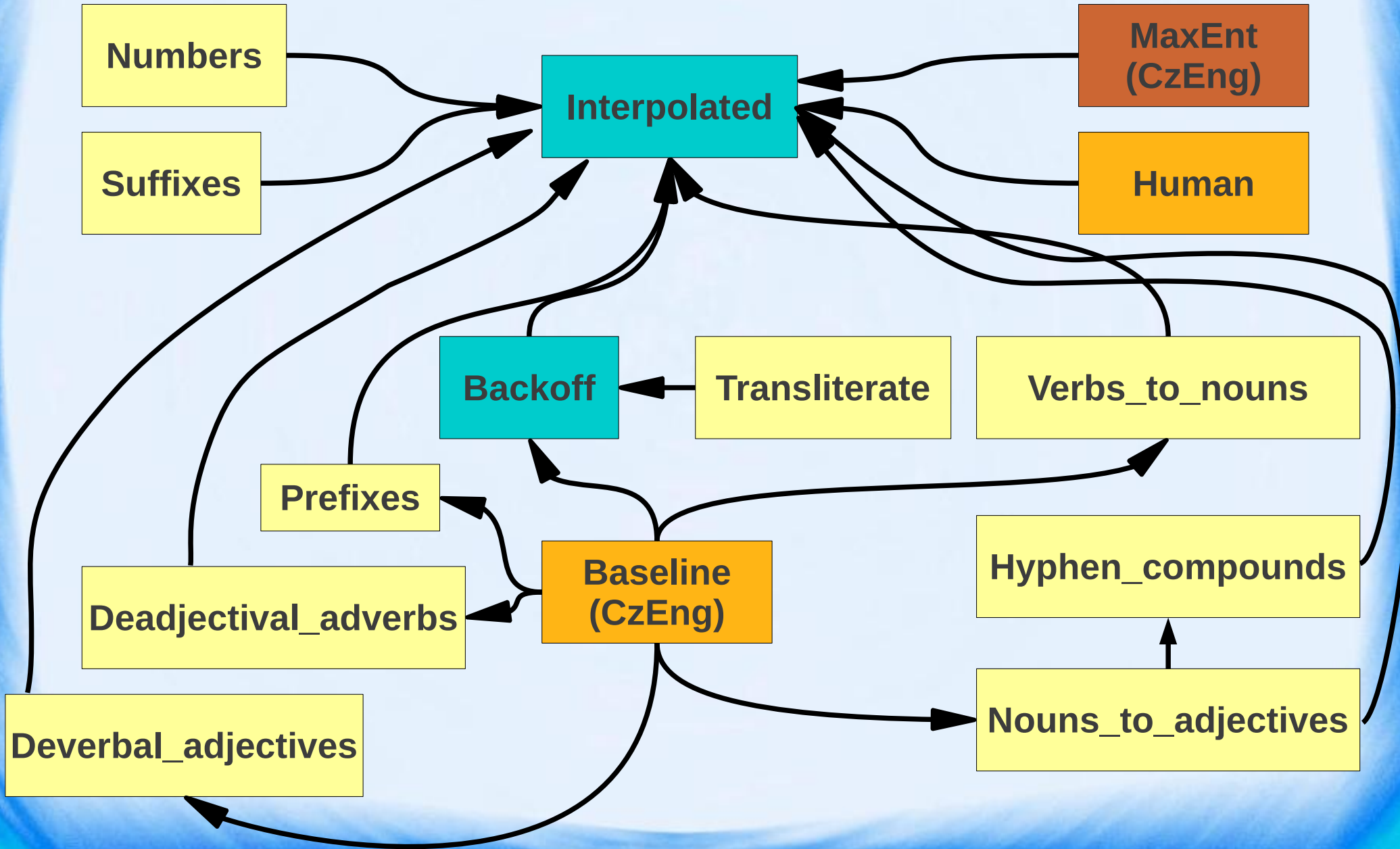
Hierarchy of lemma dictionaries



Hierarchy of lemma dictionaries



Hierarchy of lemma dictionaries



Maximum Entropy Dictionary

Baseline Dictionary

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

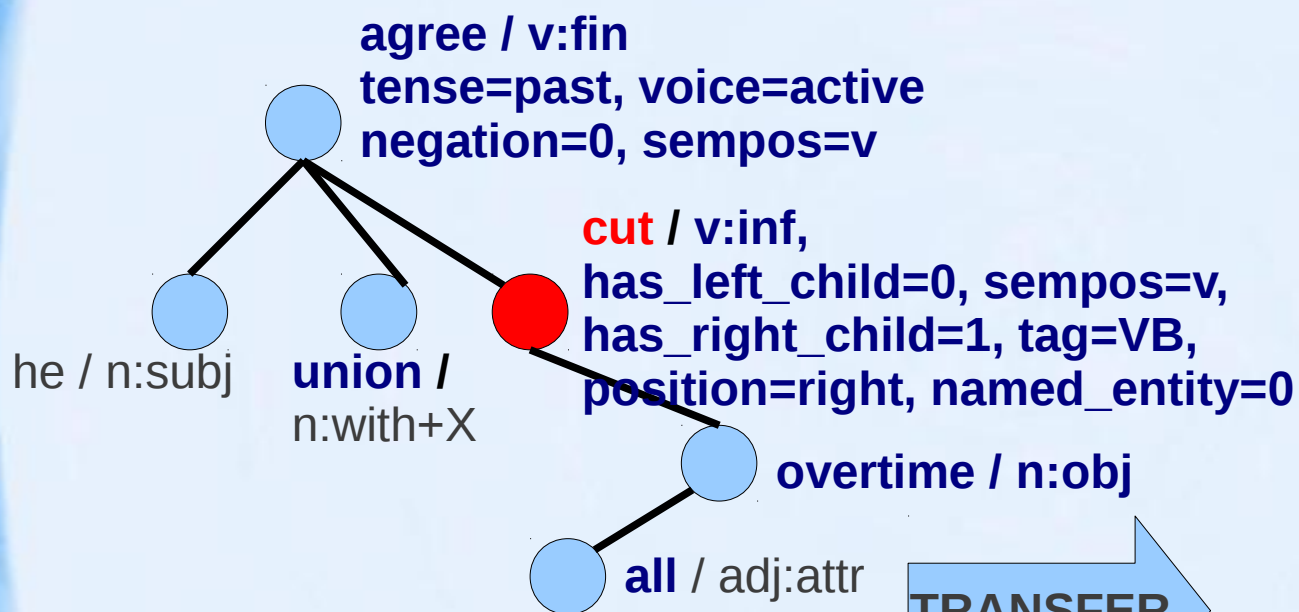
- Maximum likelihood estimates (from the training sections of CzEng 0.9)
- Pruned by thresholds on $p(x|y)$ and $p(y|x)$
- No context used
x = source lemma
y = target lemma

MaxEnt Dictionary

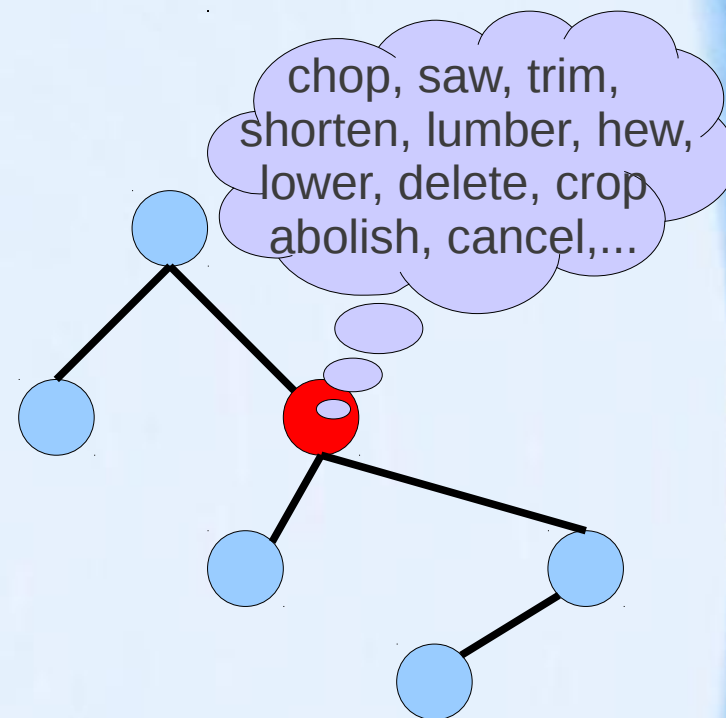
$$p(y|x) = \frac{1}{Z(x)} \exp \sum_i \lambda_i f_i(x, y)$$

- 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



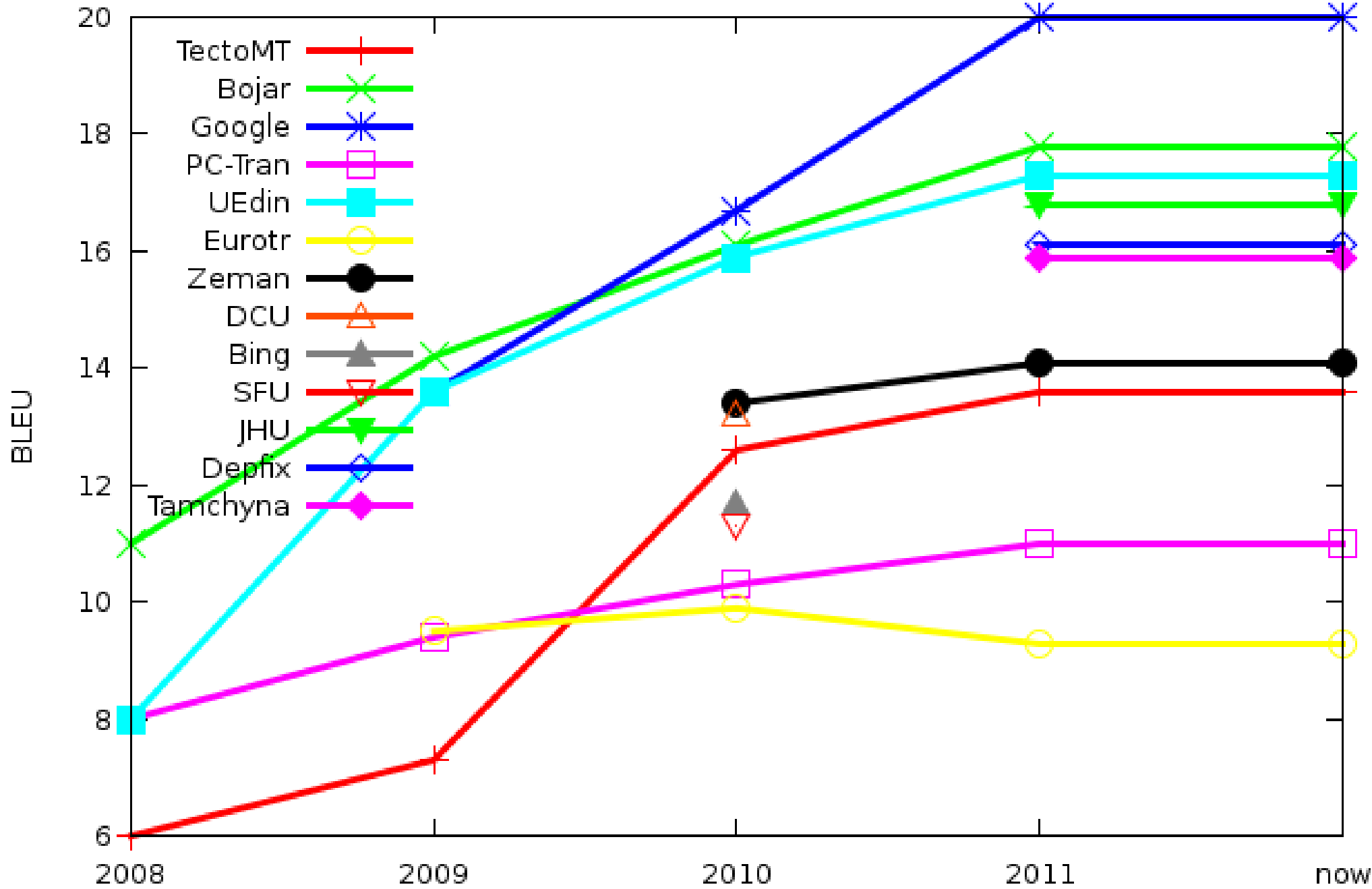
ANALYSIS

SYNTHESIS

He agreed with the unions to cut all overtime.

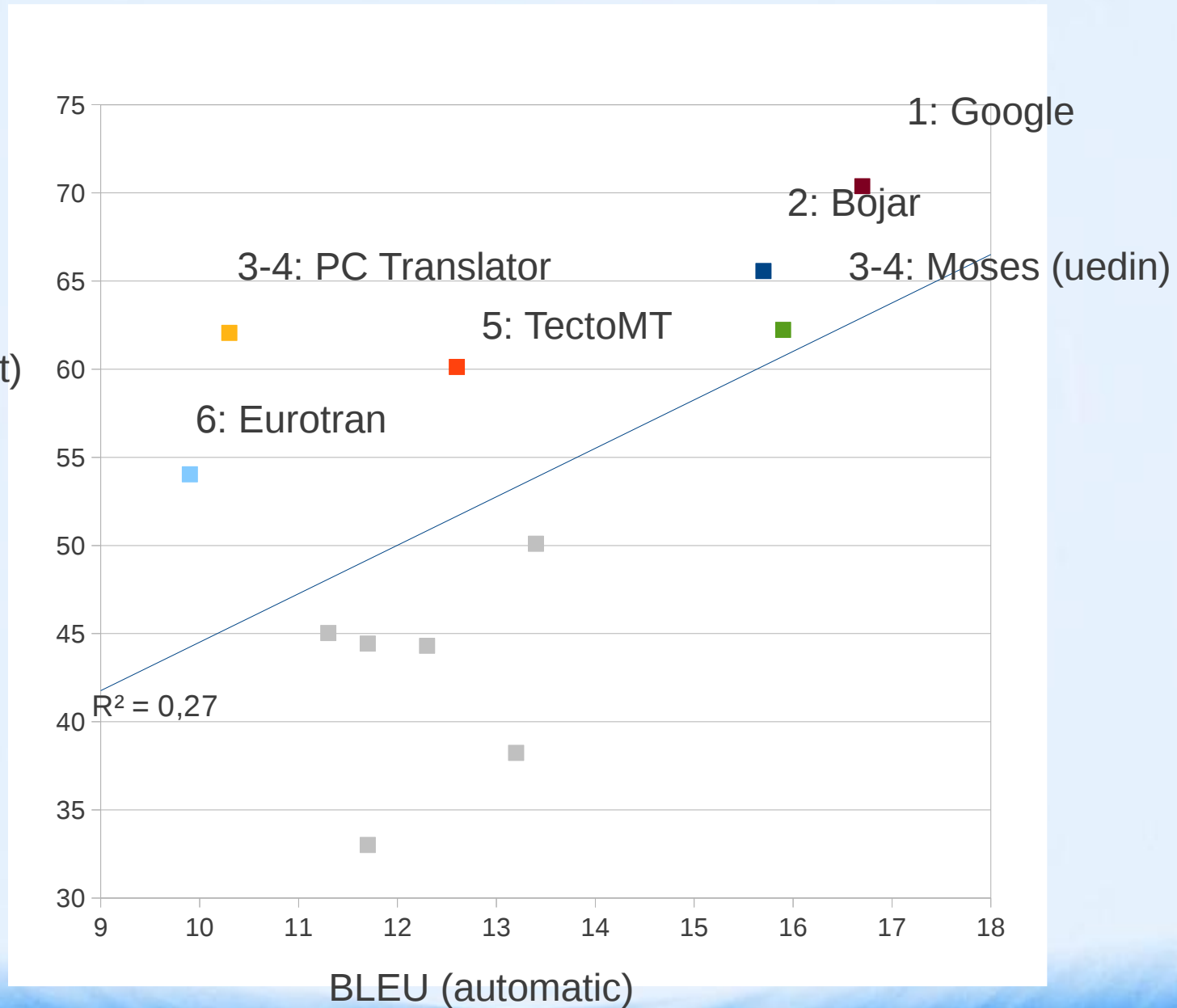
Dohodl se s odbory na zrušení všech přesčasů.

Results - BLEU

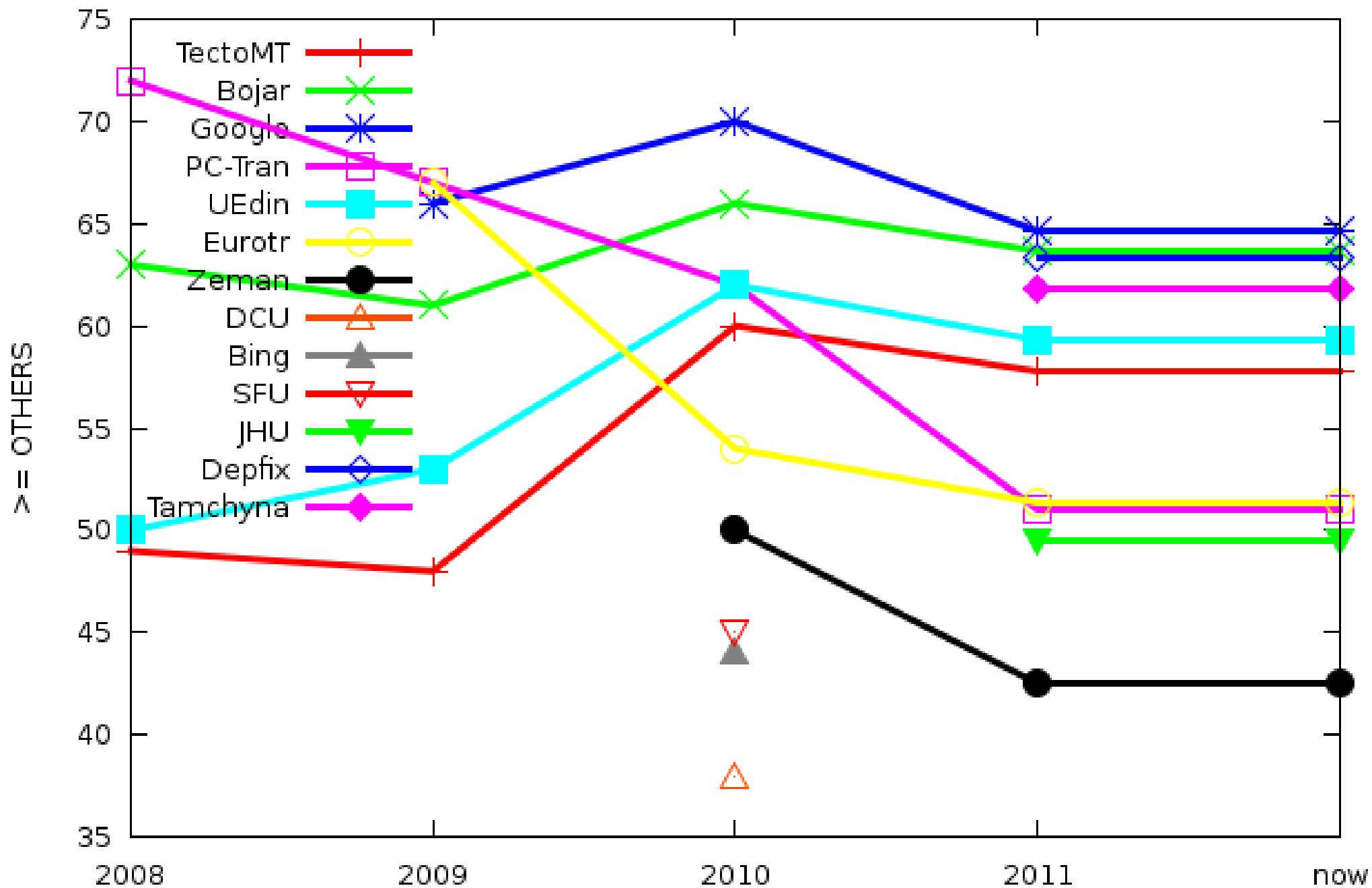


Results – BLEU vs. Ranks (WMT 2010)

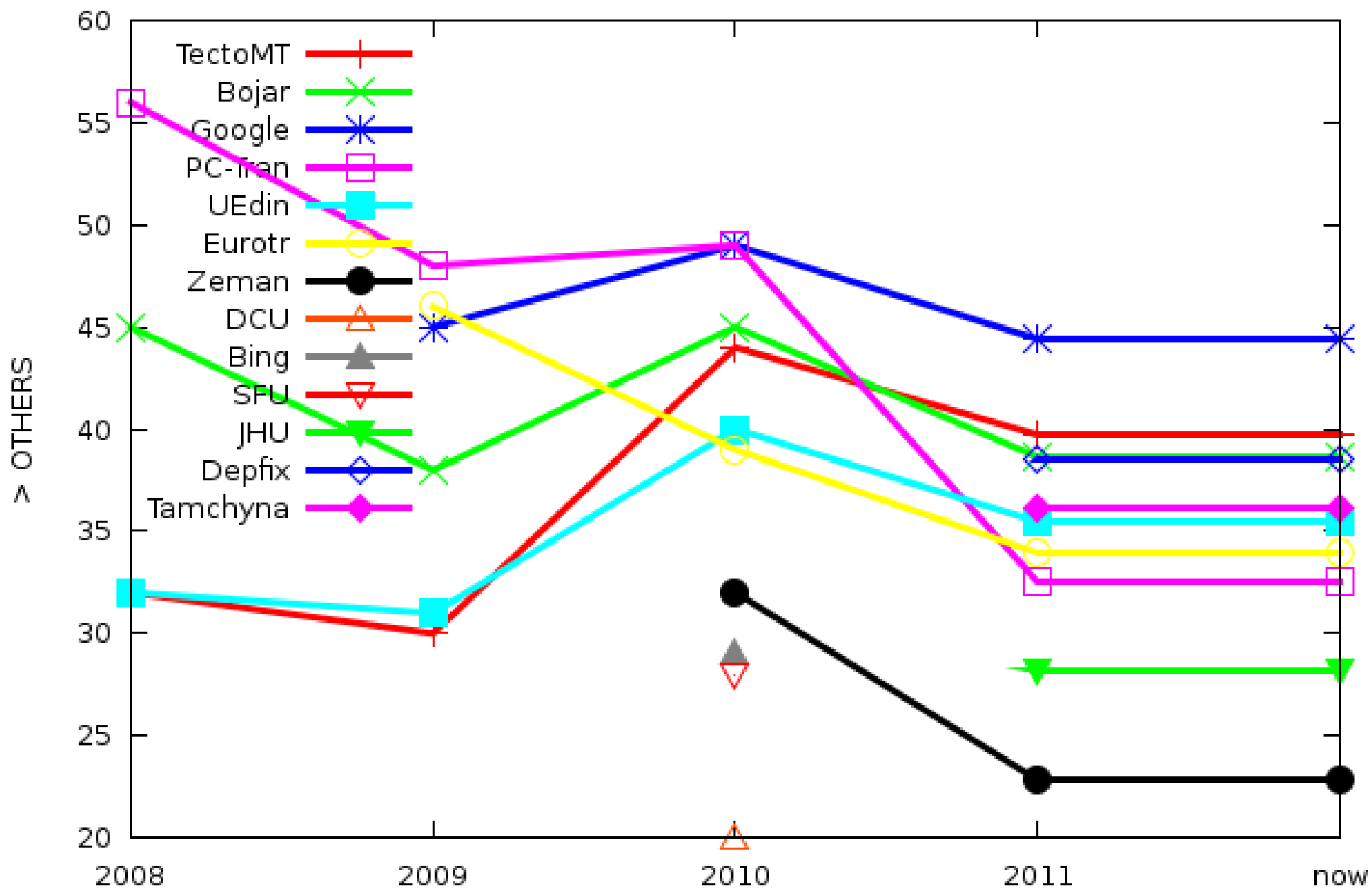
Rank
(human judgement)



Results – manual ranking (\geq others)



Results – manual ranking (> others)



Examples of Translation (2009)

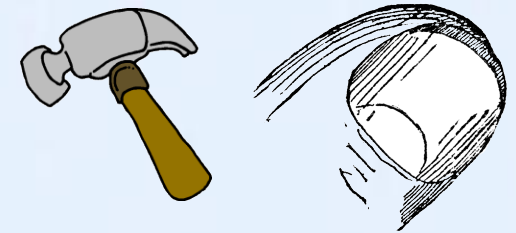
A miss by an inch
is a miss by a mile.

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.

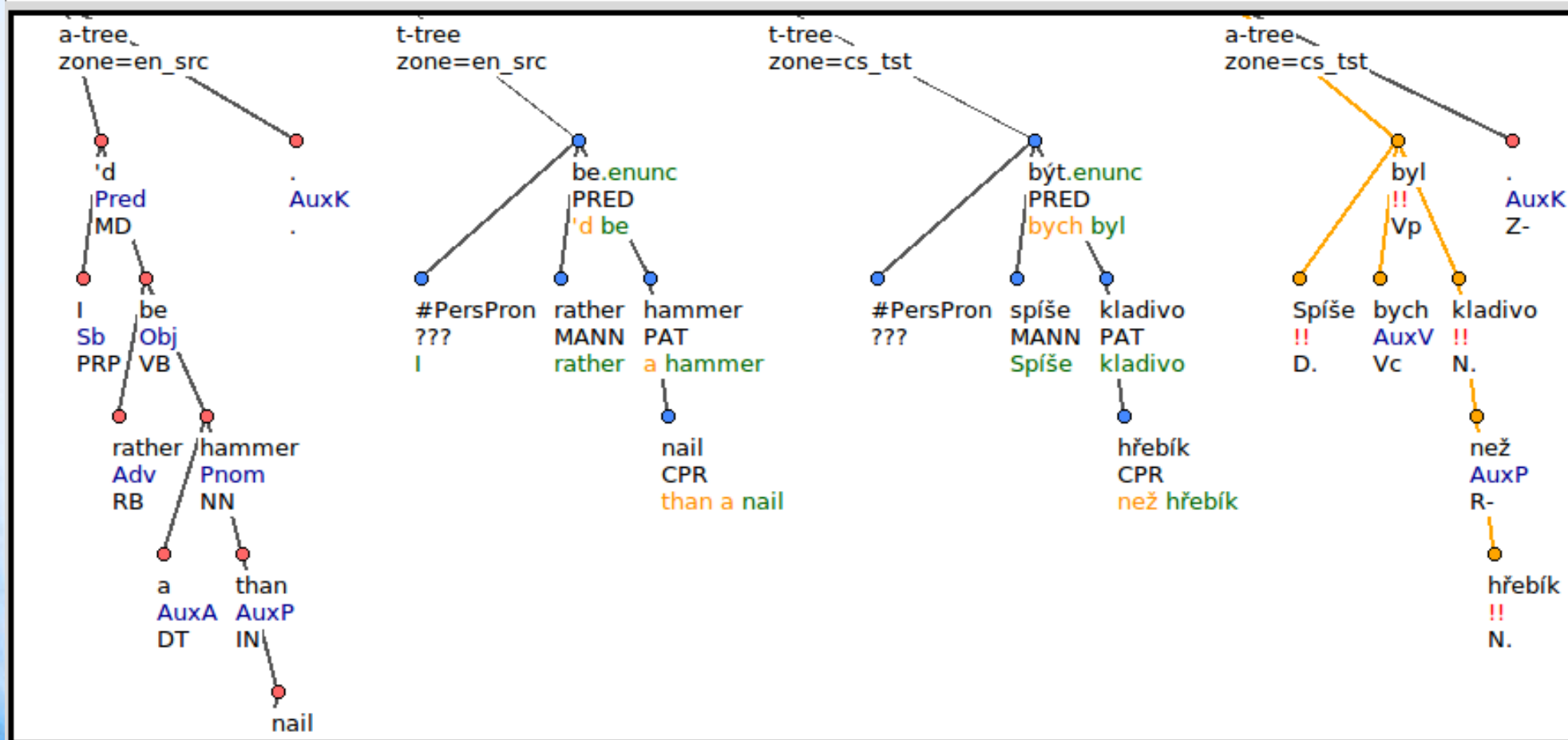
Pták v ruce je cenný
dvakrát v Bushovi.



Example of Translation (2011)

File Node Tree View Macros Setup Help

[cs_tst] Spíše bych byl kladivo než hřebík .
 [en_src] I'd rather be a hammer than a nail.



The image displays four syntax trees illustrating the translation process from English source to Czech target. The trees are labeled 'a-tree zone=en_src', 't-tree zone=en_src', 't-tree zone=cs_tst', and 'a-tree zone=cs_tst'. Each tree shows the hierarchical structure of the sentence with grammatical category labels and word alignments.

- a-tree zone=en_src:** Rooted at 'd (Pred MD), with children 'I (Sb PRP) and 'be (Obj VB). 'be' branches into 'rather (Adv RB) and 'hammer (Pnom NN). 'hammer' branches into 'a (AuxA DT) and 'than (AuxP IN). 'than' branches into 'nail (CPR).
- t-tree zone=en_src:** Rooted at 'be.enunc (PRED), with children '#PersPron ??? (I) and 'rather (MANN). 'rather' branches into 'hammer (PAT). 'hammer' branches into 'a (CPR) and 'hammer (CPR).
- t-tree zone=cs_tst:** Rooted at 'být.enunc (PRED), with children '#PersPron ??? and 'spíše (MANN). 'spíše' branches into 'kladivo (PAT). 'kladivo' branches into 'hřebík (CPR).
- a-tree zone=cs_tst:** Rooted at 'byl (Vp), with children 'Spíše (D), 'bych (AuxV Vc), and 'kladivo (N). 'kladivo' branches into 'než (AuxP R-). 'než' branches into 'hřebík (N).

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
...	

Cooperation is welcome

- Exploit English a-layer or t-layer for your project
- Try your dependency parser instead of MST (extrinsic evaluation)
- Try your machine learning (CRF, SVM, NN,...) instead of our MaxEnt dictionary

Thank you

