

**CJBB84**  
**Morfologie a korpus**

**7.30-9.00 G13**

**II.**

# Triviální a netriviální vyhledávání substantiv podle vzoru

- Morfologické značkování neobsahuje informace o skloňovacím typu
- Je možné vyhledat v korpusech lemmata skloňovaná podle určitého vzoru?

# Formální vlastnosti substantiv a vzor

- Zakončení a slovní druh (*kos, sál, pila, žeň, chudě*)
- Slovníkové informace: lemma, rod, tvar genitivu
- Rod a forma - zakončení (*kroj x zbroj, kůň x tůň*)
- 1. a 3. pozice (pražský systém)
- Atributy k a g (brněnský systém)

# Substantiva (N/k1)

- Rod (tag: pozice 3/g)
- Zakončení lemmatu

# Slovní formulace

- Najděte substantiva skloňovaná podle vzoru *žena*.
- 1. substantiva
- 2. feminina
- 3. lemma končí na *a*

# Cql dotaz

- [tag="NNF.\*" & lemma=".\*a"]

vysoko nad jejich hlavami zašustili jako badmintonový míček , sklopili	<b>letky/letka/NNFP4-----A-----</b>	a vystoupali do nějaké jir
Představoval si jeho tělo bez života a krev vytékající do	<b>vody/voda/NNFS2-----A-----</b>	s octem . A s lítostí pom
? Kolašin potřebuje leccos . Vodu , kanalizaci , zdroje	<b>elektřiny/elektřina/NNFP4-----A-----</b>	a tepla . Musíme se post
není břesk hesel , ale náročná kázeň a skromnost .	<b>Cifry/cifra/NNFP4-----A-----</b>	bych seškrstal , uklízečky
neodolají mé horlivosti a stáhnou se pryč . Slunce barví	<b>krajinu/krajina/NNFS4-----A-----</b>	rezavě rudě . Namlouvám
trochu sebrat , nemůže si přijít a po tak dlouhé	<b>době/doba/NNFS6-----A-----</b>	se vnutit k Petře na návš
krajany pracující v ČR , poslední dobou pracovala u montážní	<b>firmy/firma/NNFS2-----A-----</b>	. Měla osmiletou dceru .
. A sedmnáctka , jménem Goethova , přináší pohled na	<b>lavičku/lavička/NNFS4-----A-----</b>	vytesanou do skály , kde
zavolaly hasiče , dívku našli až policejní potápěči po několika	<b>hodinách/hodina/NNFP6-----A-----</b>	. Do bývalého kaolinových
sleva žádný vliv a jen 1 % z nich se	<b>slevám/sleva/NNFP3-----A-----</b>	záměrně vyhýbá . Špatná
až 240 km/h , 60metrový volný pád , restaurace ,	<b>prodejny/prodejna/NNFP1-----A-----</b>	suvenýrů , expozice o his

# Počet lemmat (podle frekvence) skloňovaných podle vzoru žena v korpusu SYN2015

Frekvenční limit:

**Celkem:** 20043 (401 str.)

	<b>Filtr</b>	<b>lemma</b>	<b>Frekvence</b>
1.	p/n	doba	108 001
2.	p/n	ruka	89 653
3.	p/n	strana	79 801
4.	p/n	žena	75 329
5.	p/n	hlava	72 080
6.	p/n	cesta	68 611
7.	p/n	voda	62 399
8.	p/n	hodina	56 758
9.	p/n	cena	52 219
10.	p/n	Praha	52 029
11.	p/n	řada	50 694
12.	p/n	škola	50 466
13.	p/n	firma	49 113
14.	p/n	koruna	45 414
15.	p/n	otázka	39 932
16.	p/n	kniha	39 619
17.	p/n	skupina	39 589
18.	p/n	rodina	38 284
19.	p/n	matka	37 123
20.	p/n	změna	36 652

# Lemmata rozpoznaná automatickou morfologickou analýzou

- Nerozpoznaná lemmata
- [tag="X.\*" & lemma=".\*([ayěeubdfghklmnpqrstvz] | ou | ách | á m | ami)"]



# Výsledky

Výskytů: 861 434 | i.p.m. 0: 7 134,1 (vztaženo k celému "omezeni/syn2015") | ARF 0: 363 169,21 | Výsledek je promíchán

1 / 21 536

Výběr řádků: základní

<input type="checkbox"/>	 <a href="#">Procházka amazonským pr...</a>	cesty na nějaký náhradní transport . Samotné soužití s domácími	<a href="#">castañeros/castañeros/XX</a> -----	bylo nac
<input type="checkbox"/>	 <a href="#">Aleje české a moravské kra...</a>	) v Ilford u Londýna ( 1848 ) nebo v	<a href="#">Olmstedově/olmstedově/X@</a> -----	plánu pi
<input type="checkbox"/>	 <a href="#">Sport</a>	výběrem legendy Zinedina Zidana ! Pražský tým se jmenuje FC	<a href="#">Hunters2Dreams/hunters2dreams/X@</a> -----	a je tvoi
<input type="checkbox"/>	 <a href="#">Pátý elefant</a>	přišla odpověď , rozhlédli se oba trollové kolem , spatřili	<a href="#">Tračníka/tračníka/X@</a> -----	a pomal
<input type="checkbox"/>	 <a href="#">K moři</a>	, Bára jí telefonuje a Matti jí nedovolí za ní	<a href="#">zaject/zaject/X@</a> -----	, dokud
<input type="checkbox"/>	 <a href="#">Svět kuchyní</a>	) , cena 3 320 Kč , SIKO 7/ Komfortním zařízením	<a href="#">Hot/hot/XX</a> -----	Water D
<input type="checkbox"/>	 <a href="#">Skřipavý smích Jeana Anou...</a>	. Generál přijímá návštěvy své bývalé milenky , herečky Mélusine	<a href="#">Melita/melita/X@</a> -----	, a jejíhc
<input type="checkbox"/>	 <a href="#">Reflex</a>	. Překvapila mě síla mé tety Hany , doma říkáme	<a href="#">Handulinky/handulinky/X@</a> -----	, Hnátov
<input type="checkbox"/>	 <a href="#">Křížácké zlato</a>	co to je . Jedna dávná legenda vyprávěla o obrech	<a href="#">Kablunatech/kablunatech/X@</a> -----	opásaný
<input type="checkbox"/>	 <a href="#">Ohniska napětí v postkoloni...</a>	, „ etnických čtvrtí “ , obývaných Somálci , Afary ,	<a href="#">Oromy/oromy/X@</a> -----	, Amhar
<input type="checkbox"/>	 <a href="#">Předkové</a>	pozice je v zásadě pozitivní a kritická , zdůrazňuje spíše	<a href="#">oportunismudávných/opor tunismudávných/X@</a> -----	hominin

# Frekvenční analýza

Frekvenční distribuce

strana 1

Frekvenční limit:

**Celkem:** 297434 (5949 str.)

	<u>Filtr</u>	<u>lemma</u>	<u>Frekvence</u>
1.	p/ n	the	8 288
2.	p/ n	cz	6 654
3.	p/ n	of	5 560
4.	p/ n	la	3 555
5.	p/ n	in	3 448
6.	p/ n	and	2 619
7.	p/ n	et	2 463
8.	p/ n	le	1 914
9.	p/ n	San	1 875
10.	p/ n	Los	1 711
11.	p/ n	Czech	1 617
12.	p/ n	Tour	1 405
13.	p/ n	com	1 339
14.	p/ n	Open	1 258
15.	p/ n	der	1 214

# Zjednodušení dotazu

- [tag="X.\*" & lemma=".\*([ayěeu] | ou | ách | ám | ami)"]

390 279 | i.p.m. 3 232,16 (vztaženo k celému "omezeni/syn2015") | ARF 174 380,31 | Výsledek je promíchán

řádků: základní ▾

Mladá fronta Dnes	, " řekl Váňa České televizi . V sedle třináctiletého	Tiumena/tiumena/X@-----	abso
Technik	zatížení . Panel řídicího systému RCS 05 . Kvalitní MIG	Double/double/XX-----	Pulse
Pavilon z oblaků	kouř . „ Mám hlad , “ dal se slyšet	Marume/marume/X@-----	. „ To
Mladá fronta Dnes	při každém dalším průchodu turnikety bude porovnávat obsluha . Firma	Melida/melida/X@-----	převz
Johann Sebastian Bach	novou formu ( putující chorálová melodie v kompaktní větě „	Jesu/jesu/X@-----	, mei
Lidové noviny	portrét , “ dodal Mikulka . Jakub . Vítězný portrét	Spodinka/spodinka/X@-----	POS
Velká kniha římských detek...	" Bohové ! " vykřikl Arpocras . Jediným hrozivým pohledem	Pudenta/pudenta/X@-----	umlič
Marianne	Průzračná , čistá a vonící až kovově . Laine de	verre/verre/X@-----	, Ser
Vesmír	časové období . 1 ) Bagemihl B . , Biological	Exuberance/exuberance/X@-----	. Ani
Marianne	povrch leštěný , keramický stěp , 1967 Kč/m2 , Mondo	Ceramica/ceramica/X@-----	Doko
Nymburský deník	z řecké ligy , ale i Eurocupu a Euroligy .	Zasvou/zasvou/X@-----	dosa
Maxim	a možná i o tom napíšeme . Naprosto skvělý je	Range/range/X@-----	Rove
Deník ctihodné prostitutky	, zašli si na kávičku nebo na zákusek k „	Vasilii/vasilii/X@-----	“, za
Johann Sebastian Bach	Zelenky a Giovanniho Alberta Ristoriho i koncertního mistra Johanna Georga	Pisendela/pisendela/X@-----	. Bac

# Příliš mnoho dat

Frekvenční limit:

**Celkem:** 147526 (2951 str.)

	<b>Filtr</b>	<b>lemma</b>	<b>Frekvence</b>
1.	p/n	the	8 288
2.	p/n	la	3 555
3.	p/n	le	1 914
4.	p/n	Santa	1 017
5.	p/n	Jeanie	1 000
6.	p/n	me	857
7.	p/n	du	822
8.	p/n	Zoey	780
9.	p/n	One	713
10.	p/n	die	618
11.	p/n	League	590
12.	p/n	mobile	574
13.	p/n	Daily	539
14.	p/n	Core	530
15.	p/n	Energy	527

# Ještě zjednodušíme

- `[tag="X.*" & lemma="(ách|ám|ami)"]`

(vztaženo k celému "omezeni/syn2015") | ARF [0: 1 130,19](#) | Výsledek je promíchán

1 / 61 ▶

, roste hojně v druhotných vysazených lesích , zejména v	<a href="#">akátinách/akátinách/X@-----</a>	, indikuje d
, 382 stran LEVICOVÝ INTELEKTUÁL . China Miéville patoí k	<a href="#">hvizdám/hvizdám/X@-----</a>	současné f
let života . Připojuji se k početným gratulantům a upřímně	<a href="#">blahoželám/blahoželám/X@-----</a>	. Vyprošuji
hroznu pinotu s kamenitým a nesmírně oživujícím , sametovým ,	<a href="#">kyselinkami/kyselinkami/X@-----</a>	však vyváží
módy za neuvěřitelně nízké ceny – Primark , známý českým	<a href="#">shopperkám/shopperkám/X@-----</a>	hlavně z Lo
: " Ráno s Piotrem Brožynou ve Wilanowě a v	<a href="#">Lazienkách/lazienkách/X@-----</a>	. Piotr je je
okolí , můžete si objednat „ vitaminové balíčky “ ve	<a href="#">freshbedýnkách/freshbedýnkách/X@-----</a>	, které dorá
ještě nevyprchal . Ten smrad mě dovádí k šílenství .	<a href="#">Hadami/hadami/X@-----</a>	se vrátí z p
kteří báječně osvěží váš interiér ! Tip Květináče s kvetoucími	<a href="#">pokojovkami/pokojovkami/X@-----</a>	umístíte d
: s primární aminokyselinou modrofialový produkt Ruhemanova violeť , s	<a href="#">imonoskupinami/imonoskupinami/X@-----</a>	pak vytváří
s námi sedí pod akácií před vesnicí , se jmenuje	<a href="#">Lkichami/lkichami/X@-----</a>	. U Sambur
nicméně že nehodlám odejít jen tak beze všeho . Théodore	<a href="#">Zami/zami/X@-----</a>	, kterému p
Nastoupil jsem na metro . Když jsem dorazil , vyděšená	<a href="#">Hadami/hadami/X@-----</a>	seděla na o
společně vyřeší . Už si na to zvykli . Tím	<a href="#">senám/senám/X@-----</a>	pořádalo oc

# Alespoň něco

	<b>Filtr</b>	<b>lemma</b>	<b>Frekvence</b>
1.	p/ n	Hadami	157
2.	p/ n	ami	72
3.	p/ n	ách	23
4.	p/ n	senám	18
5.	p/ n	Yami	16
6.	p/ n	ám	16
7.	p/ n	narozkám	14
8.	p/ n	Áách	14
9.	p/ n	pastrami	14
10.	p/ n	vámvámVám	12
11.	p/ n	Zami	12
12.	p/ n	klinoformami	11
13.	p/ n	mikropilotami	10
14.	p/ n	Tami	10
15.	p/ n	nami	9
16.	p/ n	serpentýnách	8
17.	p/ n	kuwách	8
18.	p/ n	Forgách	8
19.	p/ n	Zátorách	7
20.	p/ n	martenskách	7

# Vzory zjistitelné analogicky

- U kterých vzorů kombinací formy lemmatu a morfologické informace o rodu získáme jednoznačný výsledek?
- U kterých je postup složitější?

# vzory s komplikovanějším postupem

- *pán a muž*
- *hrad a stroj*
- *píseň a kost*
- *moře a kuře* POZOR na dvě grafické varianty  
*[eě]*



# Úkol na 25. 10. 2017

- Popiš postup vyhledávání substantiv skloňovaných podle vzorů *pán/muž*.
- Všímej si, jak se chovají obojetné souhlásky [*bfmpvlsz*].
- Lze mezi obojetnými souhláskami najít nějakou skupinu, která by byla typickým zakončením pouze tvrdých vzorů?
- Existují nějaké postupy, jak formálně odlišit maskulina zakončená na obojetné souhlásky, které se mohou vyskytovat v zakončení maskulin obou vzorů?