

# **TectoMT: Machine Translation System**

**Martin Popel**

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



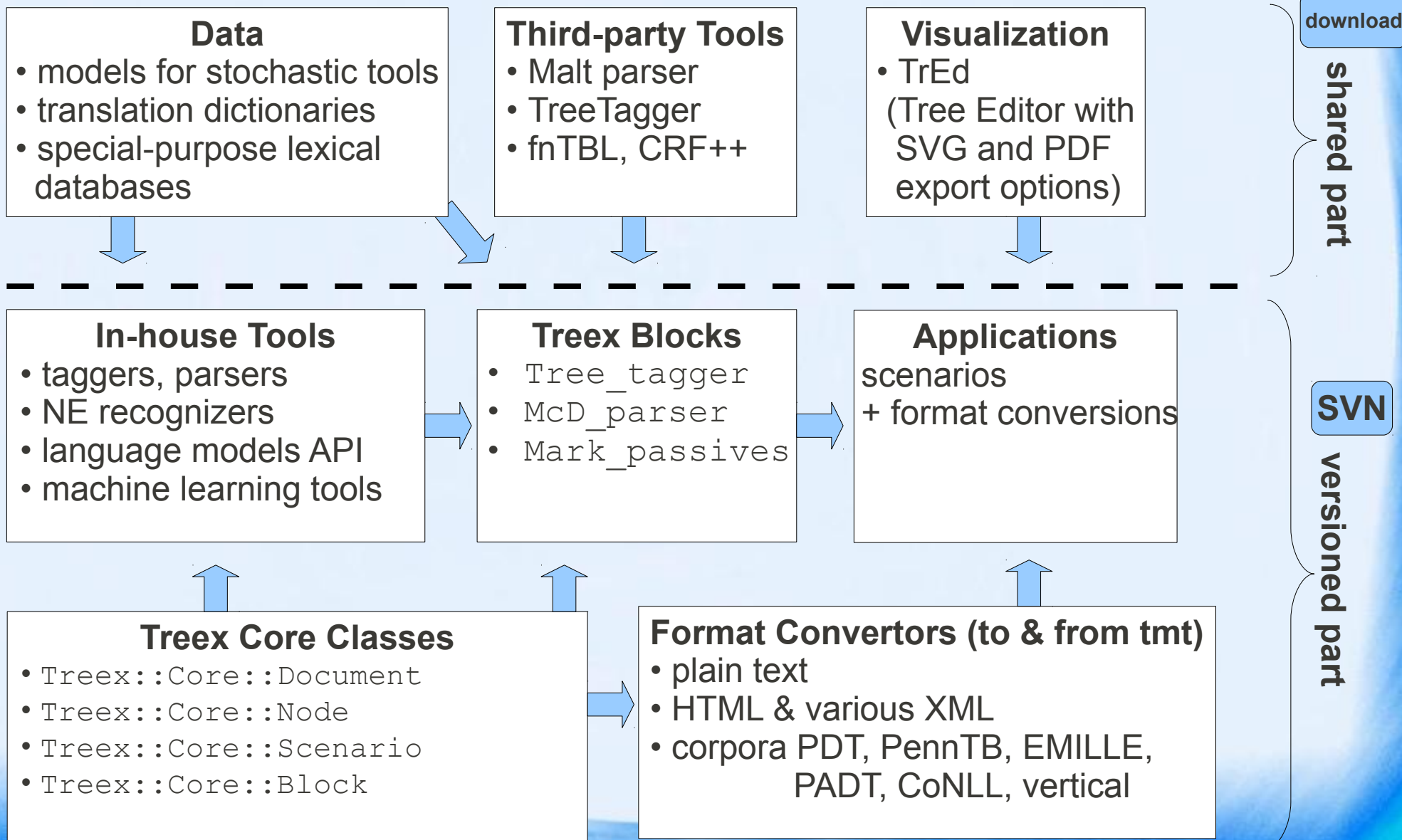
Universität des Saarlandes  
November 4<sup>th</sup> 2010, Saarbrücken, Germany

# Outline

- NLP framework Treex
- Demo translation step by step
- Annotation of translation errors
- Improvements
  - Parsing parentheses
  - Hidden Markov Tree Models (HMTM)
  - Combining dictionaries
  - Maximum Entropy dictionary
- Results – three metrics of translation quality

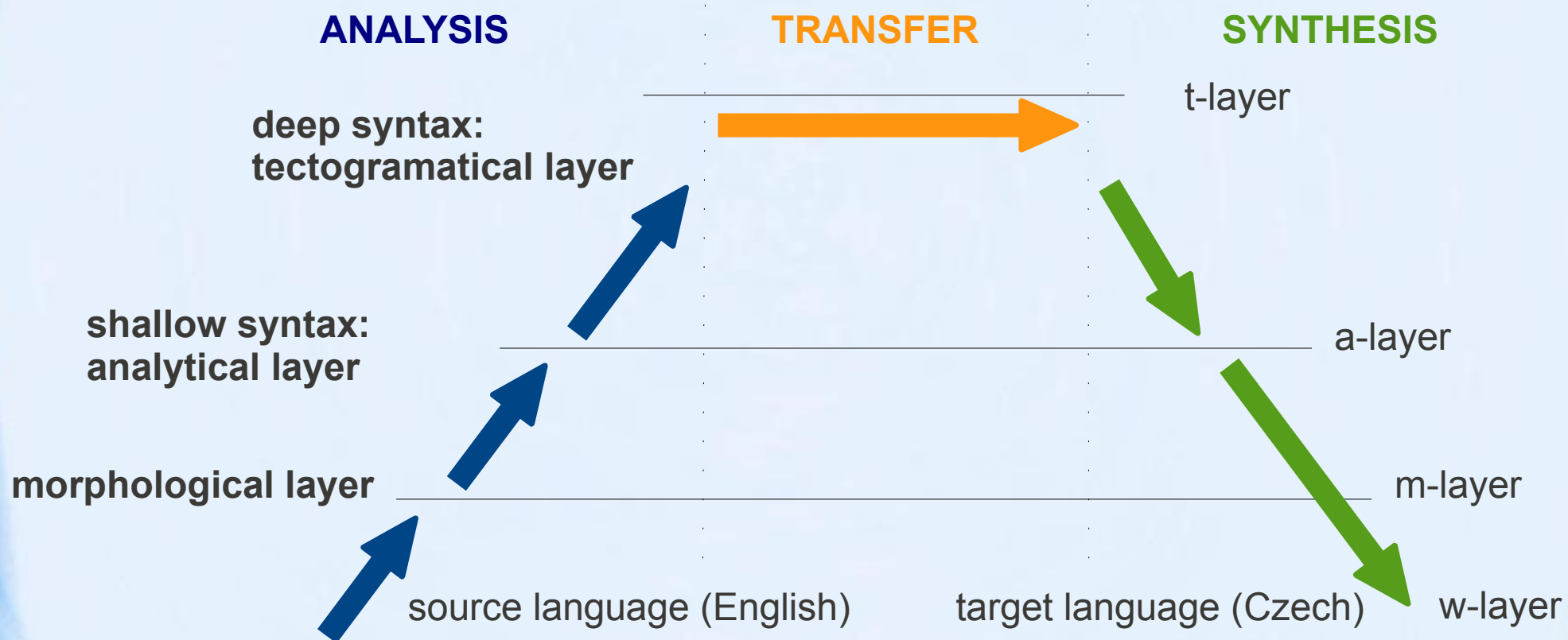
# Treex – framework for NLP

modular, open source, Perl, Linux, OOP-style



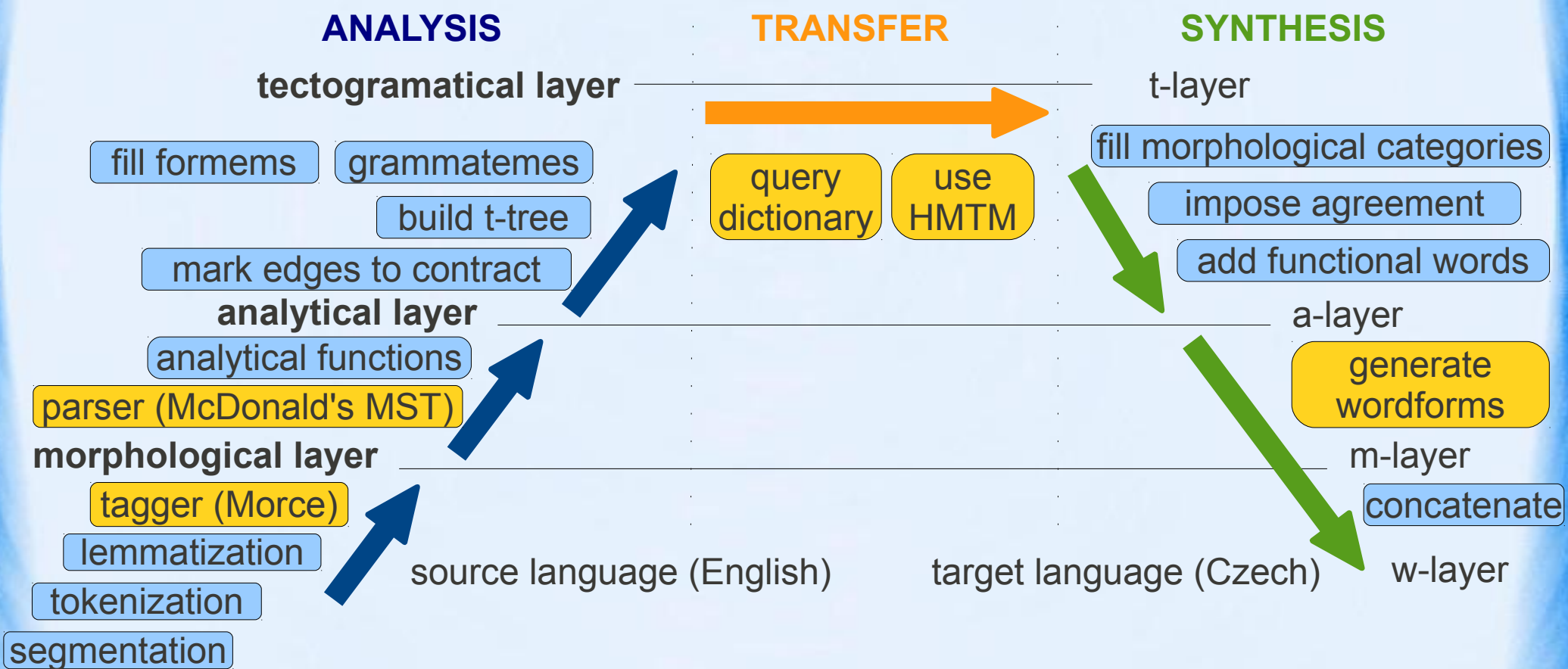
# Translation scheme

transfer over the tectogrammatical layer



# Translation scheme

rule based & statistical blocks



# Demo Translation – Analysis

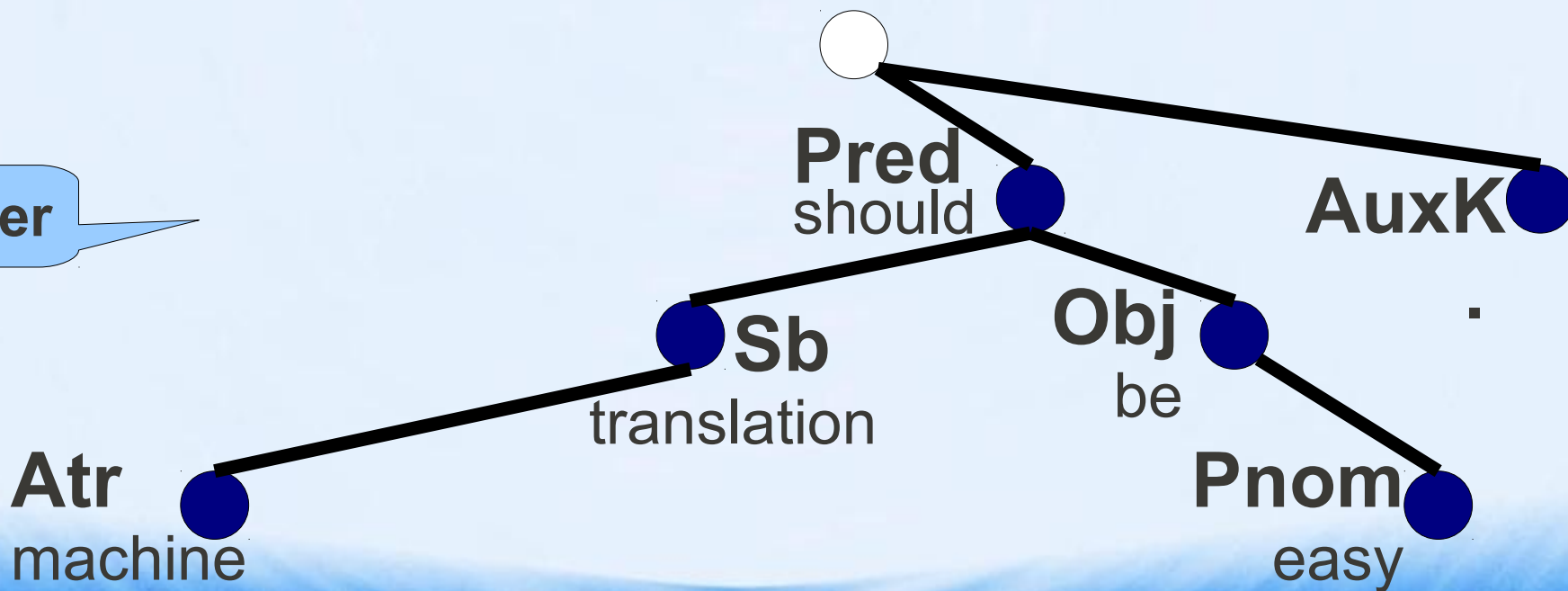
raw text

Machine translation should be easy.

m-layer

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

a-layer



# Demo Translation – Analysis

raw text

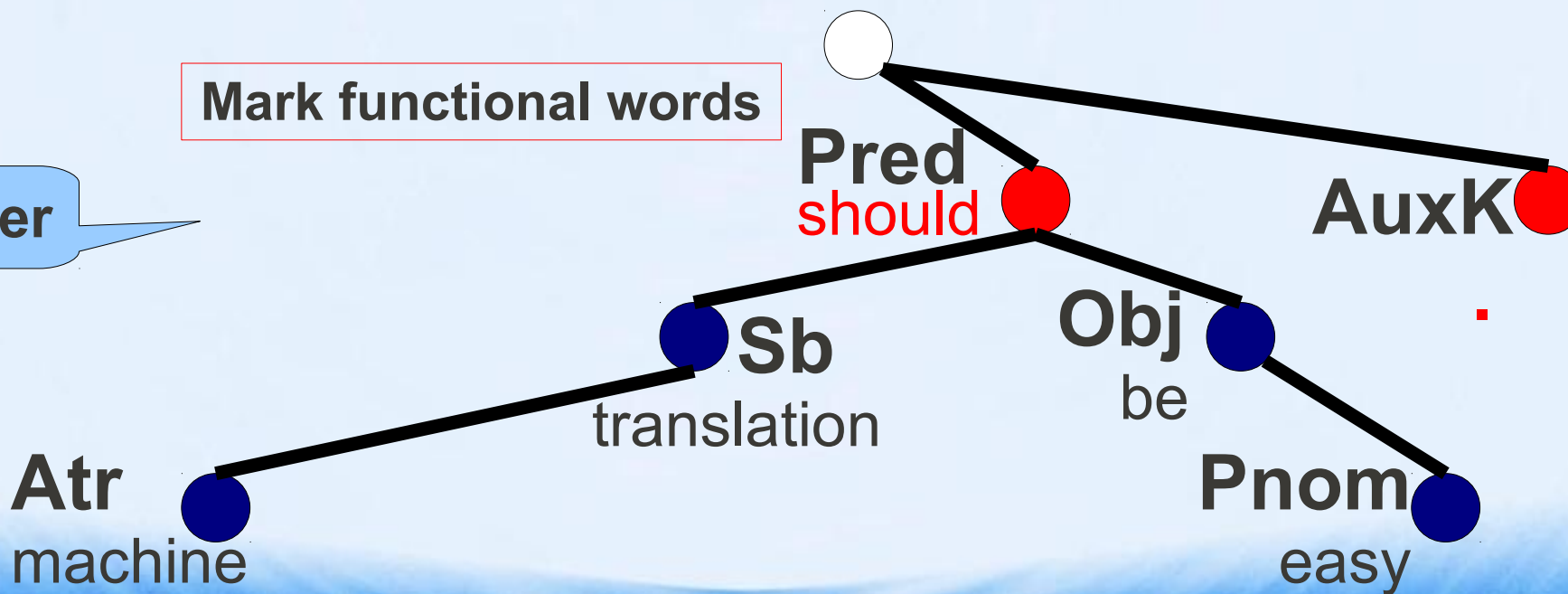
Machine translation should be easy.

m-layer

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

Mark functional words

a-layer



# Demo Translation – Analysis

raw text

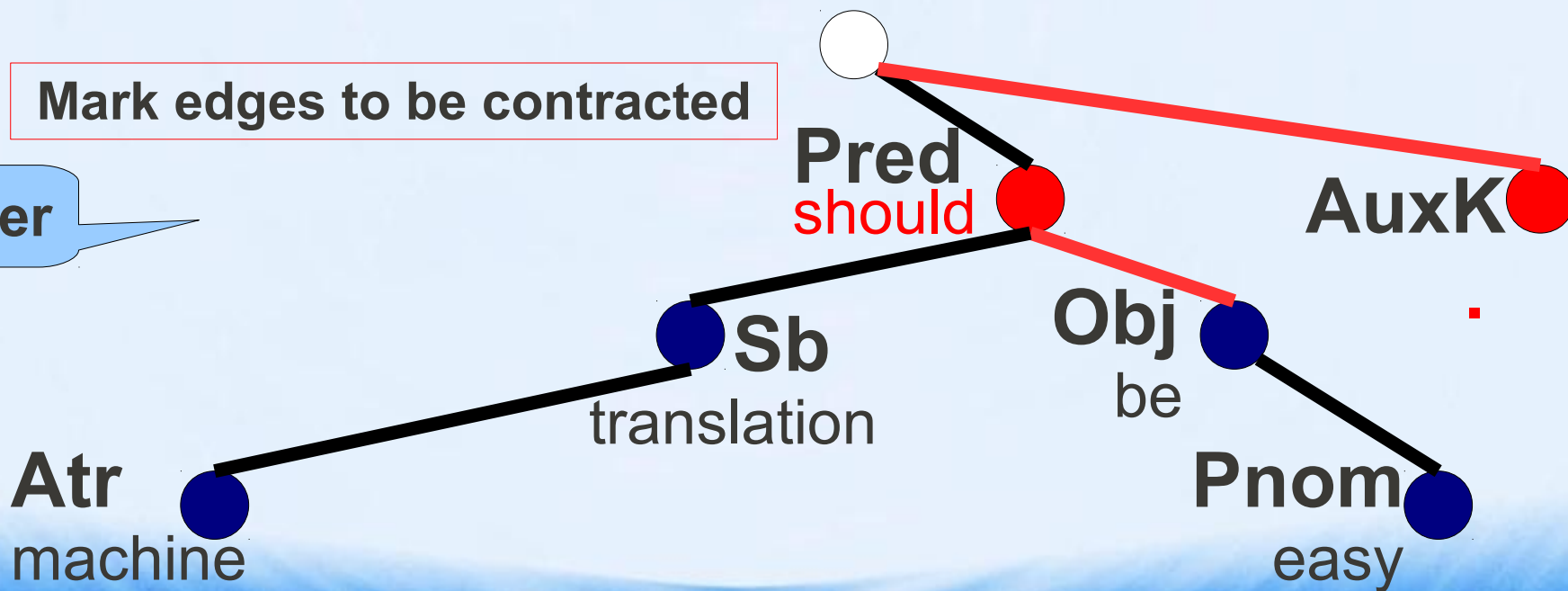
Machine translation should be easy.

m-layer

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

Mark edges to be contracted

a-layer





# Demo Translation – Analysis

raw text

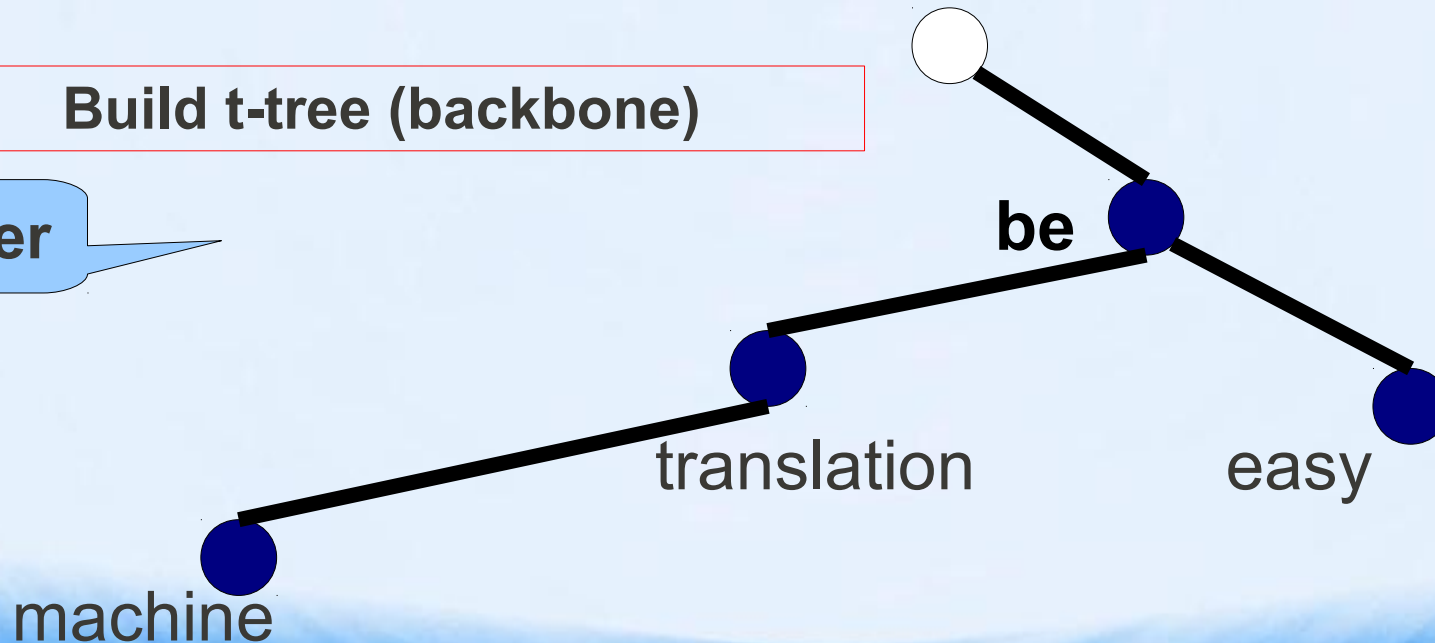
Machine translation should be easy.

m-layer

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

Build t-tree (backbone)

t-layer



# Demo Translation – Analysis

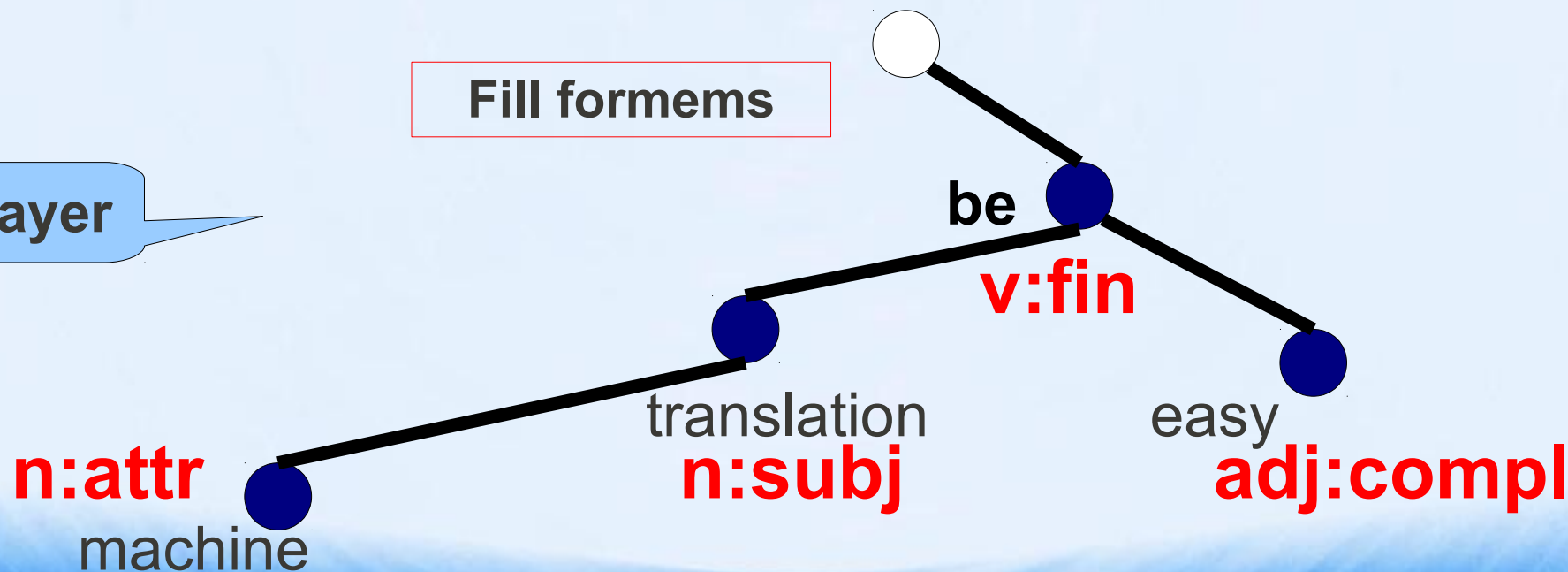
raw text

Machine translation should be easy.

m-layer

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

t-layer



# Demo Translation – Analysis

raw text

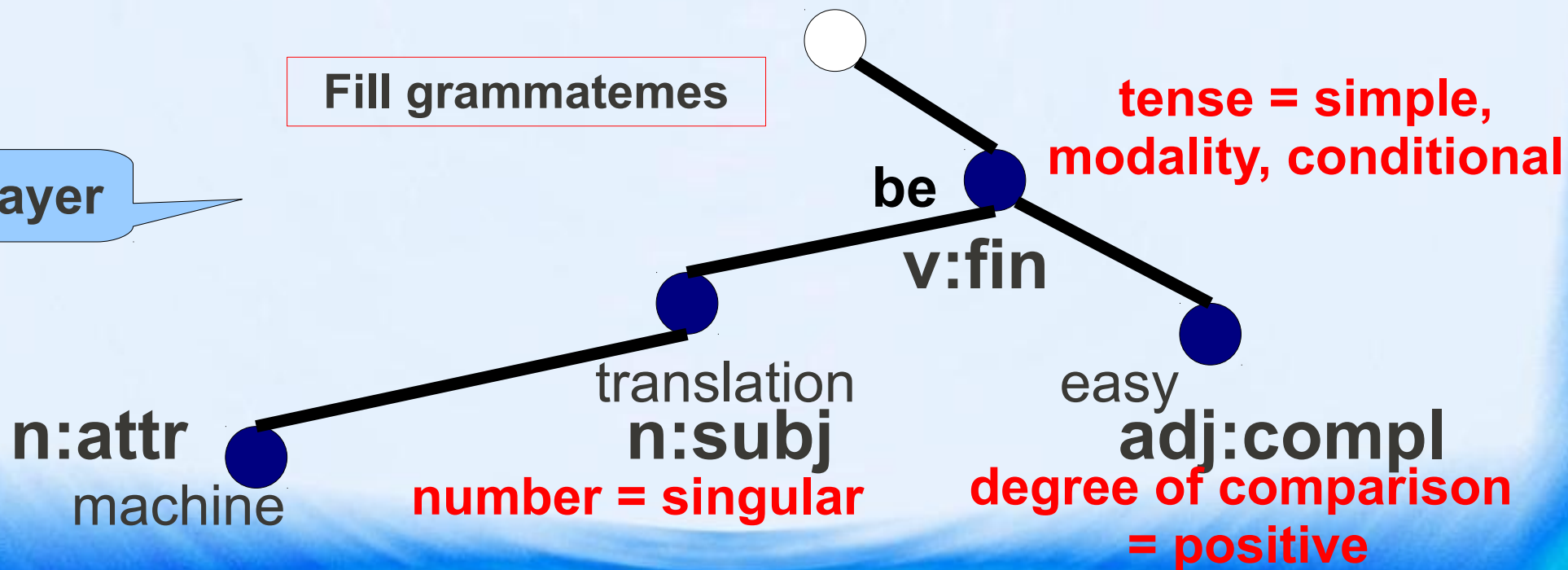
Machine translation should be easy.

m-layer

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

Fill grammatememes

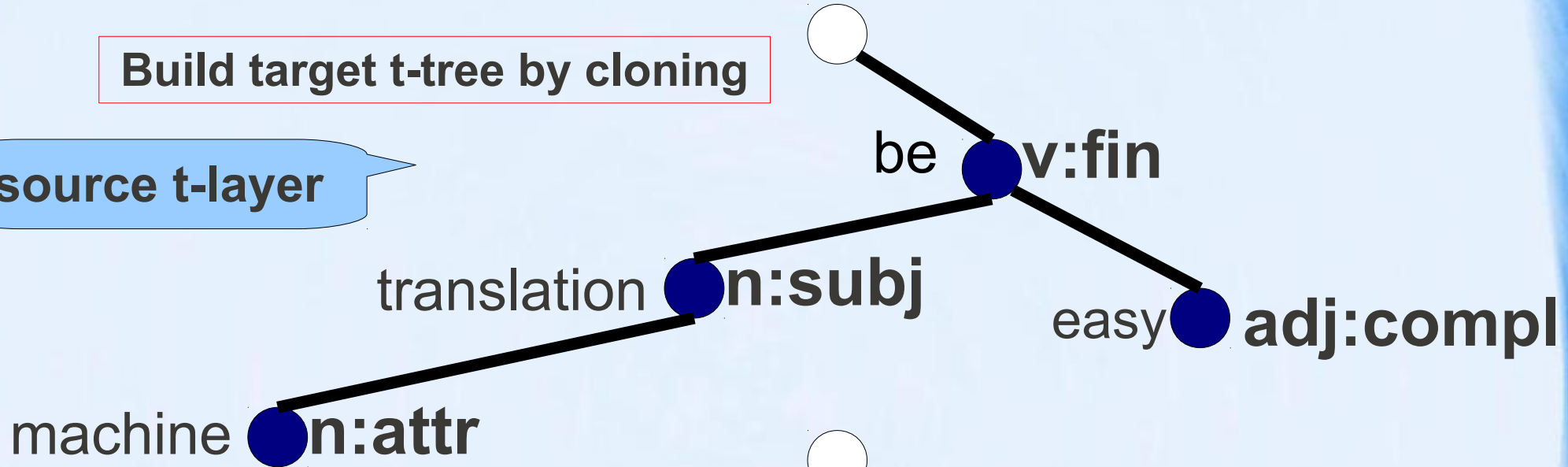
t-layer



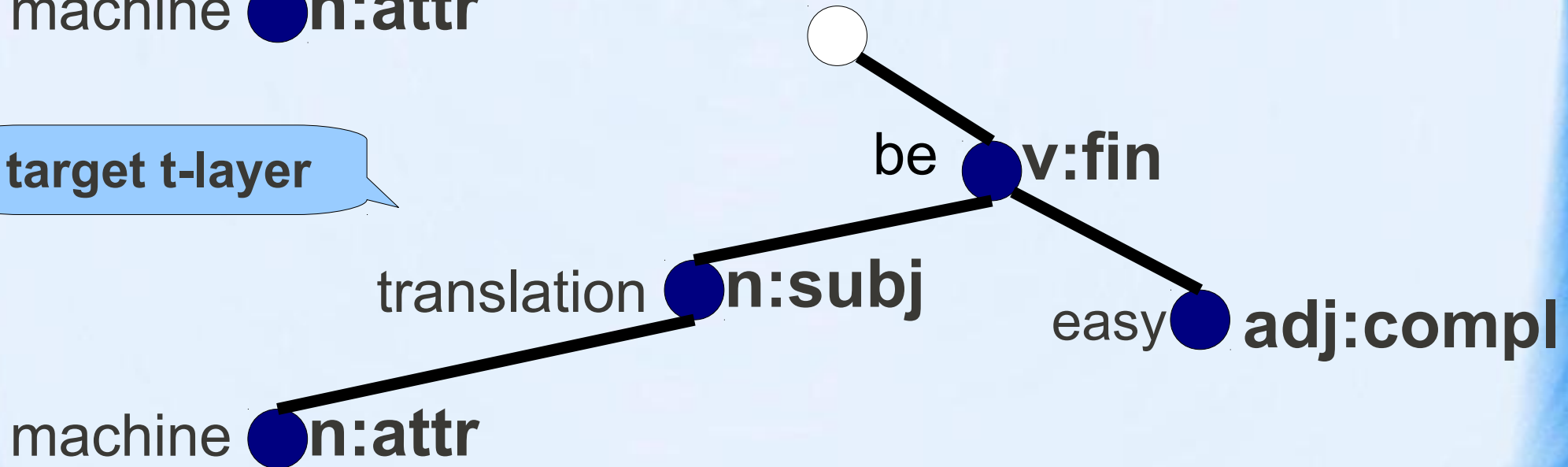
# Demo Translation – Transfer

Build target t-tree by cloning

source t-layer



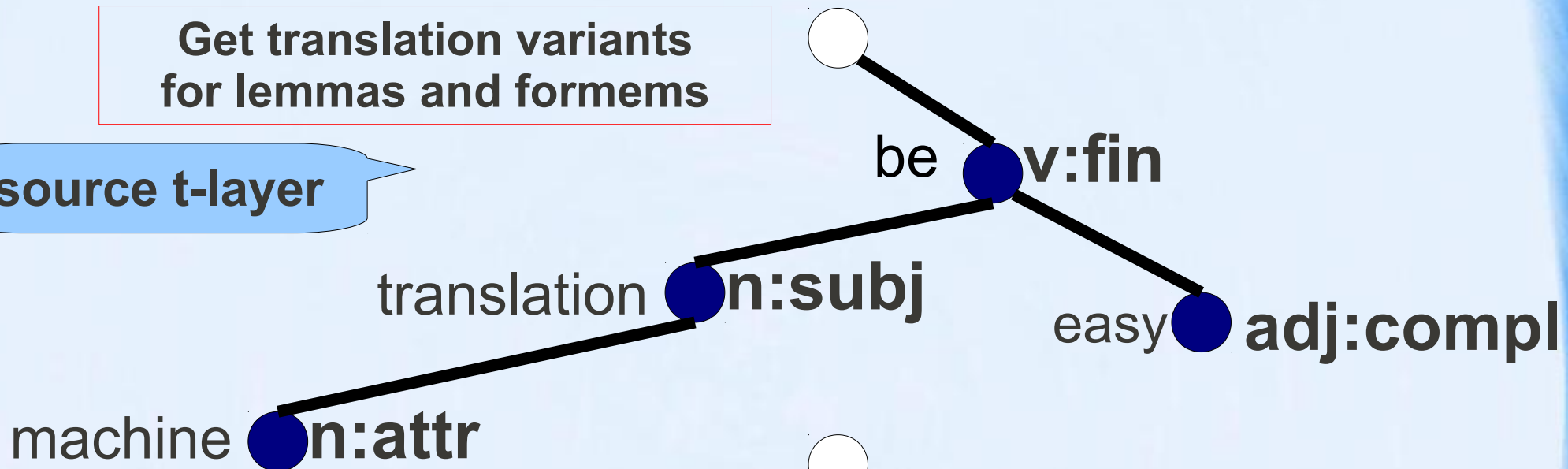
target t-layer



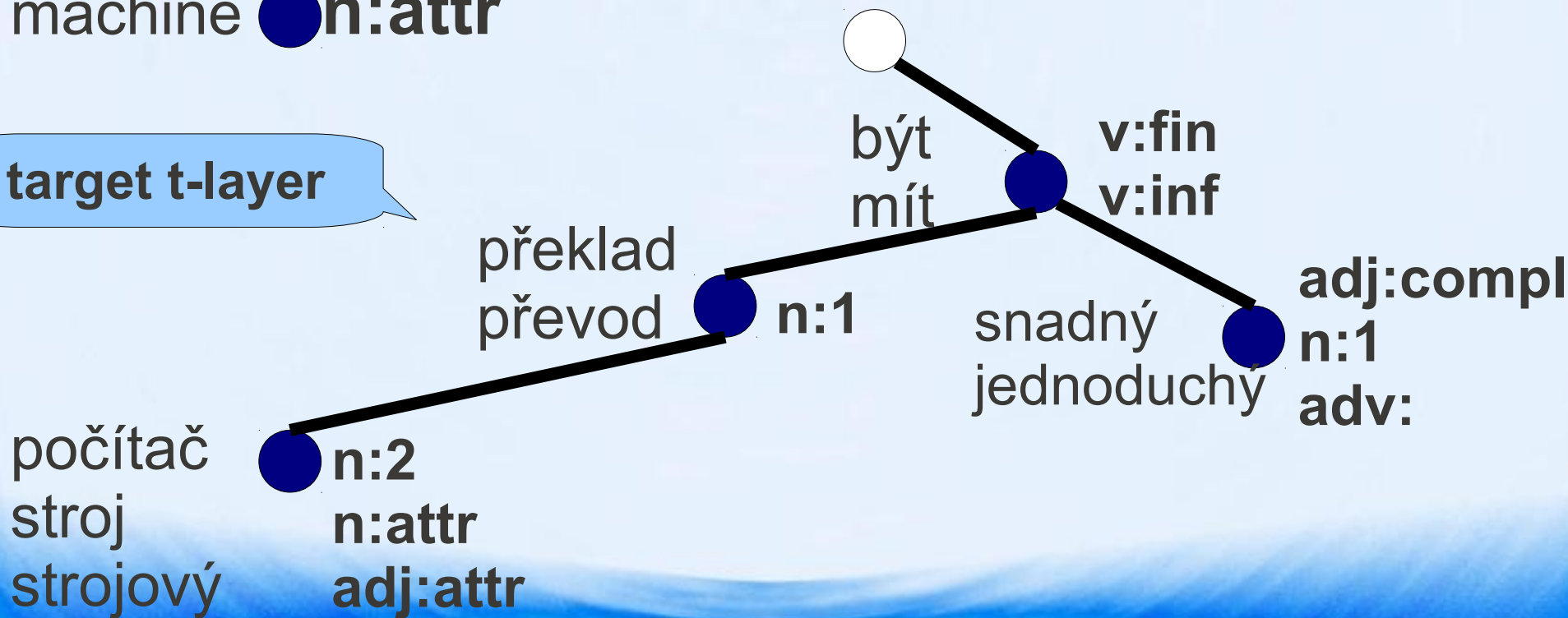
# Demo Translation – Transfer

Get translation variants for lemmas and formems

source t-layer



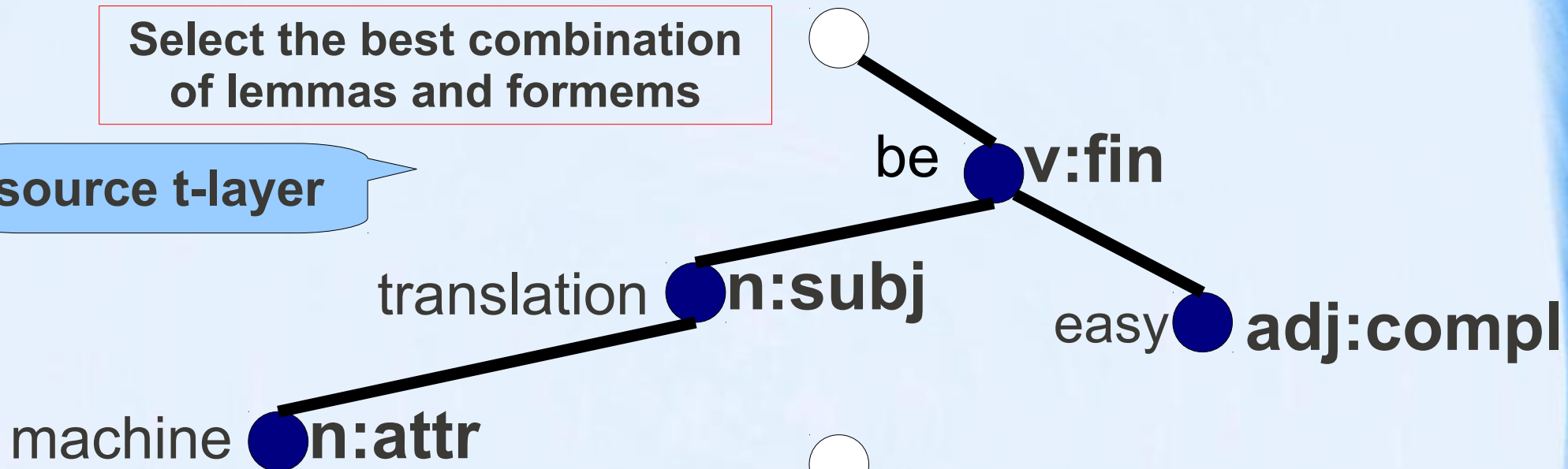
target t-layer



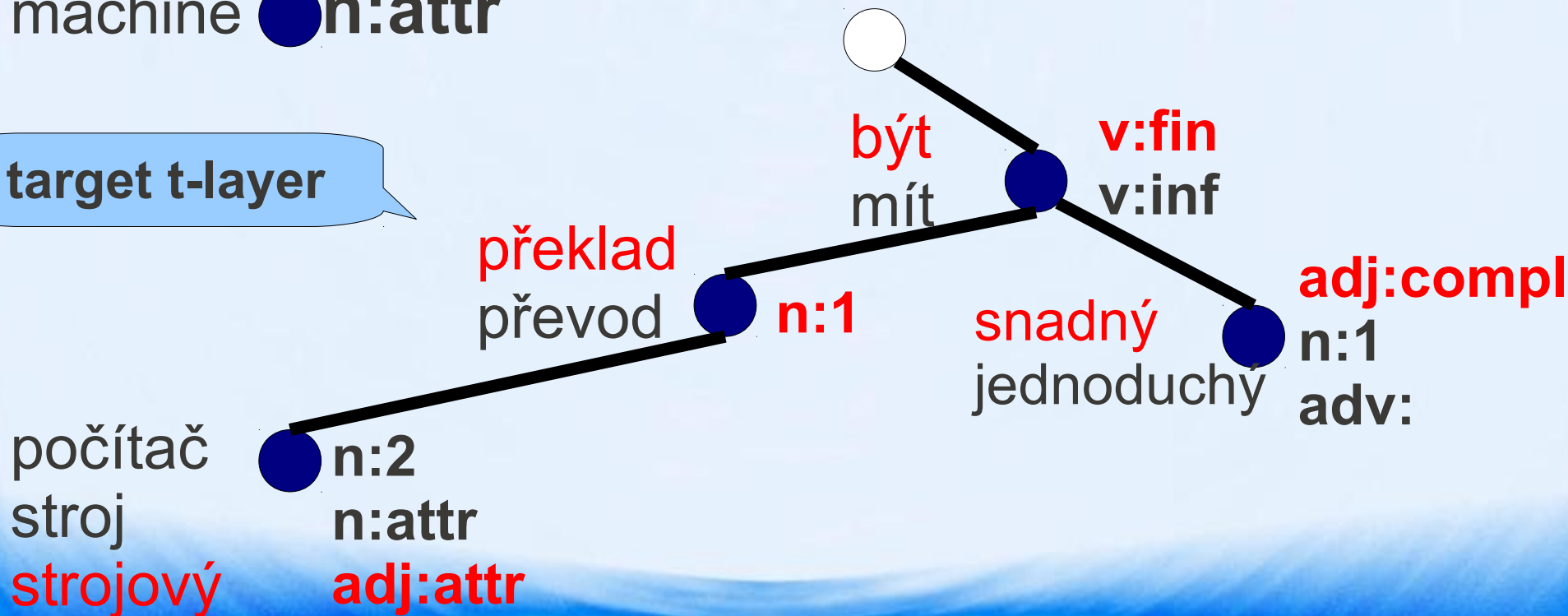
# Demo Translation – Transfer

Select the best combination of lemmas and formems

source t-layer



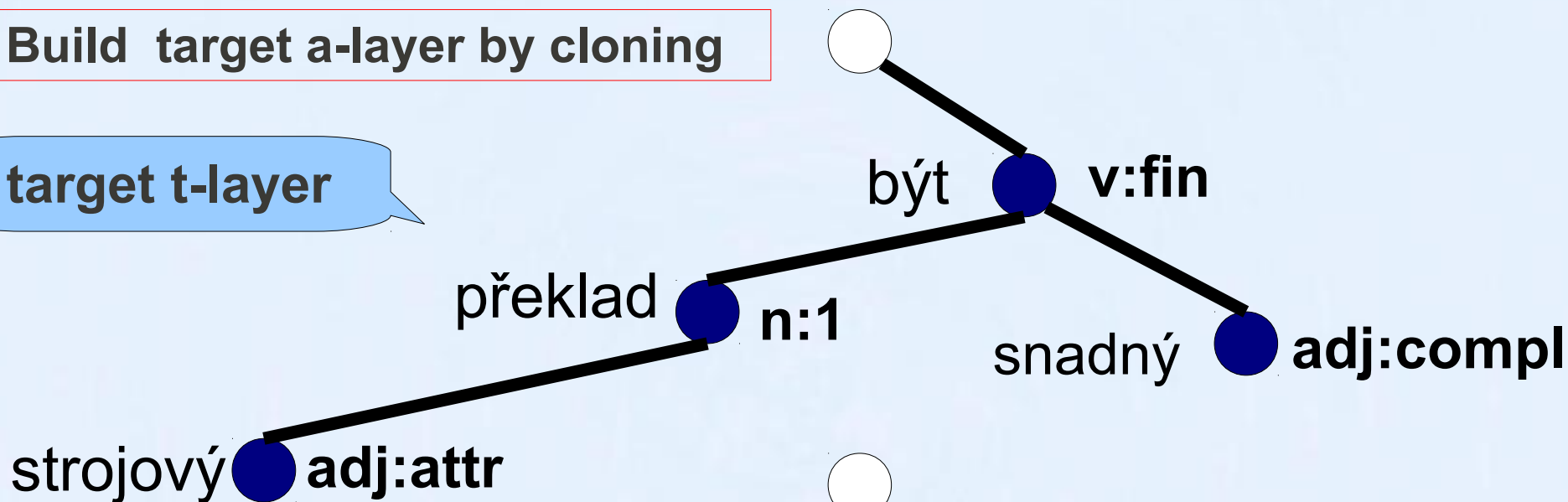
target t-layer



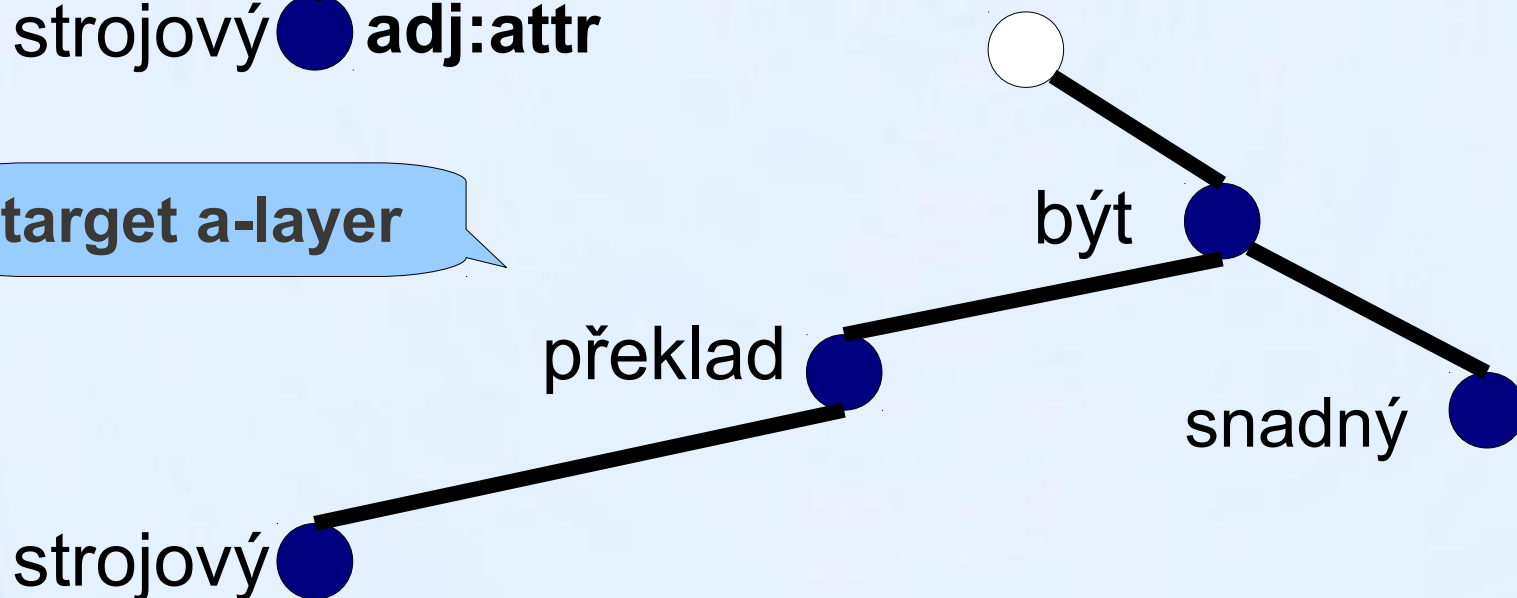
# Demo Translation – Synthesis

Build target a-layer by cloning

target t-layer



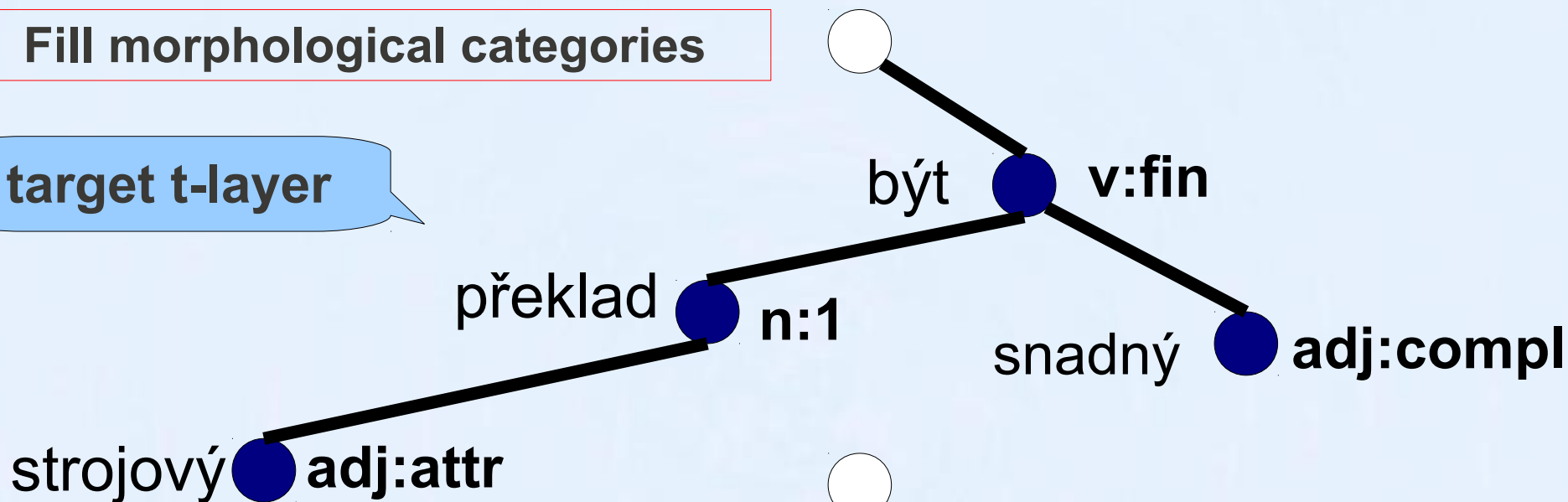
target a-layer



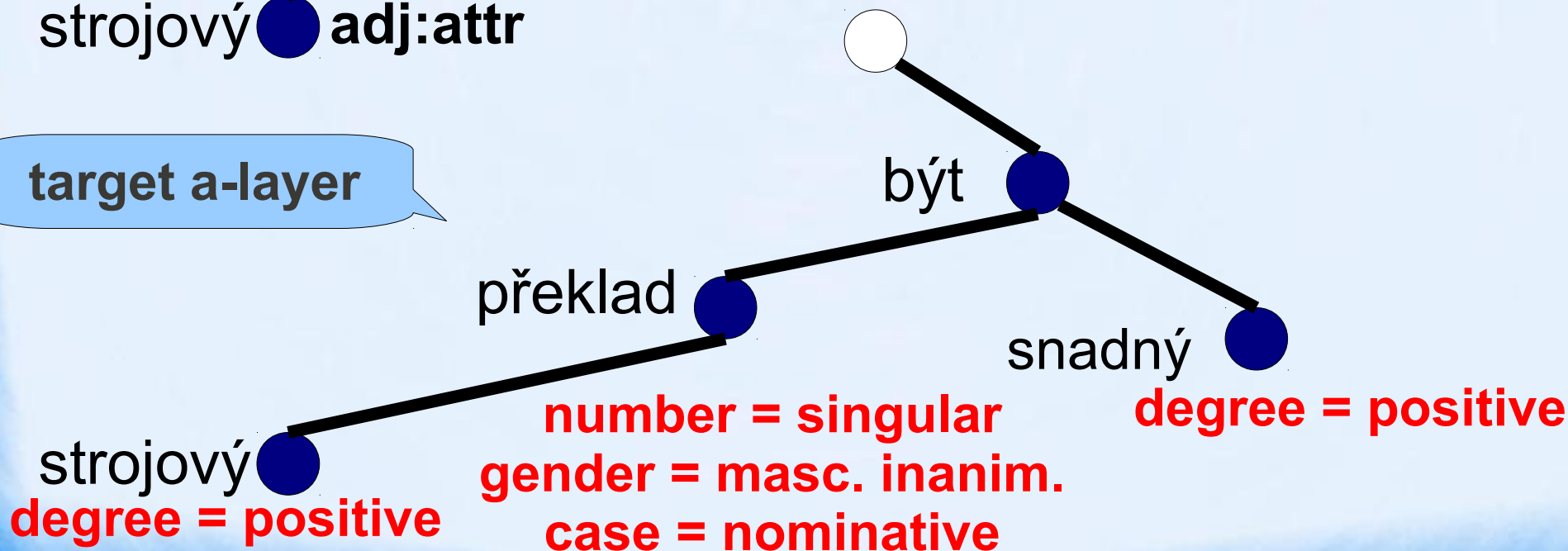
# Demo Translation – Synthesis

Fill morphological categories

target t-layer



target a-layer

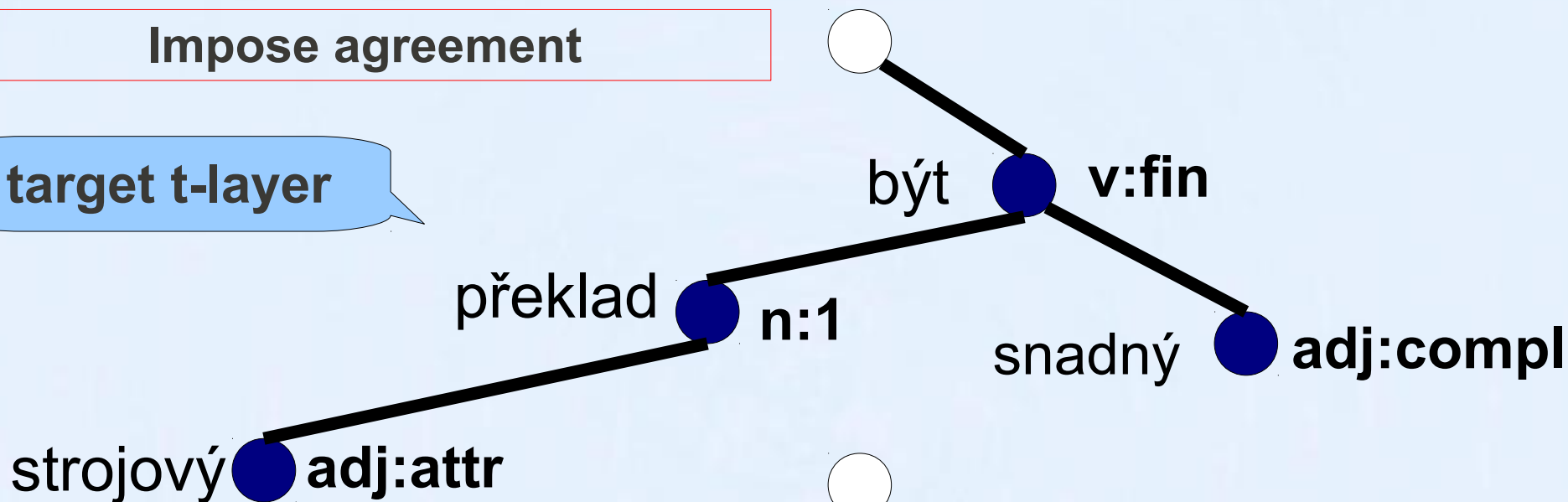




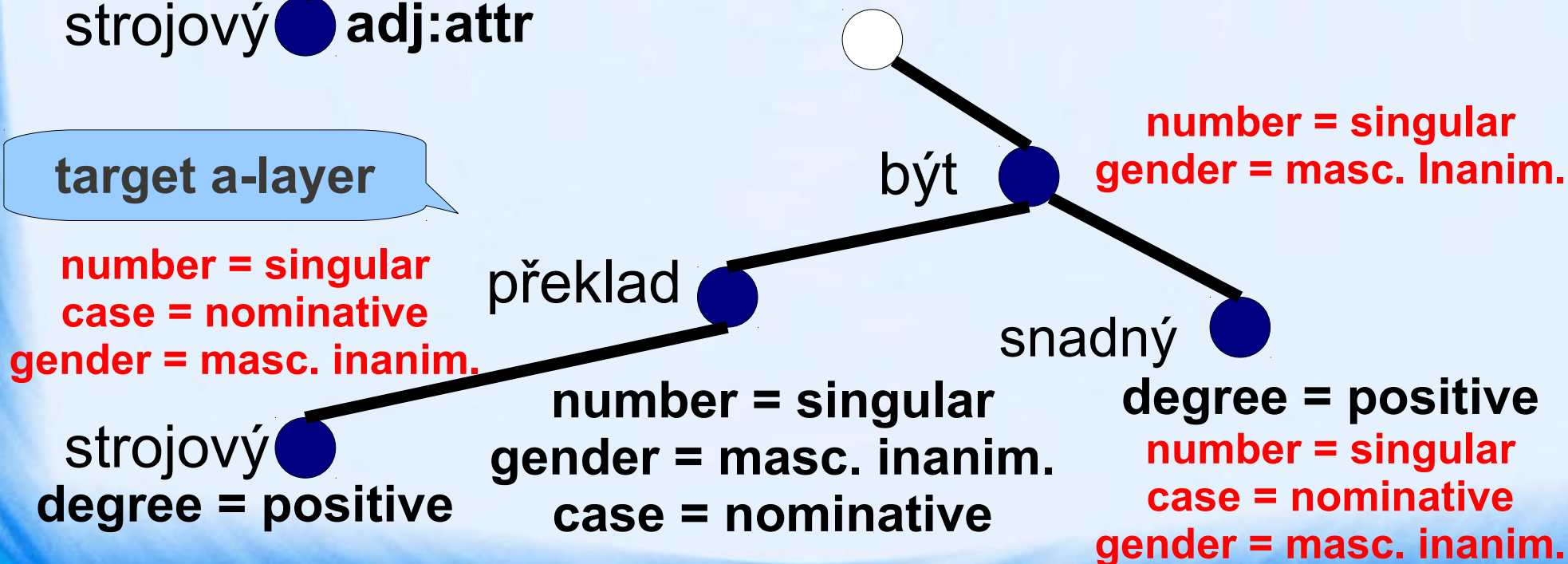
# Demo Translation – Synthesis

Impose agreement

target t-layer



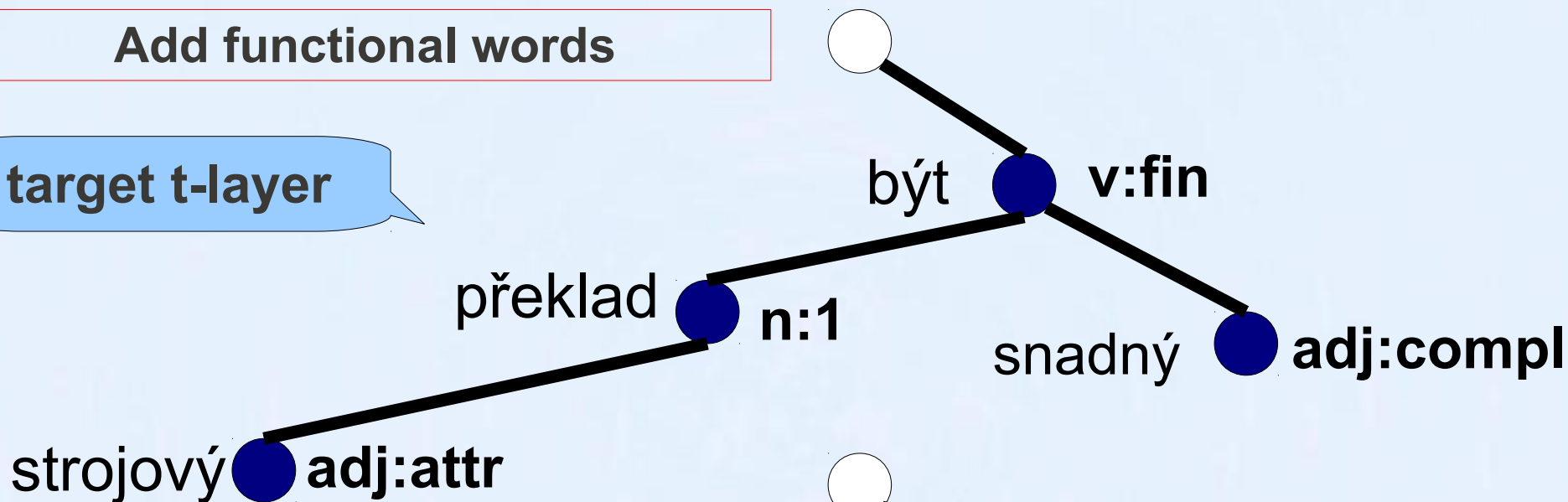
target a-layer



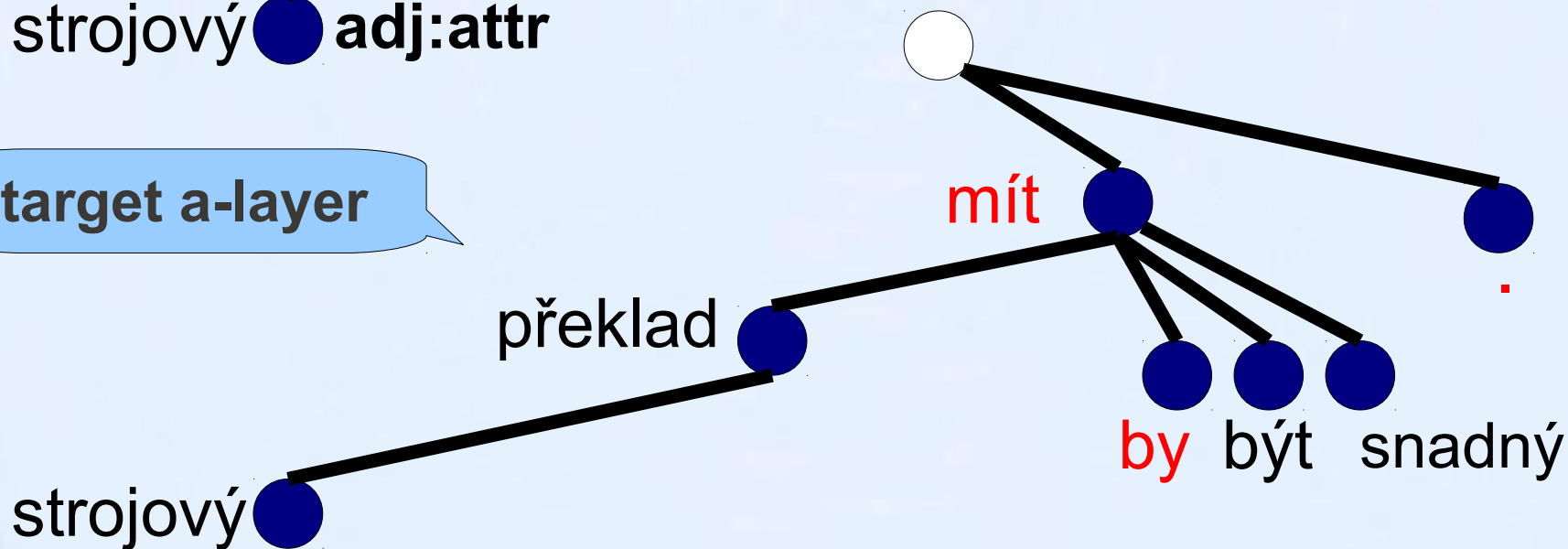
# Demo Translation – Synthesis

Add functional words

target t-layer



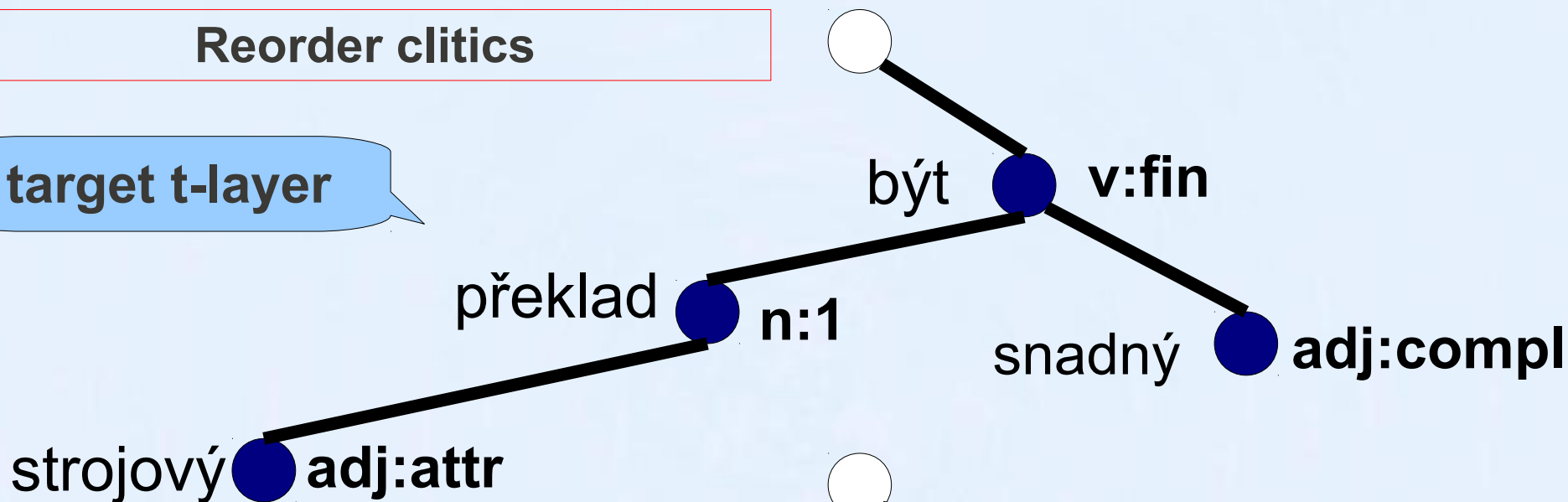
target a-layer



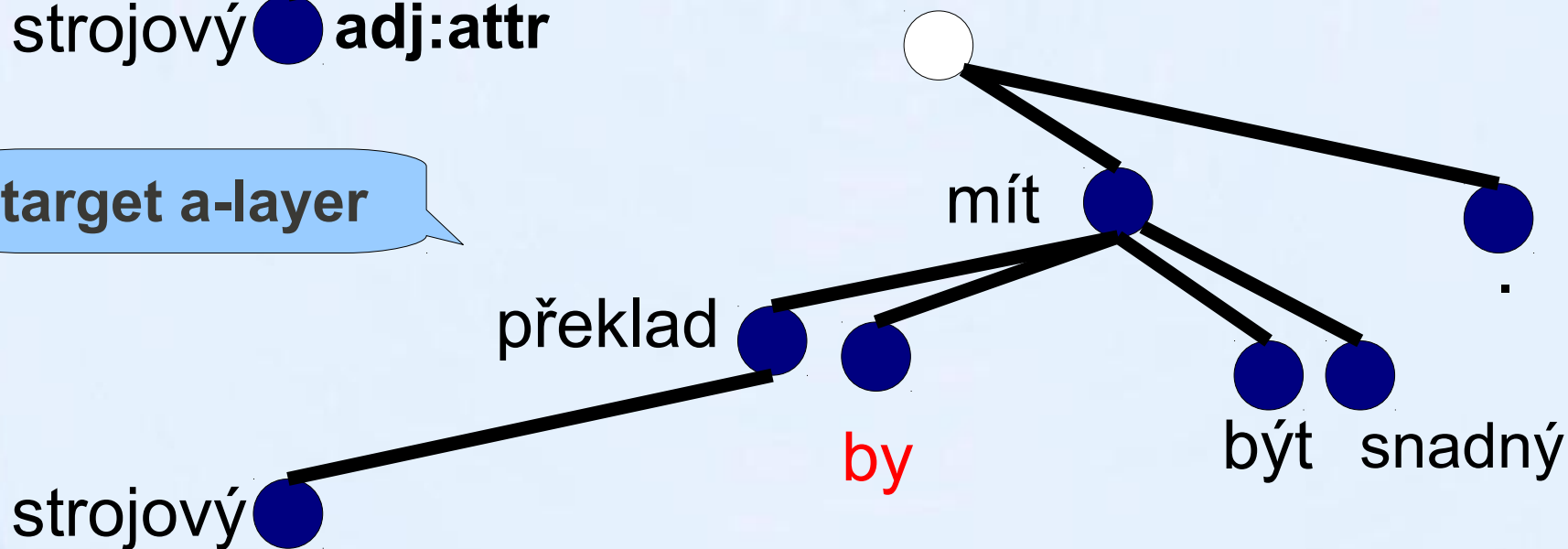
# Demo Translation – Synthesis

Reorder clitics

target t-layer



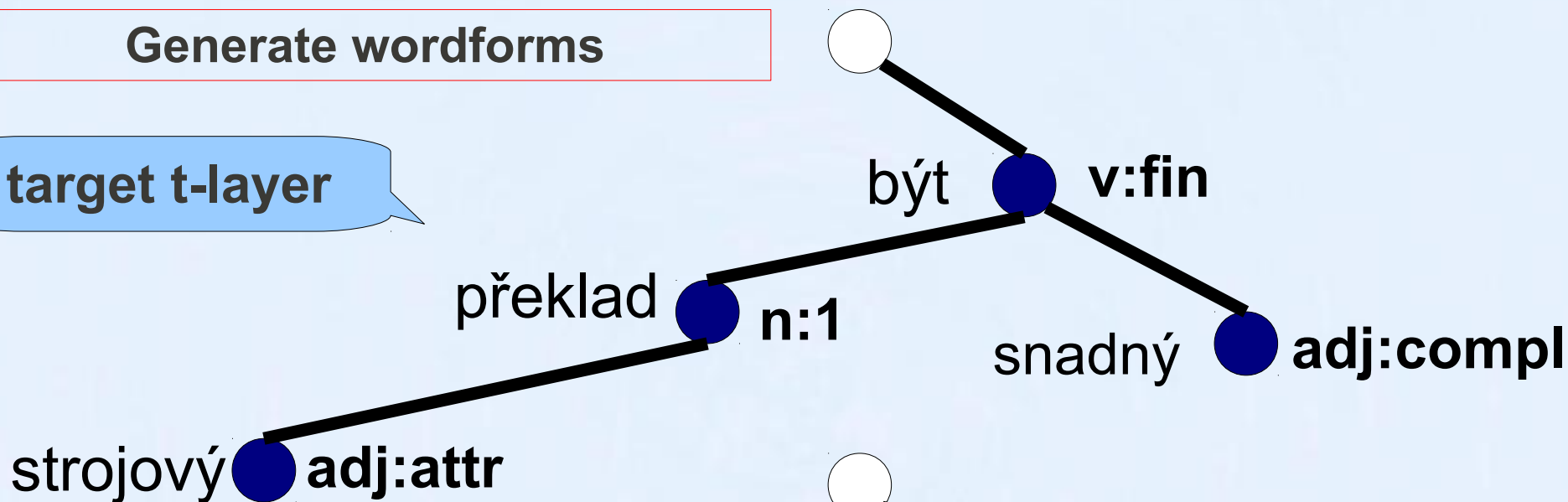
target a-layer



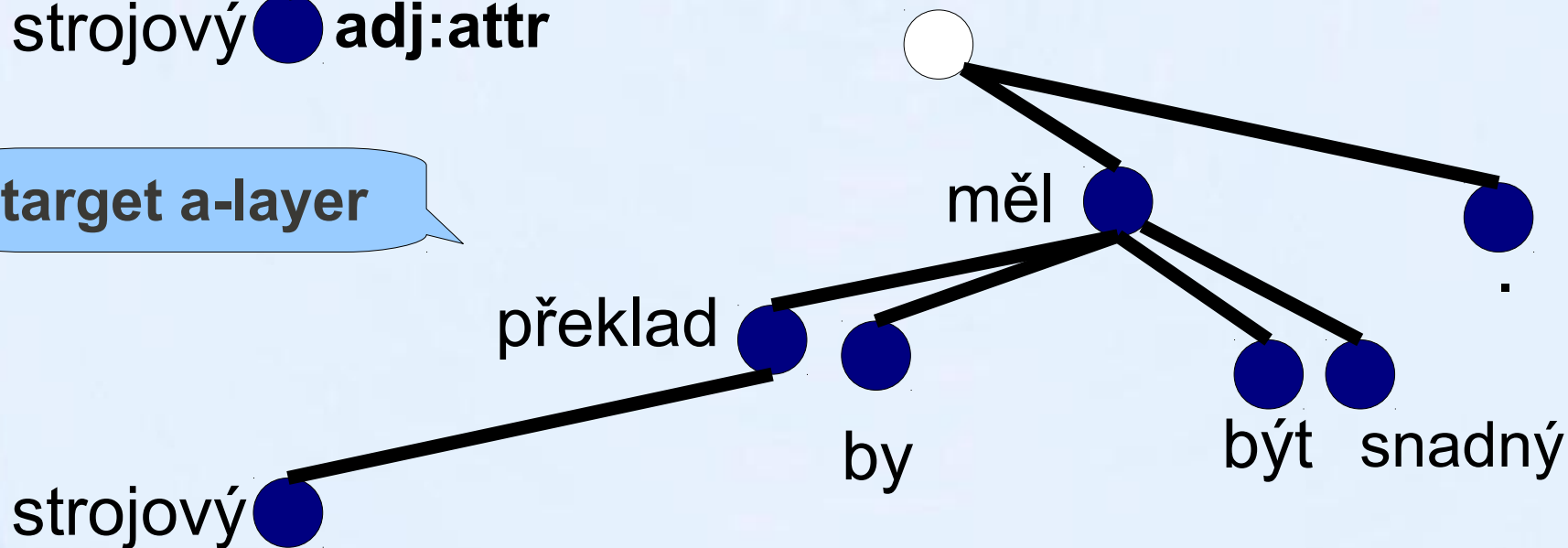
# Demo Translation – Synthesis

Generate wordforms

target t-layer



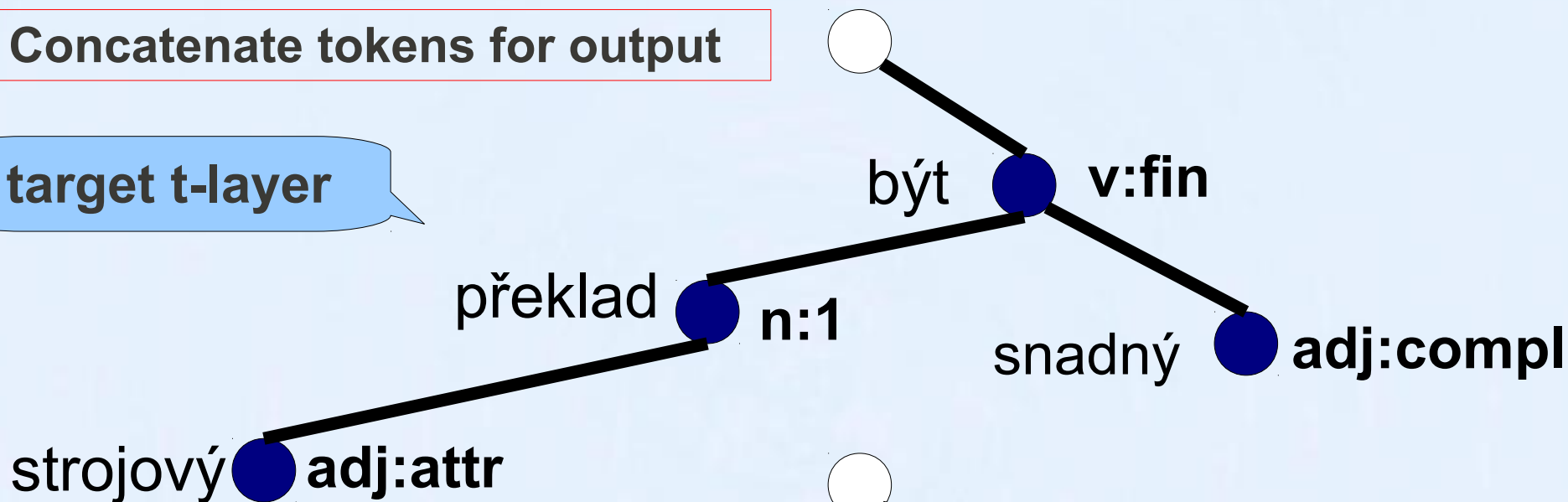
target a-layer



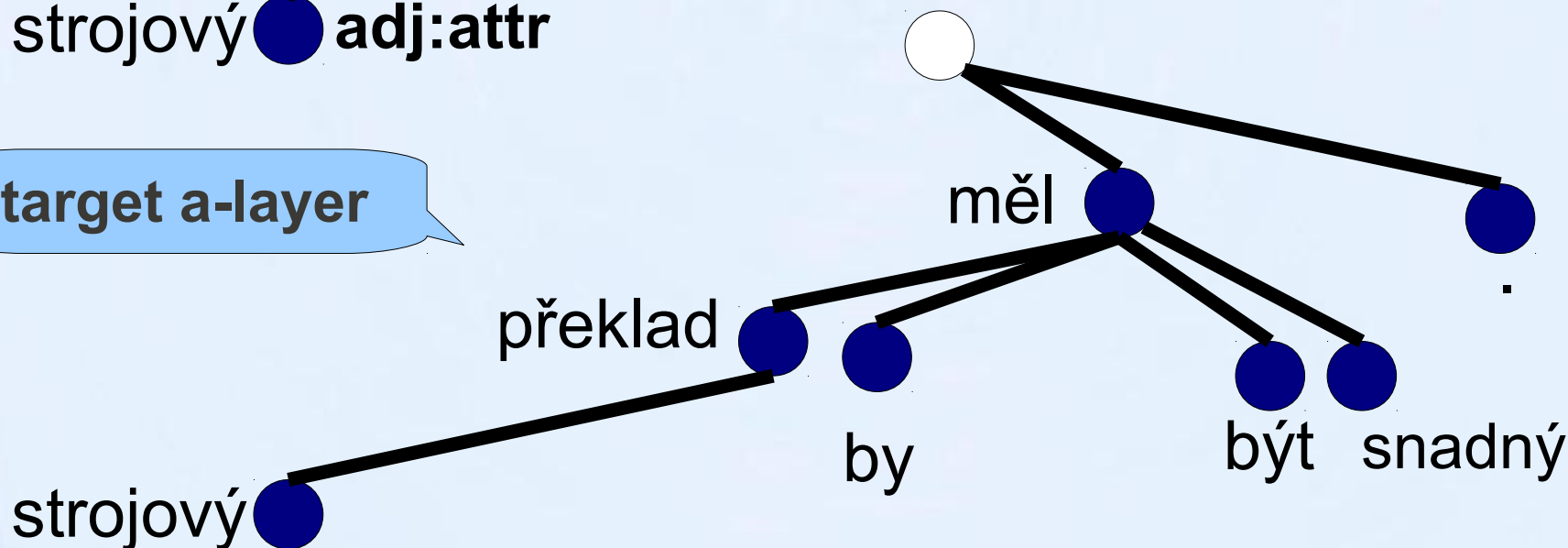
# Demo Translation – Synthesis

Concatenate tokens for output

target t-layer



target a-layer



**Strojový překlad by měl být snadný.**

# Demo Translation – Real Scenario



## **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**

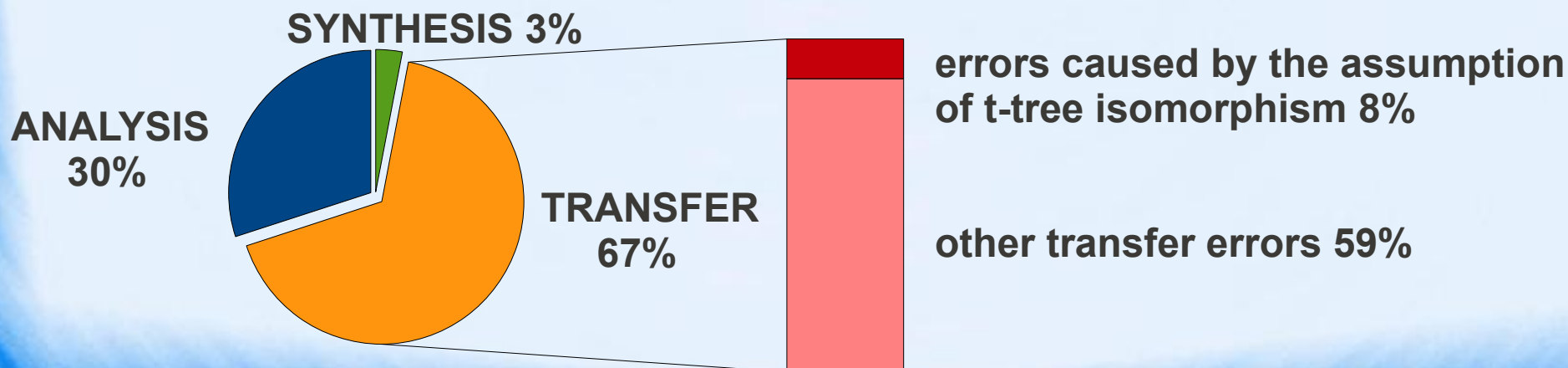
Ascii\_quotes

Remove\_repeated\_tokens

# Annotation of Translation Errors

sample of 250 sentences, 1463 errors in total

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



# Recent Improvements

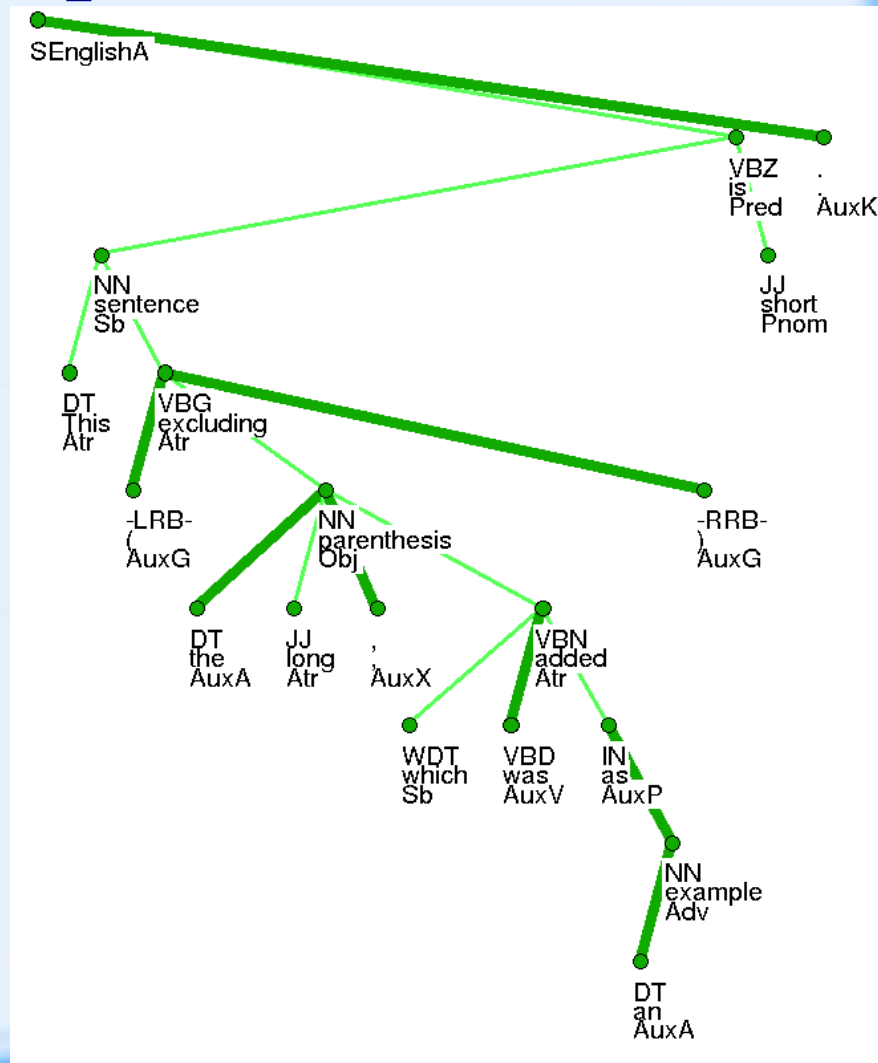
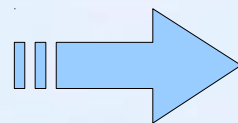
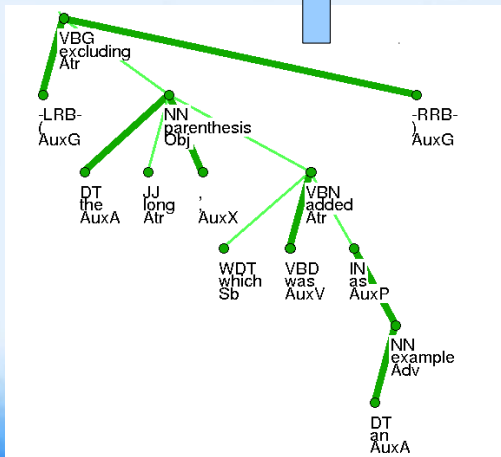
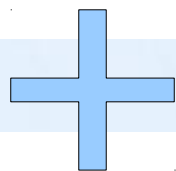
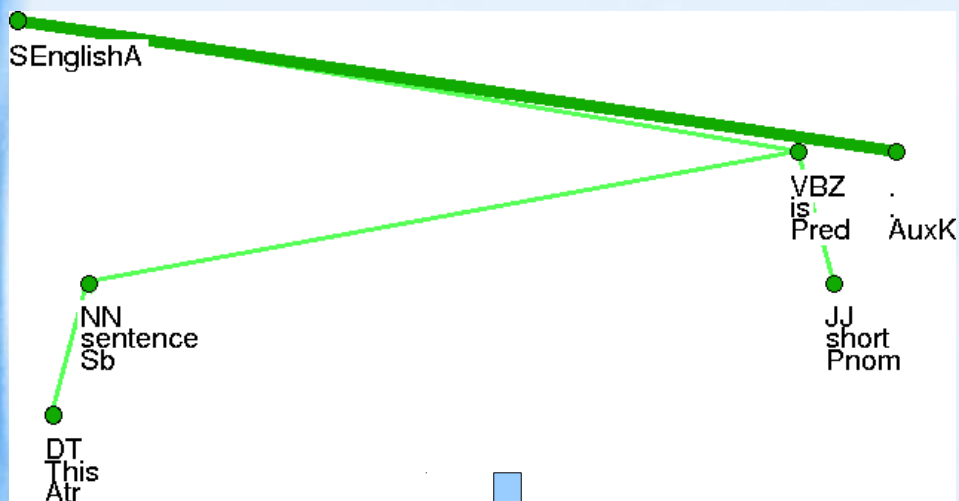
---

- Parsing parentheses
- Hidden Markov Tree Models (HMTM)
- Combining dictionaries
- Maximum Entropy dictionary



# Parsing Parentheses

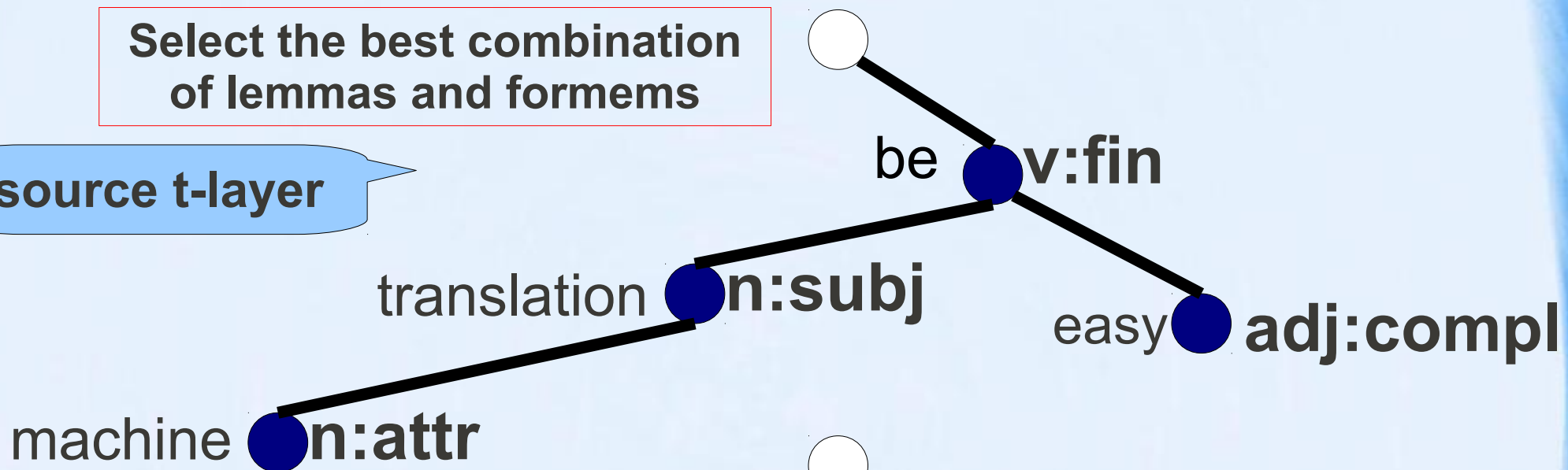
This sentence (excluding the long parenthesis, which was added as an example) is short.



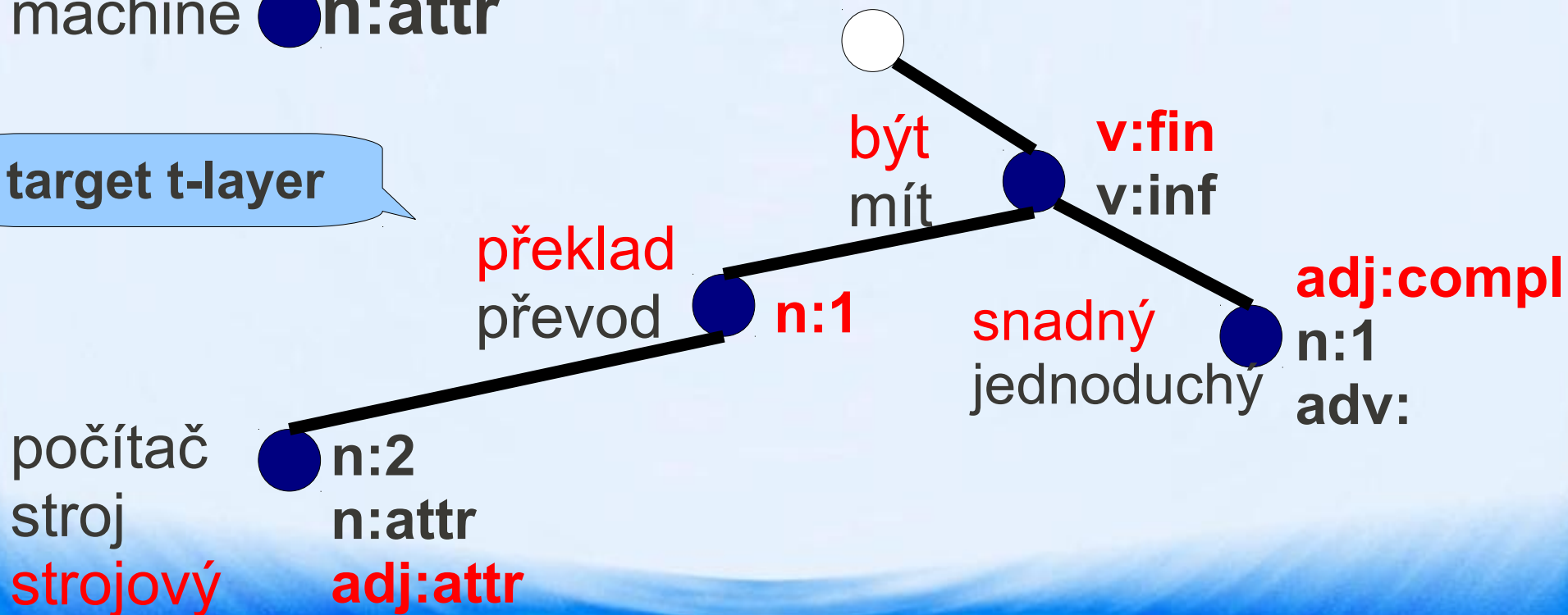
# HMTM – Motivation

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

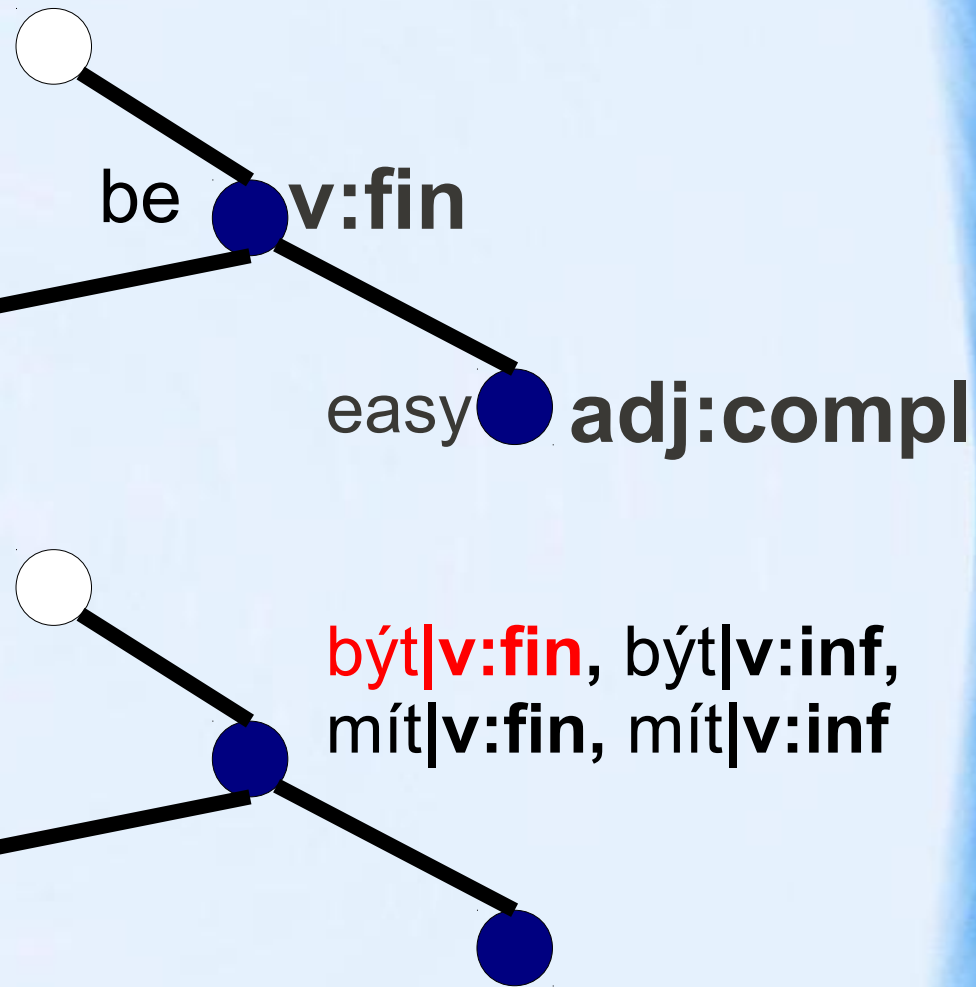
target t-layer

překlad|n:1,  
převod|n:1

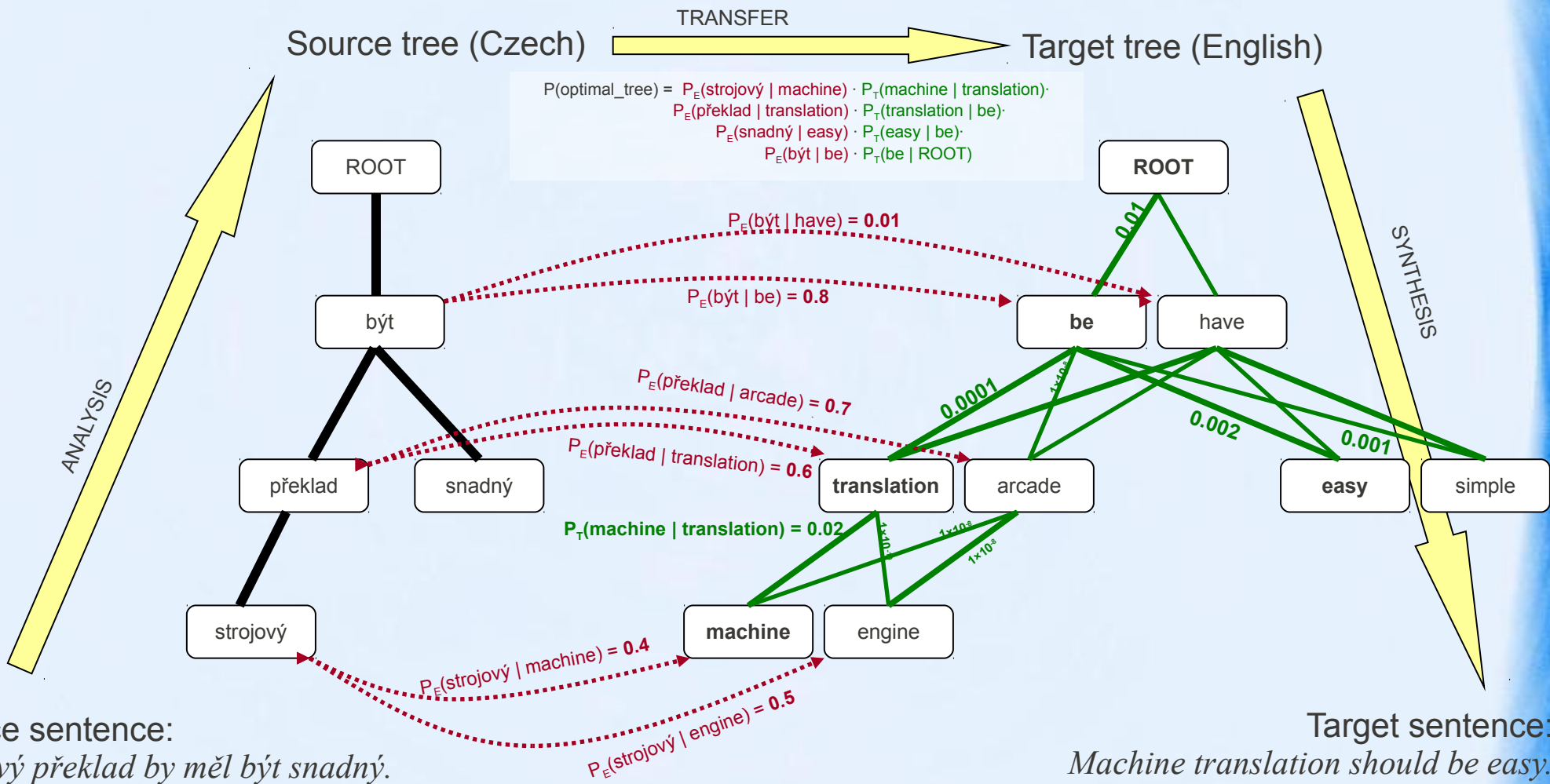
být|v:fin, být|v:inf,  
mít|v:fin, mít|v:inf

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“

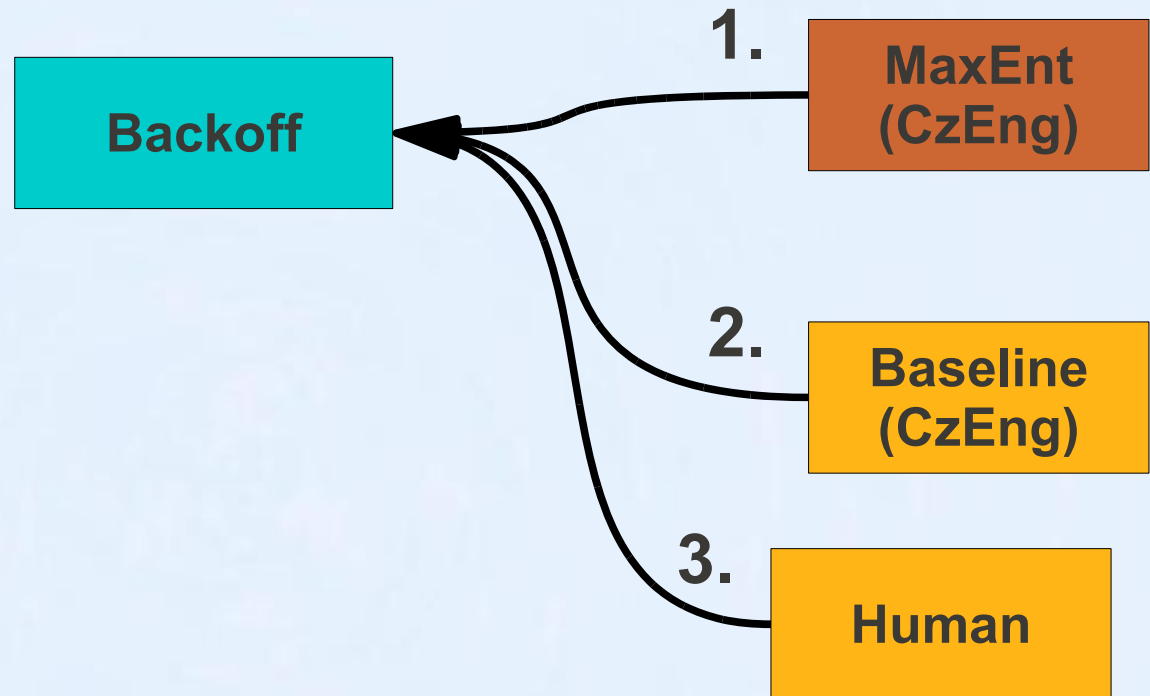
**Derivational**

translations derived dynamically, input dictionary

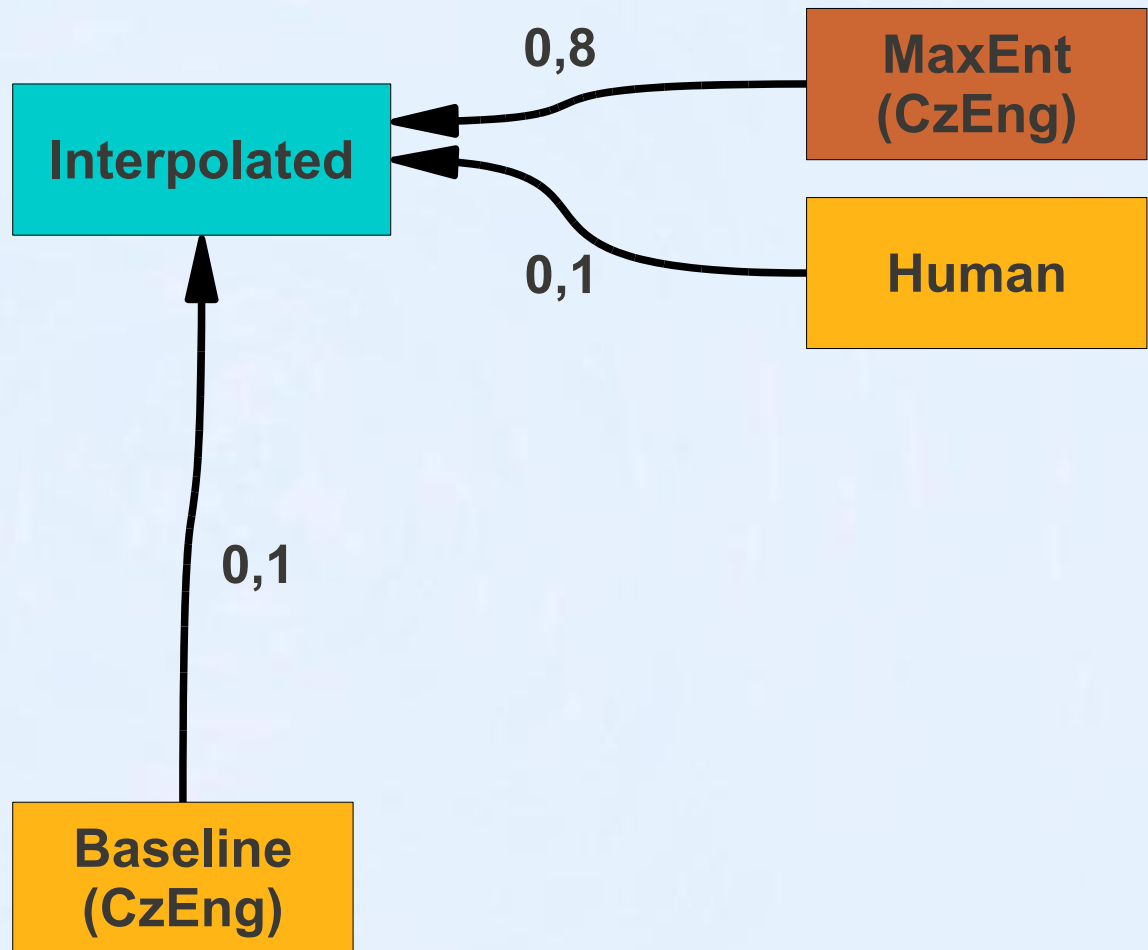
**Combinational**

combination of more input dictionaries

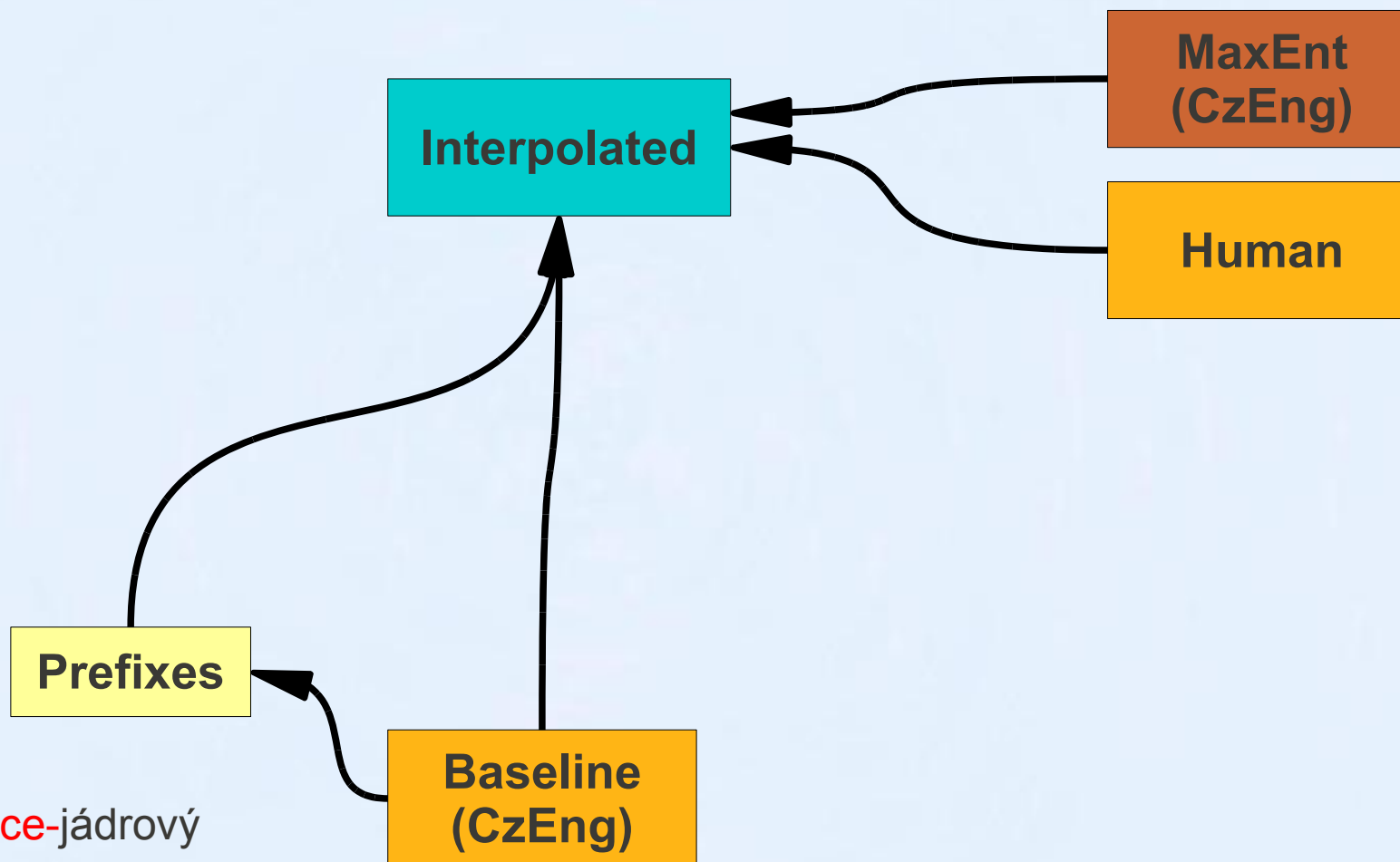
# Hierarchy of lemma dictionaries



# Hierarchy of lemma dictionaries



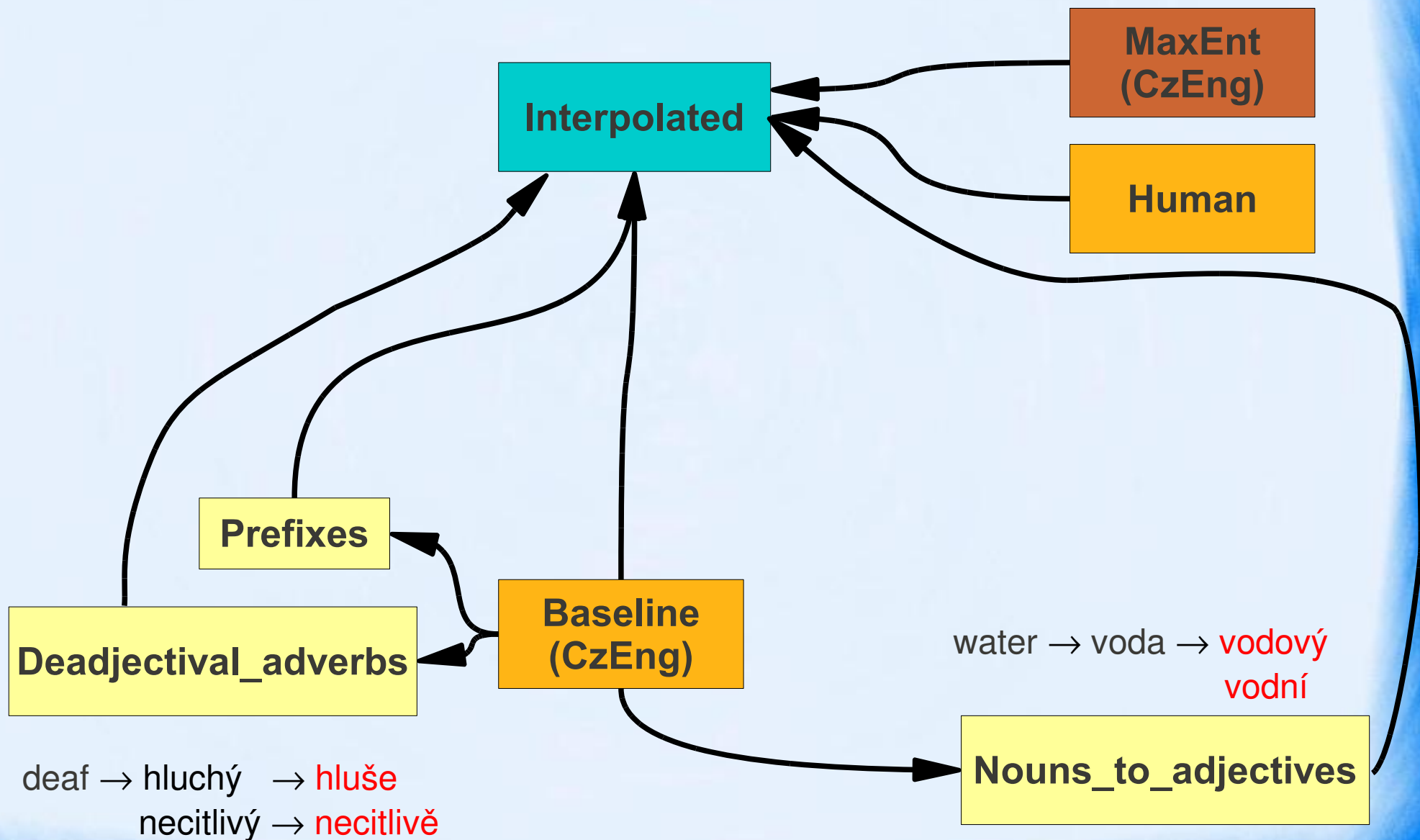
# Hierarchy of lemma dictionaries



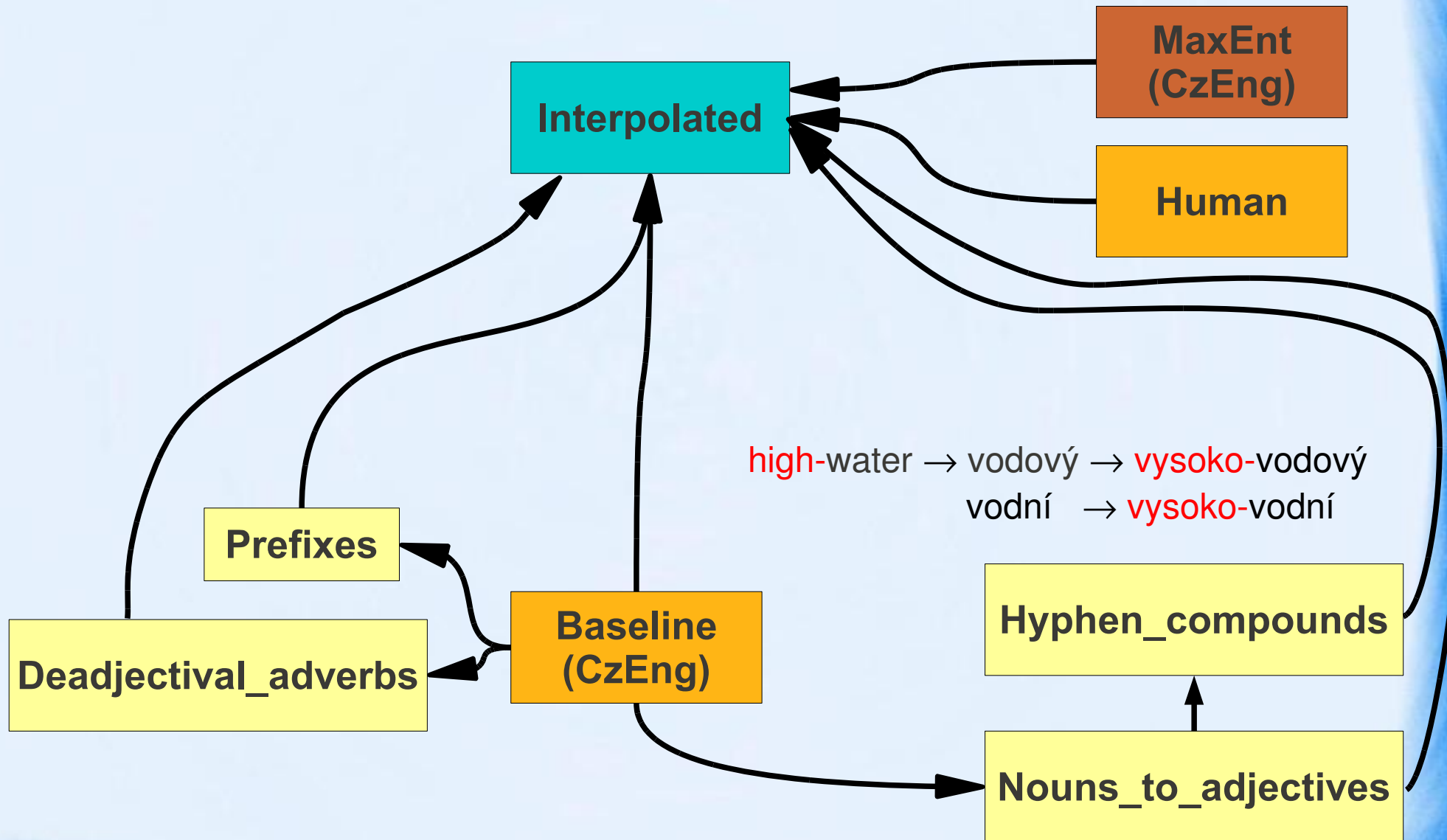
multi-core → více-jádrový  
více-jádro  
multi-jádrový  
multi-jádro



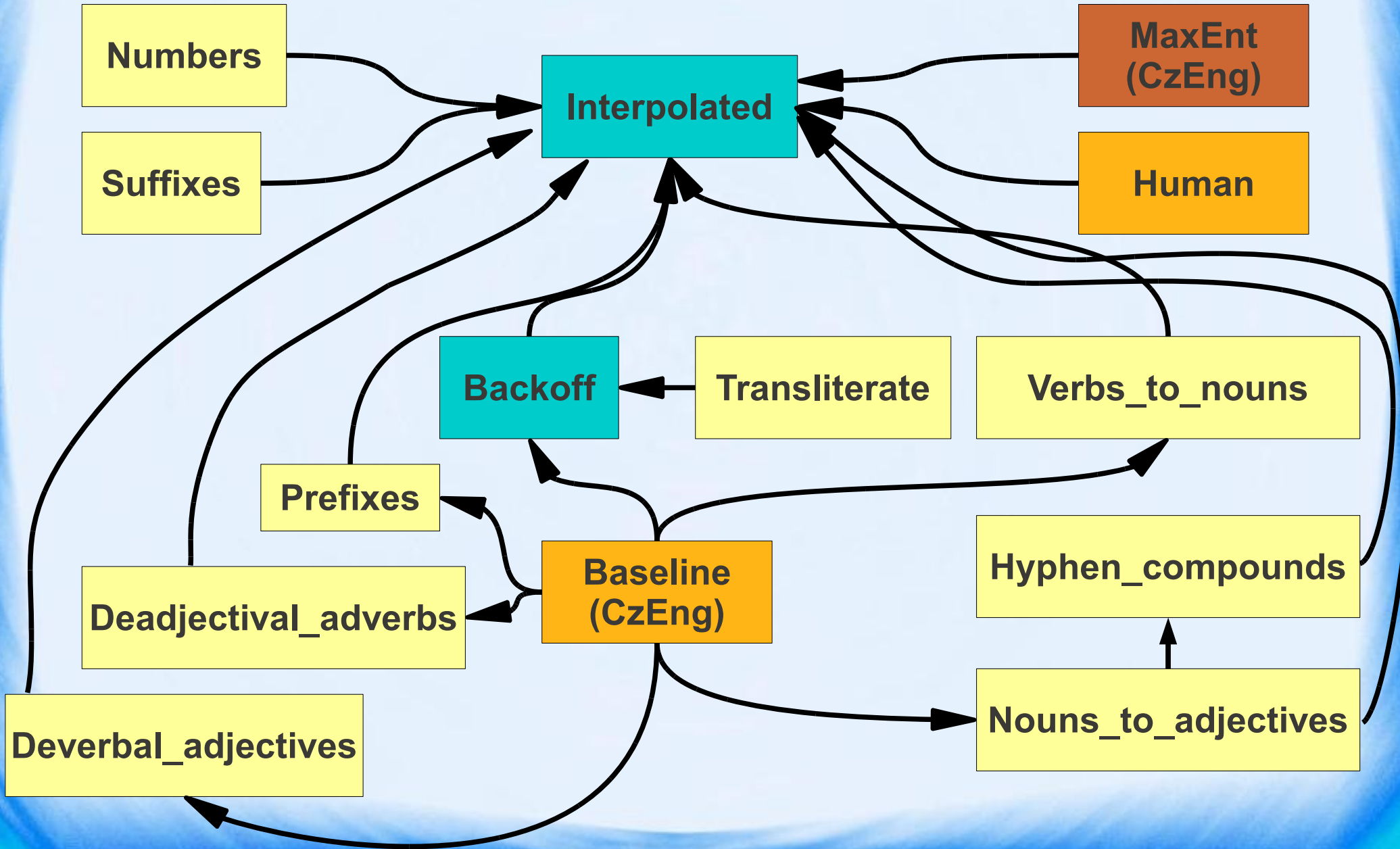
# Hierarchy of lemma dictionaries



# Hierarchy of lemma dictionaries



# Hierarchy of lemma dictionaries



# Maximum Entropy Dictionary

## Baseline Dictionary

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

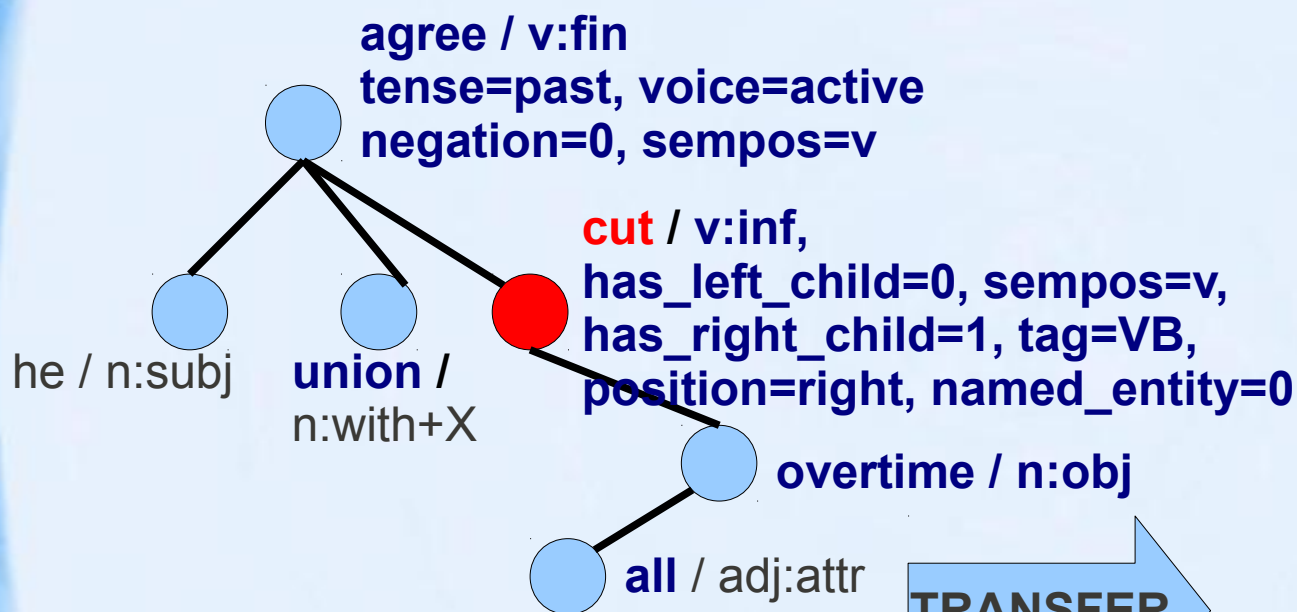
- Maximum likelihood estimates  
(from the training sections of CzEng 0.9)
- Pruned by thresholds on  $p(x|y)$  and  $p(y|x)$
- No context used  
x = source lemma  
y = target lemma

## MaxEnt Dictionary

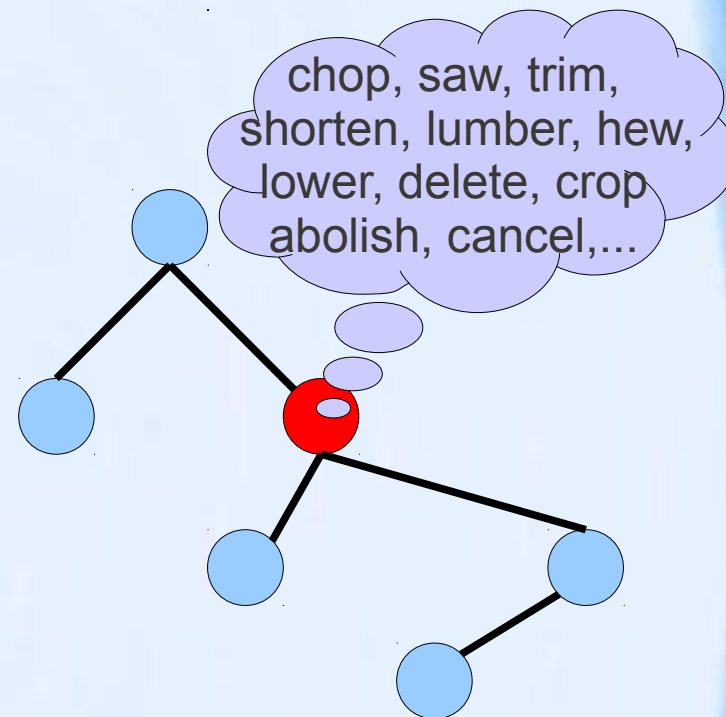
$$p(y|x) = \frac{1}{Z(x)} \exp \sum_i \lambda_i f_i(x, y)$$

- 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

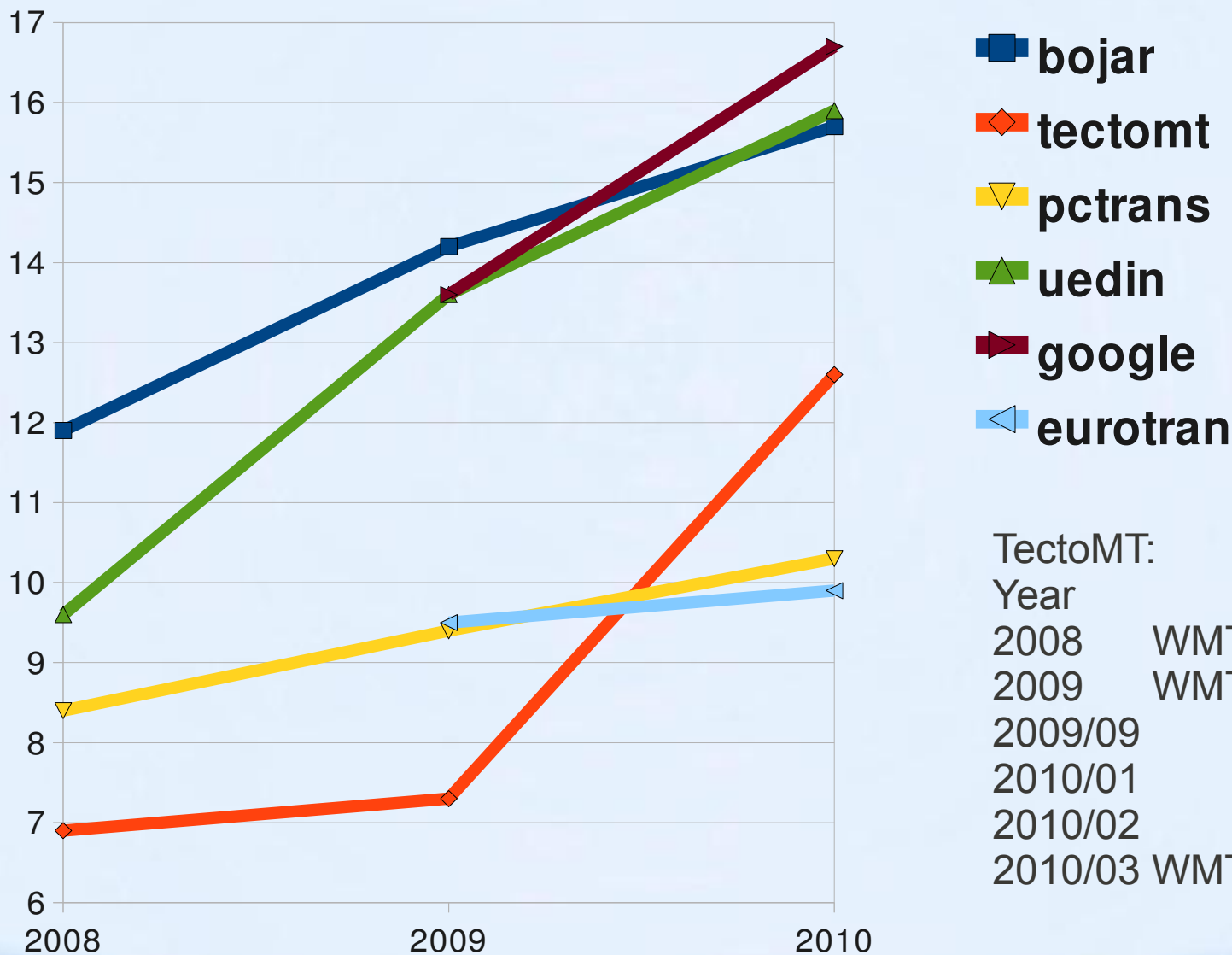
SYNTHESIS

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

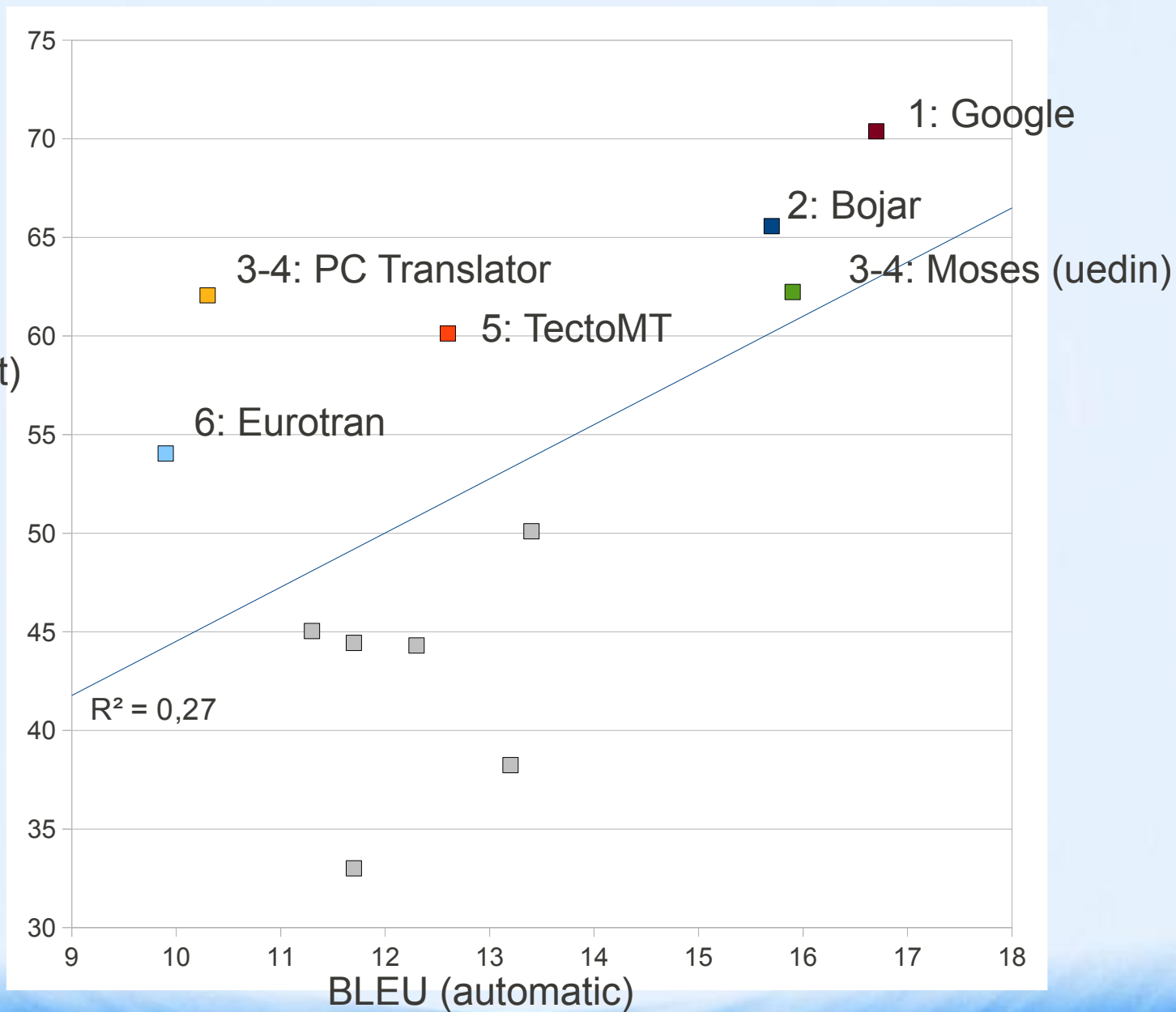


TectoMT:

| 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)



# Results – task-based evaluation

## Example

- **Translated text:**

*Bílý dům náčelník štábu Rahm Emanuel je naplánovaný k pojidany 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 |



# Examples of Translation

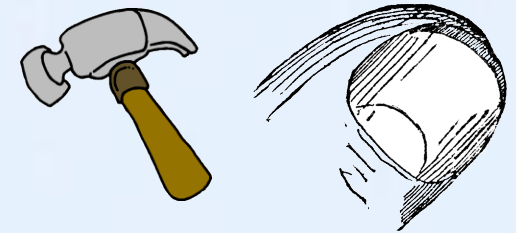
A miss by an inch  
is a miss by a mile.

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.

