

Přegenerování a podgenerování – dva problémy automatické analýzy přirozeného jazyka, konkrétně slovo tvorby

PLIN033_3

PŘEGENEROVÁVÁNÍ

- ✘ *Formální definici (algoritmu) odpovídají jednotky, které tvoří homogenní skupinu (tu, kterou se prostřednictvím formálního zadání snažíme definovat), ale i jednotky, které jsou vůči této skupině heterogenní. Tento jev spadá na vrub obecné vlastnosti přirozeného jazyka, jíž je víceznačnost (homonymie) na všech úrovních.*

PODGENEROVÁNÍ

- ✘ *Rubem téže mince je tzv. podgenerování, tedy případ, kdy formální zadání je vymezeno příliš úzce, takže nejsou zachyceny jednotky, které se jeho prostřednictvím snažíme definovat.*

PŘÍKLADY PŘEGENEROVÁVÁNÍ Z MINULÝCH CVIČENÍ

- × *Náboženství, nádeničení, ...*
- × *Klíč, míč, ...*

POMOCÍ NÁSTROJE *DERIV A MORFIO* VYHLEDEJTE KANDIDÁTY NA ČINITELSKÁ JMÉNA NA *-TEL*

- ✗ Maskulina životná s koncovým reťazcom *tel*

Hledání

značka RE

DERIV

× seznam

1	<input type="checkbox"/>	<u>badatel</u>
2	<input type="checkbox"/>	<u>bavitel</u>
3	<input type="checkbox"/>	<u>bičovatel</u>
4	<input type="checkbox"/>	<u>bořitel</u>
5	<input type="checkbox"/>	<u>bouratel</u>
6	<input type="checkbox"/>	<u>branitel</u>
7	<input type="checkbox"/>	<u>brojitel</u>
8	<input type="checkbox"/>	<u>brzditel</u>
9	<input type="checkbox"/>	<u>buditel</u>
10	<input type="checkbox"/>	<u>budovatel</u>
11	<input type="checkbox"/>	<u>burcovatel</u>
12	<input type="checkbox"/>	<u>bydlitel</u>
13	<input type="checkbox"/>	<u>cenitel</u>
14	<input type="checkbox"/>	<u>cestovatel</u>

DERIV HLEDÁNÍ DVOJIC

× $t\$/k5.*mF > tel/k1gMnSc1$

1: značka RE

2: značka RE odvodit od 1. slova: nahrazované nahrazující

DERIV

× seznam

- 1 [badat](#), [badatel](#)
- 2 [bavit](#), [bavitel](#)
- 3 [bičovat](#), [bičovatel](#)
- 4 [bořit](#), [bořitel](#)
- 5 [bourat](#), [bouratel](#)
- 6 [brojit](#), [brojitel](#)
- 7 [brzdit](#), [brzditel](#)
- 8 [budit](#), [buditel](#)
- 9 [budovat](#), [budovatel](#)
- 10 [burcovat](#), [burcovatel](#)
- 11 [bydlit](#), [bydlitel](#)
- 12 [cenit](#), [cenitel](#)
- 13 [cestovat](#), [cestovatel](#)
- 14 [cvičit](#), [cvičitel](#)
- 15 [čekat](#), [čekatel](#)
- 16 [činit](#), [činitel](#)
- 17 [čistit](#), [čistitel](#)

MORFIO

× Seznam

	vzor 1 ▲▼ (fq ▲▼)	vzor 2 ▲▼ (fq ▲▼)
1	žít (3800)	živitel (142)
2	zvěstovat (250)	zvěstovatel (49)
3	zůstavit (38)	zůstavitel (295)
4	zřizovat (555)	zřizovatel (599)
5	zprostředkovat (910)	zprostředkovatel (605)
6	zpracovat (2947)	zpracovatel (460)
7	zplodit (485)	zploditel (11)
8	znečišťovat (143)	znečišťovatel (147)
9	zmocnit (2601)	zmocnitel (29)
10	zhotovit (889)	zhotovitel (776)
11	zásobit (204)	zásobitel (15)
12	zasílat (719)	zasílatel (51)
13	zapisovat (1400)	zapisovatel (85)
14	zajistit (12412)	zajistitel (17)
15	vyzývat (1610)	vyzývatel (51)
16	vyšetřovat (2548)	vyšetřovatel (2394)
17	vystavovat (2286)	vystavovatel (477)
18	vysílat (3721)	vysílatel (40)
19	vypisovat (351)	vypisovatel (16)
20	vykořisťovat (86)	vykořisťovatel (59)
21	vychovat (1091)	vychovatel (389)
22	vyhlašovat (834)	vyhlašovatel (107)
23	vydražit (227)	vydražitel (34)
24	volit (5757)	volitel (117)
25	vnímat (10105)	vnímatel (43)
26	věznit (520)	věznitel (148)
27	věřit (29642)	věřitel (2047)

PŘEGENEROVÁVÁNÍ

× Přit/přítel

341	<input type="checkbox"/>	<u>přisvědčovat</u> , <u>přisvědčovatel</u>
342	<input type="checkbox"/>	<u>přít</u> , <u>nepřítel</u>
343	<input type="checkbox"/>	<u>přít</u> , <u>přítel</u>
344	<input type="checkbox"/>	<u>přítakat</u> , <u>přítakatel</u>

57	rusit (3234)	rusitel (41)
58	ručit (658)	ručitel (197)
59	rozhodovat (7902)	rozhodovatel (79)
60	přít (998)	přítel (23968)
61	přihlašovat (131)	přihlašovatel (30)
62	představit (22703)	představitel (6912)

DŮVODY PŘEGENEROVÁVÁNÍ

- × Příliš široké formální vymezení
- × Nemožnost užšího formálního vymezení

PODGENEROVÁNÍ

- ✘ Kde jsou slova jako *ředitel*, *uchvatitel*, *šříitel*, *majitel*, *pisatel*, ... ?
- ✘ Zahrnutí alternací do vyhledávání jakožto prostředek zúžení definice hledaných jednotek.

DERIVAČNÍ PRAVIDLA A VÝSLEDKY PRO DERIVACI SLOVESO – DĚJOVÉ JMÉNO NA *-TEL*

-tel

substituční pravidlo	příklad		dvojice	<u>přegenerování</u>
<u>(?<lí)t\$/k5.*mF</u> >tel/k1gMnSc1 ¹	<i>učit/učitel</i>	716	716	0
<u>et\$/k5.*mF</u> > <u>itel/k1gMnSc1</u>	<i>velet/velitel</i>	20	16	4 ²
<u>[[ěí]]t\$/k5.*mF</u> > <u>itel/k1gMnSc1</u>	<i>trpět/trpítel</i>	27	22	5 ³
<u>át\$/k5.*mF</u> > <u>atel/k1gMnSc1</u>	<i>znát/znatel</i>	2	1	1 ⁴
<u>át\$/k5.*mF</u> > <u>ítel/k1gMnSc1</u>	<i>přát/přítel</i>	2 ⁵	2	0
<u>á(.)[[iě]]t\$/k5.*mF</u> > <u>a\$1itel/k1gMnSc1</u>	<i>uchvátit/uchvatitel</i>	19	18	1 ⁶

VYHLEDÁVÁNÍ DVOJIC

✘ $at\$/k5.*mF > \acute{a}\check{c}/k1gMnSc1$

1: značka RE

2: značka RE odvodit od 1. slova: nahrazované nahrazující

PŘEGENEROVÁVÁNÍ

Práce s obsahem souboru PLIN033/at_áč

Soubor **PLIN033/at_áč** byl úspěšně uložen (2 s).

Poznámkou je jednotlivý znak, jeden řádek může mít více poznámek.

Jako poznámku nelze použít dvojtečku, čárku a mezeru (budou-li zadány, budou ignorovány).

- 1 [belhat](#), [belháč](#)
- 2 [brnkat](#), [brnkáč](#)
- 3 [brumlat](#), [brumláč](#)
- 4 [bukat](#), [bukáč](#)
- 5 [česat](#), [česáč](#)
- 6 [hafat](#), [hafáč](#)
- 7 [hmatat](#), [hmatáč](#)
- 8 [chmatat](#), [chmatáč](#)
- 9 [klepetat](#), [klepetáč](#)
- 10 [kokrhat](#), [kokrháč](#)
- 11 [kopat](#), [kopáč](#)
- 12 [kovat](#), [kováč](#)
- 13 [krkat](#), [krkáč](#)
- 14 [orat](#), [oráč](#)

KLEPETÁČ

× Slovník

SSJČ Slovník spisovného jazyka českého

klepetáč

-e m. expr. *velký rak*: chytil k-e za hlavu, aby nestříhl (Pujm.); --> expr. zdrob. **klepetáček**, -čka m. (mn. 1. -čci, -čkové 6. -čcích)

KRKÁČ

× slovník

ssjc Slovník spisovného jazyka českého

***krkáč**

-e m. expr. *lakomec, chamtivec, krkoun* (Herrm., Šlej.)

DŮVODY

- × Polyfunkčnost prostředku (-á-č x -áč)
- × Závisí na mimojazykových znalostech
- × Obtížně se formálně definuje

PODGENEROVÁNÍ

- × Nedostatky ve formální definici
- × Nepravidelnosti (*vozač, trubač*)
- × Jednotky nejsou zachyceny ve slovníku
- × Jednotkám nezachyceným ve slovníku chybí interpretace na úrovni lemmatu a morfologické značky

MORFIO

- × *kout/kouč, klít/klíč, sálat/salač*
- × Propria: *máchat/Machač, tykat/Tykač, dědit/Dědič, pískat/Piskač, kopat/Kopač, klapat/Klapač, kovat/Kovač, pleskat/Pleskač, bílit/Bilič*

TYPY PŘEGENEROVÁVÁNÍ

- ✘ hláskové alternace **kořenového vokálu** u derivátů od sloves **III. třídy** podle kmene přítomného (vzor *krýt*)
- ✘ hláskové alternace **kořenového vokálu** u ostatních tříd a vzorů
- ✘ hláskové alternace **kmenotvorného vokálu** u ostatních tříd a vzorů

ALTERNACE KOV U DERIVÁTŮ SLOVES PODLE **KRÝT**

- × *hrát/hráč*
- × *chcát/chcáč*
- × ? *pít/píč*
- × ? *pět/pěč*
- × ? *sít/síč*

V KORPUSECH LZE NAJÍT (SYN)

- × *pít* (čaj)/ *čajpíč*
- × *žit/žič*
- × *! šít/šič*

" Po kafi mám blijavku , soudruzi ! Máte čaj . . . ? Já jsem . . . **čajpíč** ! Kdo je tady eště . . . čaj . . . ? "

Reflex, č. 45/2002	, že spadáme do kategorie	žičů /žičů/X@-----	. Jako je někdo optimista
Reflex, č. 45/2002	flegmatik , tak my jsme	žiči /žiči/X@-----	. A když už se
Reflex, č. 45/2002	už se jednou člověk takovým	žičem /žičem/X@-----	narodí , tak jím nikdy
Blesk, 17. 1. 2000	měsíc . " Rekvafikaci na	šiče /šiče/X@-----	mi nabídli , když jsem
Blesk, 17. 1. 2000	Loni bylo na rekvafikační kurs	šiče /šiče/X@-----	, který zajiřtuje zmíněná firma
Blesk, 18. 5. 2000	zájem o práci šiček a	šičů /šičů/X@-----	, kterých tady pracuje sto
Udá ní franta DNES 20. 9. 2000	uř více než sto řtuřet	šičů /šičů/X@-----	šiček. Těto nátek

A KROMĚ TOHO U NEŽIVOTNÝCH MÁME

- × *bít/bič*
- × *rýt/rýč*

VŠIMNĚME SI DVOJIC

- × vyprá*á*vět/vypr*a*věč | vyprá*á*věč
- × vyjedná*á*vat/vyjedn*a*vač | vyjedná*á*vač

IJP

- × <http://prirucka.ujc.cas.cz/?id=730#nadpis14>
- × **2 Střídání krátkých a dlouhých samohlásek při tvoření slov**
- × Příklady nikoli pravidla (?seznamy výjimek)

LITERATURA

- ✘ OSOLSOBĚ, Klára. *Morfologie českého slovesa a tvoření deverbativ jako problém strojové analýzy češtiny*. 1. vyd. Brno: Masarykova univerzita, 2011. 220 s. Spisy FF MU v Brně č. 401. ISBN 978-80-210-5565-0.
- ✘ CVRČEK, Václav: *Co je nového v ČNK II. KORPUS – GRAMATIKA – AXIOLOGIE 7/2013*, 95-97.

ÚKOL NA 30.10. 2013

- ✘ Pomocí nástrojů *Deriv* a *morfio* vyhledejte kandidáty na trojice sloveso-činitelské jméno na -č – ženský protějšek na -čka (sloveso – {jméno prostředku na -č} – jméno prostředku na -čka).
- ✘ Popište případy přegenerování popř. podgenerování