

Staročeská (hyper)lemmatizace



*Vznik příspěvku byl podpořen projektem Ministerstva školství, mládeže a tělovýchovy
č. [LM2015081](#) „Výzkumná infrastruktura pro diachronní bohemistiku“ (akronym RIDICS) v rámci
Projektu velkých infrastruktur pro VaVal.*



 RIDICS

Ondřej Svoboda, ÚJČ AV ČR, svoboda@ujc.cas.cz

Staročeská (hyper)lemmatizace a (nejen) morfologické značkování

- I. **určení**: staročeská textová banka (primárně)
- II. možnosti stč. lemmatizace/značkování
- III. jak vzniká stč. morfologická databáze
 - historie + generování tvarů
- IV. ukázky

Ve třech krocích

- I. generování slovních tvarů (od všech hyperlemmat)
- II. morfologická databáze (word, hyperlemma, tag, `is_lemma`, `variant_tag`)
- III. (víceznačně a částečně) anotovaný korpus

Pojmy

- „kuoním“ ⇒ **lemma** „kuoň“ (fonologicky odpovídá tvaru)
- **hyperlemma** „kóň“ (odpovídá fonologickému systému češtiny kolem r. 1300, tak jako staročeské slovníky: GbSlov, StčS, MSS, ESSČ)
- hyperlemma je zároveň jedním z lemmat: „**kóň**“, „kuoň“, „kũň“

Určení: stč. textová banka

- vzniká v oddělení vývoje jazyka od roku 2006
- korpus (ručně) transkribovaných českých textů cca. do roku 1500
- transkripce **standardizuje pravopis** (včetně interpunkce), ale zachovává fonologicky relevantní rysy jazyka
- dostupný ve Vokabuláři webovém a nově i prostřednictvím KonTextu
- více než 200 pramenů různého rozsahu, téměř 5,5 milionu tokenů (asi 200 tisíc typů)

Možnosti (automatické) anotace

- (baseline) moderní data a nástroje
- „postaršení“ novoč. morfologické databáze, „ponovočeštění“ stč. textu (Hana et al., 2012)
- lemmata pro 2000 nejčastějších tvarů (ibid.)
- **formální popis** stč. lexika (Jínová/Synková et al., 2014–)
- morf. databáze (word, (hyper)lemma, tag) a vedlejší aplikace
- (v budoucnu) značky pro **varianty/mutace** (à la NovaMorf) a trénovací desambiguovaný korpus

Stč. morfologická databáze

- Pavlína Synková: formální popis stč. apelativní deklinace (disertační práce, 2017):
 - Gebauerova **Hist. mluvnice jazyka českého** + ověřování v interní textové bance, rukopisech a jejich edicích
 - vzory (kmenotvorné přípony + koncovky) podle (psl.) kmene a rodu
 - seznam (**hyper**)lemmat (ze slovníků): vzor, alternace tvarotvorného základu, omezení paradigmatu

Stč. morfologická databáze

- Pavlína Synková: formální popis stč. apelativní deklinace (disertační práce, 2017):
 - vzory (kmenotvorné přípony + koncovky) podle (psl.) kmene a rodu
 - seznam (hyper)lemmat: vzor, alternace tvarotvorného základu, omezení paradigmatu
 - hláskové a pomocné/formální změny
 - systém + výjimky
 - vzory ve webové aplikaci

Funkce generátoru

I. (hyper)lemma **hadačka**, vzor **žena**, alternace

EO bezkonečkové tvary K_1eK_2 tvary s koncovkou K_1K_2

II. odtržení koncovky Nsg **-a** \Rightarrow tvarotvorný základ <hadačk>

III. koncovce **-a** odpovídá alternaci K_1K_2

IV. "hadačk" = "(hada) K_1K_2 " $\Rightarrow K_1=\check{c} K_2=k$

V. přiřazení alternací TZ ke koncovkám paradigmatu:

$-\emptyset$ = hada**ček**-, ostatní hada**čk**-

VI. spojení tvarotvorných základů a koncovek a aplikace

„hláskových“ změn: Lsg *hada**čk-ě** > hada**čcě** > hada**čce**

Vývoj lemmatizace

2017

2018

* apelativní substantiva (číslo a pád)

* hyperlemmata (bez desambiguace)

* atributivní (brněnské) morfologické značky

* aplikace stč. slovní tvary a KonText interně v „alfaverzi“

+ neohebné slovní druhy, příslovce a 6. slovesná třída

+ lemmata (bez desambiguace)

+ poziční (pražské, Hajičovy) značky à la korpusy SYN*

⇒ obě aplikace na webu (s manuály)

Výhled na rok 2019 a dál

- práce na slovesech (vzory **pracěvati** (substráty), **býti**, negace, 5. a 4. třída – **prošu**), později na adjektivech, zájmenech, číslovkách, propriích
- anotace agregátů (nejdříve **-s/-ť/-ž**)
- webové rozhraní k vývoji českého hláskosloví (s možností aplikace změn)

- samostatné variantní/ mutační tagy (NovaMorf) – ***kuoním:**
globalMutation**U**flexiveVariant**IEM**soundChange**í**
- desambiguace: asistovaná / pravidlová / ML

Ukázky

- kózle, dobřě, na, peskovati
- slovesa = více morfémů/částí: tematický sufix, prefixy (negace), flexe participií...
- zdrojová data: vzory. (2), parametry lemmat, hláskové změny

Skloňování

Lemma	Vzor	Omezení	
mládě ▾	subst.n.nt-kmen.kuře ▾	žádné (plné paradigma) ▾	
SG	DU	PL	
NOM	<ul style="list-style-type: none">• mládě◦ mládě	<ul style="list-style-type: none">• mládětě◦ mládětě	<ul style="list-style-type: none">• mláďata
GEN	<ul style="list-style-type: none">• mláděte◦ mláděte	<ul style="list-style-type: none">• mláďatú◦ mláďatou	<ul style="list-style-type: none">• mláďat
DAT	<ul style="list-style-type: none">• mláděti◦ mláděti	<ul style="list-style-type: none">• mláďatma• mláďatoma	<ul style="list-style-type: none">• mláďatóm◦ <u>mláďatuom</u><ul style="list-style-type: none">▪ mláďatům◦ mláďatom
ACC	<ul style="list-style-type: none">• mládě◦ mládě	<ul style="list-style-type: none">• mládětě◦ mládětě	<ul style="list-style-type: none">• mláďata
VOC	<ul style="list-style-type: none">• mládě◦ mládě	<ul style="list-style-type: none">• mládětě◦ mládětě	<ul style="list-style-type: none">• mláďata
LOC	<ul style="list-style-type: none">• mláděti◦ mláděti	<ul style="list-style-type: none">• mláďatú◦ mláďatou	<ul style="list-style-type: none">• mláďětech◦ mláďětech

Časování

kupovati	verb.6.kupovati	žádné (plné paradigma)	sl
participle			
nt			
	SG	DU	PL
NOM	<ul style="list-style-type: none">• kupujě<ul style="list-style-type: none">◦ kupuje• kupující<ul style="list-style-type: none">◦ kupující	<ul style="list-style-type: none">• kupujúce<ul style="list-style-type: none">◦ kupující• kupujúc<ul style="list-style-type: none">◦ <u>kupujíc</u>	<ul style="list-style-type: none">• kupujúce<ul style="list-style-type: none">◦ kupující• kupujúc<ul style="list-style-type: none">◦ kupujíc
slovní tvar	kupujíc		
tvarotvorný základ	kup		
zakočnění	-ujíc (PAR.NT.DU.NOM.m)		
hlásková změna	jú > jí (u-i-fronting, 2. až 3. čtvrtina 14. století)		
kde	na konci tvarotvorného základu (na posledním grafému + možná i v zakončení)		
identifikátor	kupujíc (PAR.NT.DU.NOM)		
odvozeno z	kupujúc (PAR.NT.DU.NOM)		
(atributivní)	k5gMnDc1almS		

Upravený KonText

Korpus: Staročeská textová banka 1.1.9.1 | Dotaz: kůň (1 802 výskytů) ▶ Pozitivní filtr: .*\\.* (210 výskytů)

Výskytů: 210 | i.p.m.: 38,4 (vztaženo k celému korpusu) | ARF: 89,85 | Výsledek je seříděn

1 / 2

Výběr řádků: základní ▾

1	<input type="checkbox"/>	MastDrk ↗ 3r ↗ 242	vonie a jako samé hovno	konie /koň kuoň kóň kůň/N-MP2-----	. Rubienus
2	<input type="checkbox"/>	HradSat ↗ 133v ↗ 64	, aby koval rozhlédaje a	koniem /koň kuoň kóň kůň/N-MP3-----	nezajímaje
3	<input type="checkbox"/>	PasMuzA ↗ 229	nazajtřie kázal ciesař svatého Jiříe	kuoňmi /koň kuoň kóň kůň/N-MP7-----	po všem r
4	<input type="checkbox"/>	PasMuzA ↗ 500	, stříelejí , bodú ,	koňmi /koň kuoň kóň kůň/N-MP7-----	tlačie , har
5	<input type="checkbox"/>	ŠtítKlem ↗ 8v	, a potom sú někteří	koňmi /koň kuoň kóň kůň/N-MP7-----	vláčení i v
6	<input type="checkbox"/>	ŠtítKlem ↗ 60r	svatú Štěnán nežť hv	koňmi /koň kuoň kóň kůň/N-MP7-----	vládl ; and
7	<input type="checkbox"/>	DaIV ↗ 3r			jezditi oby
8	<input type="checkbox"/>	DaIV ↗ 3r			dievky jez
9	<input type="checkbox"/>	DaIV ↗ 4v			podpěchu
10	<input type="checkbox"/>	DaIV ↗ 28r			říšskému a

...
Datace pramene	2. polovina 14. století	Předlohou je	rukopis
Zkratka pramene	MastDrk	Ediční poznámka	https://vokabular.uj...
Digitální edice	https://vokabular.uj...	Heslo kóň ve slovnících	https://vokabular.uj...

1 / 2

Hláskové změny

- **cítit** > **cejtit**, **vozík** > **vozejk**, ...
- **Kí** > **Kej**; **K** = **c**, **z**, ...
- od konce **14.** století a ve století 15. (cca 1386–1500)

Vybraná bibliografie

- Hana, Lehečka, Feldman, Černá, Oliva: **Building a Corpus of Old Czech**. LREC 2012, pp. 9–15
- Jínová/Synková, Lehečka, Oliva: **Describing Old Czech Declension Patterns for Automatic Text Analysis**. Mundo Eslavo 2014, pp. 7–17
- Synková: **Popis staročeské apelativní deklinace (...)**. FF UK, disertační práce
- Synková, Lehečka, Svoboda: **Na cestě k lemmatizaci staročeských textů: data, software, aplikace**. SALI 2018, pp. 66–84