

st. 10.50-12.25 G13

CJBB75 ZÁKLADY VYUŽITÍ KORPUSU PRO PRAXI

RŮZNÉ KORPUSY A ROZDÍLY V ANOTAČNÍCH SCHÉMATECH

- ✘ Tokenizace, automatická anotace, d[e]sambiguace
- ✘ Anotace velkých synchronních korpusů ČNK
- ✘ Anotace mluvených korpusů
- ✘ Anotace KSK
- ✘ Pražský a brněnský systém anotací
- ✘ Specifika anotací SYN2005
- ✘ Co se skrývá za označením slovní druh X.*

TOKENIZACE, AUTOMATICKÁ ANOTACE, D[EI]SAMBIGUACE

- ✘ Tokenizace – rozdělení textu na jednotky, s nimiž se dále pracuje při strojovém zpracování PJ.
- ✘ Automatická anotace – automatická morfologická analýza – slovník (word/lemma/tag) – je obecně víceznačná.

TOKENIZACE

- × grafické slovo
- × slova se spojovníkem
- × spřežky
- × zkratky

MORFOLOGICKÁ HOMONYMIE – VÍCEZNAČNOST FORMY

× Zdraví je velký dar.

ZDRAVÍ

- × *zdraví/zdraví/NNNS1.**
- × *zdraví/zdraví/NNNS2.**
- × *zdraví/zdraví/NNNS3.**
- × *zdraví/zdraví/NNNS4.**
- × *zdraví/zdraví/NNNS5.**
- × *zdraví/zdraví/NNNS6.**
- × *zdraví/zdraví/NNNP1.**
- × *zdraví/zdraví/NNNP2.**
- × *zdraví/zdraví/NNNP3.**
- × *zdraví/zdraví/NNNP4.**
- × *zdraví/zdraví/NNNP5.**

ZDRAVÍ

- ✘ *zdraví/zdravý/AAMP1.**
- ✘ *zdraví/zdravý/AAMP5.**
- ✘ *zdraví/zdravit/VB-S---3P.**
- ✘ *zdraví/zdravit/VB-P---3P.**
- ✘ *zdraví/zdravět/VB-S---3P.**

JE

× *je/být/VB-S---3P.**

× *je/on/PPXP4—3.**

× *je/on/PPNS4—3.**

VELKÝ

- × *velký/velký/AAIS1.**
- × *velký/velký/AAIS4.**
- × *velký/velký/AAIS5.**
- × *velký/velký/AAMS1.**
- × *velký/velký/AAMS5.**

DAR

- × *dar/dar/NNIS1.**
- × *dar/dar/NNIS4.**
- × *dar/dar/NNIS5.**

D[EI]SAMBIGUACE

- × Zjednotnění – volba kontextově správné varianty.
- × stochastické metody
- × pravidlové metody
- × hybridní metody

ZDRAVÍ JE VELKÝ DAR.

- × *zdraví/zdraví/NNNS1.**
- × *je/být/VB-S---3P.**
- × *velký/velký/AAIS1.**
- × *dar/dar/NNIS1.**

ANOTACE VELKÝCH SYNCHRONNÍCH KORPUSŮ

ČNK :

- × Tzv. pražský systém založený na morfologické analýze (slovníku) J. Hajiče
- × Stochastické metody disambiguace
- × Pravidlové metody disambiguace
- × Guessery/hadače

STRUKTURA ZNAČKY

[HTTP://UCNK.FF.CUNI.CZ/BONITO/ZNACKY.PHP](http://ucnk.ff.cuni.cz/bonito/znacky.php)

- ✘ Každá značka je řetězcem 16 znaků (16. pozice chybí pouze v korpusech SYN2000 a ORWELL).
- ✘ Značka je konstruována tak, aby každá pozice odpovídala jedné morfologické kategorii podle víceméně tradičního lingvistického pojetí.
- ✘ Každé hodnotě v dané kategorii odpovídá jeden znak, převážně písmeno velké abecedy (např. 'P' pro plurál, neboli množné číslo), výjimečně i jiný znak (např. 'f' pro infinitiv, nebo ',' pro podřadicí spojky).
- ✘ Hodnota, která nedává smysl (např. pád u sloves), je reprezentována znakem '-' (pomlčka).

ANOTACE MLUVENÝCH KORPUSŮ

- × Ruční
- × Není široce přístupná

ANOTACE KSK

- ✘ Upravená verze morfologického slovníku (Osolsobě 1996) a morfologického analyzátoru *ajka* (Sedláček 2004).
- ✘ Ruční disambiguace.

PRAŽSKÝ A BRNĚNSKÝ SYSTÉM ANOTACÍ

- × Projekt nové národní morfologie

SPECIFIKA ANOTACÍ SYN2005

- × Testování guesserů

CO SE SKRÝVÁ ZA OZNAČENÍM SLOVNÍ DRUH X.*

- × Slova, kterým nelze na základě morfologického slovníku přiřadit žádnou interpretaci.
- × Méně obvyklá slova.
- × Méně obvyklé tvary.
- × Překlepy.

DOPORUČENÁ ČETBA PRO ZÁJEMCE O PROBÍRANOU PROBLEMATIKU:

- ✘ Jelínek, T.: Nové značkování v Českém národním korpusu. *Naše řeč* 91, 2008, s. 13–20.
- ✘ Jelínek, T., Petkevič, V.: Systém jazykového značkování korpusů současné psané češtiny. In Petkevič, V. – Rosen, A. (eds.) 3. *Gramatika a značkování korpusů*, Praha : Nakladatelství Lidové noviny/Ústav Českého národního korpusu, 2011, s. 154–170.

DOPORUČENÁ ČETBA PRO ZÁJEMCE O PROBÍRANOU PROBLEMATIKU:

- ✘ Osolsobě, K.: Popis gramatických významů (hodnot) jednoduchých slovesných tvarů v anotacích českých (slovenských) korpusů. *SPFFBU A 55*, Brno : FF MU, 2007, s. 201–218.
- ✘ Petkevič, V.: Reliable Morphological Desambiguation of Czech: Rule-Based Approach is Necessary. In: Šimková, M. (ed.), *Insight into the Slovak and Czech Corpus Linguistics*, Bratislava : Veda, 2006, s. 26–44.

DOPORUČENÁ ČETBA PRO ZÁJEMCE O PROBÍRANOU PROBLEMATIKU:

- ✘ Petkevič, V.: Využití vidu ke zkvalitnění automatického značkování češtiny. In Bičan, A. – Klaška, J. – Macurová, P. – Zmrzliková, J. (eds.), *Karlík a továrna na lingvistiku. Prof. Petru Karlíkovi k životnímu jubileu*, Host : Brno, 2010, s. 368–387.