

# PLIN009 – Strojový překlad

## Pravidlový strojový překlad

**Vít Baisa**

jaro 2012

10. dubna 2013











## Přístup pomocí interlinguy

- předpokládá, že je možné SL konvertovat do sémanticko-syntaktické reprezentace, která je (částečně) nezávislá na jazyku
- interlingua musí být jednoznačná (unambiguous)
- z této podoby (interlingua) je generován TL
- analýza SL je jazykově závislá, ale nezávislá na TL
- analogicky syntéza TL
- SL a TL nepřijdou do styku
- pro překlad mezi  $N$  jazyky potřebujeme  $2 \times N$  modulů







# Tokenizace

- 1 Úvod
- 2 Tokenizace**
- 3 Morfologická rovina
- 4 Lexikální rovina
- 5 Syntaktická rovina
- 6 Sémantika a logika

# Tokenizace

## Co to je?

- rozdělení vstupního řetězce do tokenů
- token = řetězec znaků
- výstup *tokenizace* = seznam tokenů
- slouží jako vstup pro další zpracování
- označení hranic vět

## Problémy

- don't: do\_n't, do\_n\_’t, don\_’t, ?
- červeno-černý: červeno\_-černý, červeno-černý, červeno\_černý
- Zeleninu jako rajče, mrkev atd. ¶Petr nemá rád.
- Složil zkoušku a získal titul Mgr. ¶Petr mu dost záviděl.

# Tokenizace – jak se to dělá?

V drtivé většině případů heuristika. (`unitok.py`)

## Dělení na tokeny

- pro jazyky používající hlásková písmata: dělení podle mezer
- a podle dalších interpunkčních znamének
- `? ! . , - ( ) / : ;`

## Dělení na věty

- MT v naprosté většině případů pro věty
- u plaintextu: podle seznamu interpunkčních znamének
- problém: Měl jsem 5 (sic!) poznánek.
- výjimky: zkratky (aj., atd., etc.), tituly (RNDr., prof.)
- někdy (HTML) lze využít strukturní značky

# Morfologická rovina

- 1 Úvod
- 2 Tokenizace
- 3 Morfologická rovina**
- 4 Lexikální rovina
- 5 Syntaktická rovina
- 6 Sémantika a logika









# Problém s víceznačností

- v mnoha případech: více morfologických značek
- víceznačnost mezi slovními druhy (více lemmat)  
*jednou* → k4gFnSc7, k6eAd1, k9  
*ženu* → k1gFnSc4, k5eAaImIp1nS
- víceznačnost v rámci slovního druhu
- typicky (čeština): nominativ = akuzativ  
*víno* → k1gNnSc1, k1gNnSc4, ...  
*odhalení* → 10 značek



# Statistická disambiguace

- nejpravděpodobnější posloupnost značek  
*Ženu je domů.*  
 k5 | k1, k3 | k5, k6 | k1  
*Mladé muže*  
 gF | gM, nS | nP
- těžká situace: *dítě škádlí lvíče*
- strojové učení na ručně značkových datech
- různé metody: Brill, TreeTagger
- pro češtinu: Desamb (hybridní)
- je nutné mít k dispozici trénovací data (korpus)



# Morfologická segmentace

- proč místo lemmatu (např. infinitiv) nepoužít kořen slova?
- existují i systémy, které provádí segmentaci automaticky na základě seznamu slov pro daný jazyk
- problém: *mít, měj, mám, měl, mívá, ...* – různé podoby téhož morfému
- problém: *i, ové, a, y* – stejná gramatická funkce, různé morfémy
- *bychom* → bych?
- gramatické kategorie mají konkrétní formu (gramémy)  
*nad-měr-ný, ne-patr(n)-ně, vid-ím, ne-chci, čtyř-i-cet, po-po-sun-out, u-dě-l-al-i*
- nutné pokud nemáme morfologický analyzátor k dispozici

Morfologie – závěrem

<b>slovo</b>	<b>analýzy</b>	<b>disambiguace</b>
Pravidelné	k2eAgMnPc4d1, k2eAgInPc1d1, k2eAgInPc4d1, k2eAgInPc5d1, k2eAgFnSc2d1, k2eAgFnSc3d1, k2eAgFnSc6d1, k2eAgFnPc1d1, k2eAgFnPc4d1, k2eAgFnPc5d1, k2eAgNnSc1d1, k2eAgNnSc4d1, k2eAgNnSc5d1, ... (+ 5)	k2eAgNnSc1d1
krmení	k2eAgMnPc1d1, k2eAgMnPc5d1, k1gNnSc1, k1gNnSc4, k1gNnSc5, k1gNnSc6, k1gNnSc3, k1gNnSc2, k1gNnPc2, k1gNnPc1, k1gNnPc4, k1gNnPc5	k1gNnSc1
je	k5eAalmp3nS, k3p3gMnPc4, k3p3gInPc4, k3p3gNnSc4, k3p3gNnPc4, k3p3gFnPc4, k0	k5eAalmp3nS
pro	k7c4	k7c4
správný	k2eAgMnSc1d1, k2eAgMnSc5d1, k2eAgInSc1d1, k2eAgInSc4d1, k2eAgInSc5d1, ... (+ 18)	k2eAgInSc4d1
růst	k5eAalMF, k1gInSc1, k1gInSc4	k1gInSc4
důležité	k2eAgMnPc4d1, k2eAgInPc1d1, k2eAgInPc4d1, k2eAgInPc5d1, k2eAgFnSc2d1, k2eAgFnSc3d1, k2eAgFnSc6d1, k2eAgFnPc1d1, k2eAgFnPc4d1, k2eAgFnPc5d1, k2eAgNnSc1d1, k2eAgNnSc4d1, k2eAgNnSc5d1, ... (+ 5)	k2eAgNnSc1d1

# Universal POS tags

Počet značek se v různých jazycích značně liší → snaha o zjednodušení.

TAG	význam
VERB	verbs (all tenses and modes)
NOUN	nouns (common and proper)
PRON	pronouns
ADJ	adjectives
ADV	adverbs
ADP	adpositions (prepositions and postpositions)
CONJ	conjunctions
DET	determiners
NUM	cardinal numbers
PRT	particles or other function words
X	other: foreign words, typos, abbreviations
.	punctuation

Vytvořeno mapování pro cca 25 jazyků s *tree banks*.

# Odhadování POS na základě gramémů

EN	CZ	význam
-s	-á	3. os., j. č., přít.
-ed	-al, -l, -en.	minulý čas
-ing	-(ov)ání	průběhový čas
-en	-en(.)	příčestí minulé
-s	-y, -i, -ové, -a	množné číslo
-’s	ov(o, a, y)	přivlastňování
-er	-ší	komparativ
-est	nej-, -ší	superlativ

Problém: *myší, west, fotbal, . . .*



## Tomáš Hanák – Sám v lese II

Když jsi sám v lese,  
ano, sám-li v lese's,  
však skutečně, v lese sám's-li.  
Zkrátka v lese sám-li's.

Však kde vlastně vzal ty tu's?  
Z meze-li v les's vlez?  
Či z nebes v les se snesl's?

Pověz, ach, tvář tvá perlí přívalem se slz.  
Teď rud's, zas bled's, co pivoňka's  
Snad tedy autem's tu, či kolmo's?

Mlčíš a slza tvá dál  
sama malá padá v mechu číš.

Ano, teď teprve snad poprvé sám svět's.

# Brillův tagger

- učení z trénovacích dat
  - transformation-based, error-driven
  - úspěšnost přes 90 %
- 1 inicializuj značkování (nejčastější značka)
  - 2 porovnej s trénovacími daty
  - 3 vytvoř sadu pravidel pro změnu značek
  - 4 ohodnot' pravidla
  - 5 aplikuj pravidlo a opakuj od 2. dokud je co zlepšovat







# Slova a slovníky ve strojovém překladu

- 1 Úvod
- 2 Tokenizace
- 3 Morfologická rovina
- 4 Lexikální rovina**
- 5 Syntaktická rovina
- 6 Sémantika a logika







# Problém s víceznačností

- slovům odpovídají významy
- co je ale význam? pro počítač potřebujeme formální popis
- počítač je diskrétní, význam je zřejmě spojitý
- *muž* – dospělý člověk mužského pohlaví
- co 17letý člověk mužského pohlaví?

# Spojitost významu



špalek



?



židle





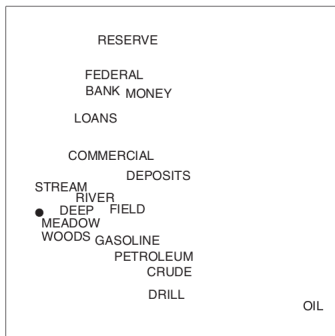
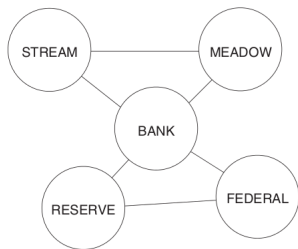






# Reprezentace významu

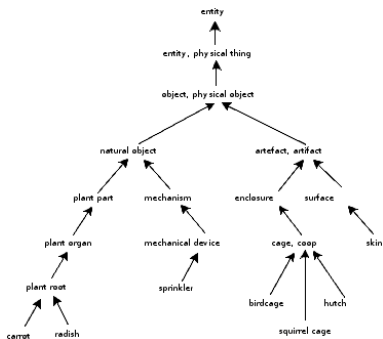
- nejčastější způsob: banka významů
- graf: významy jsou uzly, sémantické relace jsou hrany
- prostor: významy jsou body, podobné významy jsou prostorově blízko





# Sémantická síť – WordNet

- **literál** dát:8, **synset** louže:1, kaluž:1, tratoliště:1
- sémantické relace: hypero-, hypo-, holo-, meronymum
- 150k slov, 117k synsetů: n, adj, v a adv
- WN používán jako referenční banka významů



# VerbaLex

- WordNet neobsahuje syntaktické vazby, morfosyntaktické omezení
- synsety (6 256)  
atakovat:1, útočit:2, dorážet:3, napadnout:6
- valenční rámce (mačkat:1) a sloty (19 247)  
AG<sup>kdo1</sup><sub>person:1</sub> + VERB + OBJ<sup>co4</sup><sub>object:1</sub> + (PART<sup>v</sup><sub>hand:1</sub> čem6)
- sémantické role  
I: ABS, ISUB, AG, KNOW, PAT, VERB, ... (29)  
II: abstraction:1, person:1, artifact:1, body part:1, ... (10<sup>3</sup>)
- další omezení:  
předložkové pády, životnost, slovní druhy, obligatornost
- synsety napojeny na WordNet

# Word Sense Disambiguation

- nalezení významu slova v daném kontextu
- pro člověka triviální, pro PC těžké
- jde o klasifikační úlohu
- potřebujeme konečný inventář významů
- při použití WN: pro dané slovo určit konkrétní synset
- kvalita se těžko vyhodnocuje (SensEval, SemEval)
- přesnost kolem 90 %

Pro strojový překlad zásadní:

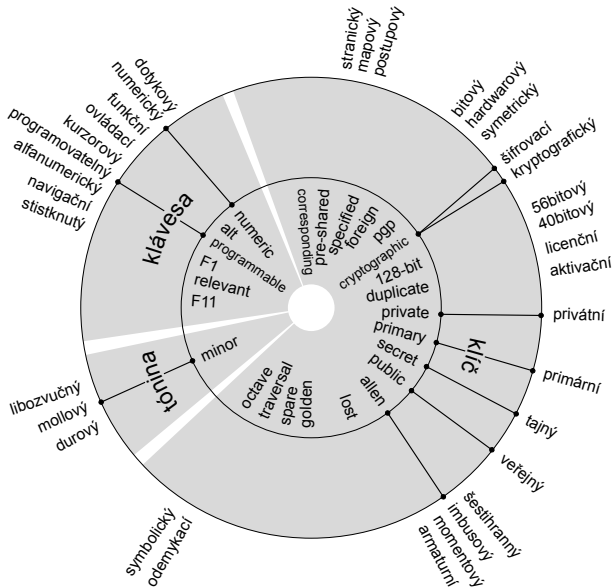
*Ludvig dodávka Beethoven, kiss me honey, . . .*

# WSD – metody

Problém: jak přeložit **box in the pen** (Bar-Hillel).

- hloubkové (deep)
  - využívají znalosti o světě (common sense)
  - nejsou vhodné pro obecný jazyk (spíše omezené domény)
  - znalosti typu: ptáci umí létat, jablka rostou na stromě, . . .
  - metody založené na reprezentaci znalostí, na slovníku
  - Leskův algoritmus: shoda slov z okolí se slovy ze slovníku patřícími ke konkrétnímu významu
  - algoritmy s využitím valenčních slovníků (BP)
- povrchové (shallow)
  - využívají slova z kontextu
  - levnější, rychlejší implementace
  - různé metody strojového učení (klasifikační problémy)
  - učení s učitelem (supervised), bez učitele (unsupervised)
  - možné použít varianty Brillova algoritmu

# Jak určit vhodnou granularitu?



# Lexika – shrnutí

- význam hlavně na úrovni slov (překladové slovníky)
- WSD zcela klíčový pro pravidlové systémy
- počet slov se mezi jazyky řádově liší
- na přesnost WSD má nejvíce vliv požadovaná granularita
- lexikální víceznačnost je bottleneck (RB)MT

# Syntaktická analýza

- 1 Úvod
- 2 Tokenizace
- 3 Morfologická rovina
- 4 Lexikální rovina
- 5 Syntaktická rovina**
- 6 Sémantika a logika

# Syntaktická analýza I

- další patro v MT trojúhelníku
- snaha o konečný popis nekonečného množství frází, vět
- konečným způsobem = gramatikou
- vstup (většinou): morfologicky označovaná data
- výstup: syntaktický strom, les, graf



# Syntaktická analýza II

- úkol SA: pro danou gramatiku a vstupní větu vrať všechny možné derivační stromy
- potenciálně milióny různých analýz (viz Synt)
- pro analýzu je potřeba:
  - výběr formalismu
  - napsání gramatiky
  - implementace algoritmu analýzy
- v současnosti většina **parserů** využívá statistiky

# Syntaktická analýza III

## Gramatické formalismy

- bezkontextová gramatika: na levé straně mohou být pouze jednoduché **neterminály**
- regulární gramatika: bezkontextová + pravidla pouze typu  $N \rightarrow \epsilon \mid A \mid bB$
- tree-adjoining: podobné bezkontextovým, přepisují se stromy nikoli znaky (řetězce)

## Typy analýz

- top-down analýza (shora): hledá se taková nejlevější derivace, která generuje analyzovaný řetězec
- bottom-up analýza (zdola): hledají se pravidla, která přepíšou vstupní řetězec na výslednou posloupnost pravidel

# K čemu je syntaktická analýza?

- sémantická interpretace zdrojového kódu (informatika)
- mezistupeň k sémantické reprezentaci věty
- transferové systémy: konečný počet transferových pravidel pro nekonečný počet možných frází
- WSD: zachycení vztahů na větší vzdálenosti (širší kontext)
- jaká slova k sobě patří a jaká ne

# Syntaktická víceznačnost I

- *I saw a man with a telescope.*  
Uzřel jsem muže (s) dalekohledem.
- *I'm glad I'm a man, and so is Lola.*  
Jsem rád, že jsem muž a Lola také.
- *Someone ate every tomato.*  
Někdo snědl všechna rajčata.  
Každé rajče bylo někým sněženo.
- *Lvíče škádlí dítě.*  
A child teases a lion cub.  
A lion cub teases a child.

## Syntaktická víceznačnost II

- *Letadlo spadlo do pole za lesem.*
- *Ženu holí stroj.*  
*Ženu holý stroj.*
- *Zabít ne propustit.*  
*Ibis, redibis nunquam per bella peribis.*
- *Rodiče by mu mohli závidět.*
- *Neboť každý, kdo prosí, dostává a kdo hledá, nalézá a tomu, kdo tluče, bude otevřeno. (Lk: 11,10)*

# Částečná synt. víceznačnost – garden path

- *The man returned to his house ... was happy.*
- *The man whistling tunes ... pianos.*
- *Time flies like an arrow; fruit flies like a banana.*
- *Ženu krávy ... nezajímají.*

# Vyhodnocení kvality syntaktické analýzy

- jaká analýza je nejlepší?
- vyhodnocení kvality je obtížné a interpretace je sporná
- nejlepší analyzátoři dosahují přesnosti cca 85 %

# Frázová struktura

- jeden z nejstarších formalismů
- gramatika obsahuje přepisovací pravidla
- nejčastěji bezkontextová gramatika
- zachycuje, jak se skládají fráze: **konstituenty**

S → NP VP

VP → ADV V | V ADV

NP → DET N

DET → the | a | an

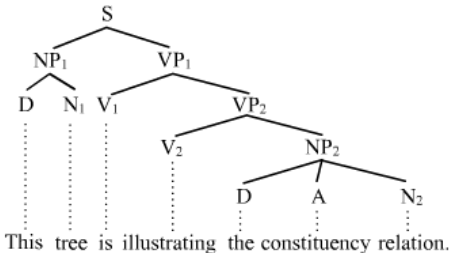
N → cat | dog

...

Analýza: **the dog runs fast** (shora a zdola)



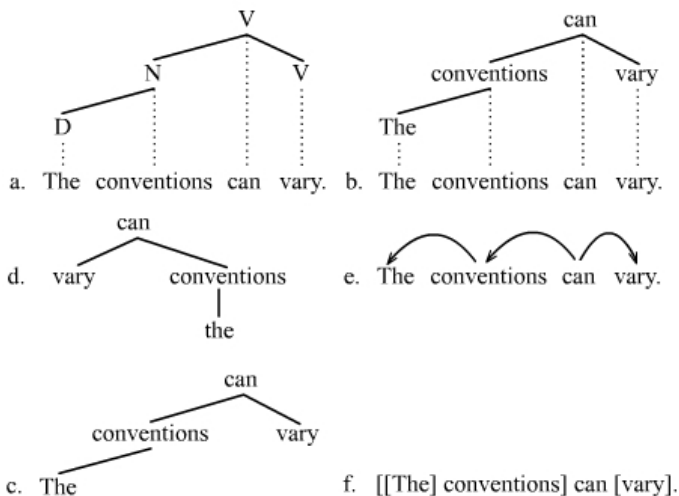
# Frázový strom



# Závislostní struktura

- zachycuje závislosti mezi slovy
- strom neobsahuje **neterminály**
- hlava a závislá slova
- vhodné pro jazyky s volným slovosledem (čeština)

# Závislostní strom



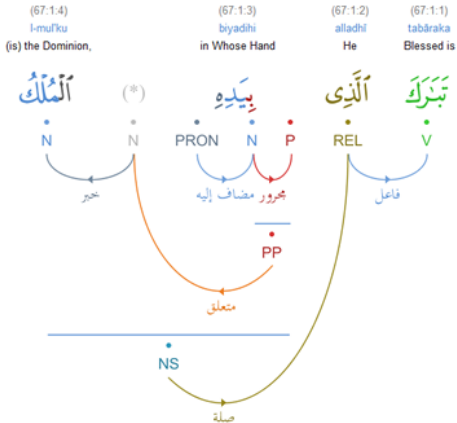
# Constituency vs. Dependency

- každé paradigma vhodné pro něco jiného
- složky: pevný slovosled, koordinace
- nevýhoda: neschopnost zachytit neprojektivitu souvislým složkovým stromem  
neprojektivní závislost = závislost mezi dvěma slovy oddělenými ve větě třetím slovem, které nezávisí na žádném z nich  
*I saw a man with a dog yesterday which was a yorkshire terrier.*
- závislosti: volný slovosled, morfosyntaktická shoda
- nevýhoda: neschopnost zachytit doplněk (dvojitá závislost)  
*Babička seděla u stolu shrbená.* (doplněk)  
Babička seděla u stolu shrbeně. (PUZ)
- lze převádět mezi sebou nebo kombinovat: hybridní stromy

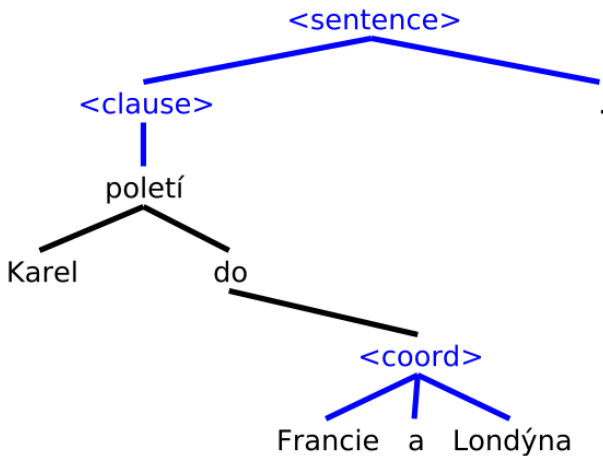
## Hybridní strom I

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

## Chapter (67) sūrat l-mulk (Dominion)



# Hybridní strom II



## Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“



# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor
- slovo obsahující 3 po sobě jdoucí diakritická znaménka

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor
- slovo obsahující 3 po sobě jdoucí diakritická znaménka jednodušší

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor
- slovo obsahující 3 po sobě jdoucí diakritická znaménka jednodušší
- věta obsahující 4x po sobě jdoucí „se“

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor
- slovo obsahující 3 po sobě jdoucí diakritická znaménka jednodušší
- věta obsahující 4x po sobě jdoucí „se“ nesnese se se sestrou

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor
- slovo obsahující 3 po sobě jdoucí diakritická znaménka jednodušší
- věta obsahující 4x po sobě jdoucí „se“ nesnese se se sestrou
- slovo, 5 písmen, význam i retrográdně

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor
- slovo obsahující 3 po sobě jdoucí diakritická znaménka jednodušší
- věta obsahující 4x po sobě jdoucí „se“ nesnese se se sestrou
- slovo, 5 písmen, význam i retrográdně tokej, jelen

# Hledání slov a vět splňujících podmínku

Slovní tvary jako ve scrabble.

- slovo obsahující 3x „r“ reproduktor
- slovo obsahující 3 po sobě jdoucí diakritická znaménka jednodušší
- věta obsahující 4x po sobě jdoucí „se“ nesnese se se sestrou
- slovo, 5 písmen, význam i retrográdně tokej, jelen
- slovo, které má význam i v češtině i v angličtině





















# Sémantické role

- syntaxe umožňuje odhalit sémantické vztahy
- konstituenty vět odpovídají **sémantickým rolím**
- vztah **predikátu** a ostatních větných členů
- také **semantic case, thematic role, theta role**
- agent, causer, instrument, manner, patient, result, time, source
- různé množiny rolí, viz např. VerbaLex (29 rolí)

**Dítě      škádlí      lvíče.**

AG/SUBJ   PRED/V   PAT/OBJ

**A child** (SUBJ) **teases** (PRED/V) **a lion cub** (PAT/OBJ).

**A lion cub** (SUBJ) **teases** (V) **a child** (OBJ).

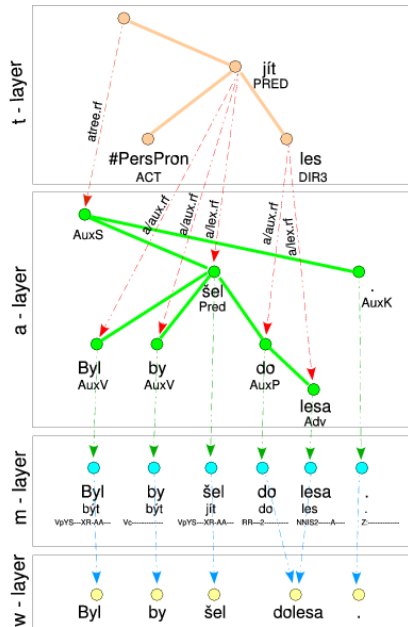




# Prague Dependency TreeBank 2.0

- aplikace teorií Pražského lingvistického kroužku
- funkční generativní popis jazyka
- rovina: fonologická a fonetická, morfonologická, morfematická, povrchová syntax a
- tektogramatická rovina – rovina významu jazyka
- *nižší rovina je formou vyšší a vyšší rovina funkcí nižší*
- 2M morfologicky, 1,5M syntakticky a 800k sémanticky označovaných slov z novinových článků v ČNK
- koreference a aktuální členění větné  
*Petr dal Petře kytici. Pak ji vzal a dal do vázy.*
- uzly pro nevyjádřená slova
- vazby mezi uzly na různých úrovních

## Prague Dependency TreeBank



# TectoMT – systém

- vysoká modularita
- maximální rozložení úkolů do série bloků – scénáře
- bloky jsou Perl moduly, komunikují přes API
- struktura systému odpovídá struktuře PDT
- vnitřní reprezentace jazyka: stromy v tmt formátu odvozeném od PML pro PDT
- bloky umožňují masivní zpracování dat, paralelizace
- bloky mohou implementovat pravidlové, stochastické či hybridní metody
- zpracování:
  - ① konverze do formátu tmt
  - ② aplikace scénáře
  - ③ konverze do výstupního formátu

# TectoMT – jednoduchý blok

Převod anglických negativních částic na příznaky sloves.

```

sub process_document {
  my ($self,$document) = @_ ;

  foreach my $bundle ($document->get_bundles()) {
    my $a_root = $bundle->get_tree('SEnglishA');

    foreach my $a_node ($a_root->get_descendants) {
      my ($eff_parent) = $a_node->get_eff_parents;
      if ($a_node->get_attr('m/lemma')=~/(not|n\`t)$/
        and $eff_parent->get_attr('m/tag')=~/^V/ ) {
        $a_node->set_attr('is_aux_to_parent',1);
      }
    }
  }
}

```

## Transferový systém TectoMT

