



► MOGUĆNOSTI UPOTREBE TXM
► SOFTVERA ZA ANALIZU I VIZUELNO
► PREDSTAVLJANJE DIGITALNIH
► TEKSTUALNIH KORPUSA

Jelena Jaćimović
Stomatološki fakultet Univerziteta u Beogradu

Beograd, 9. maj 2019. godine

► O ČEMU ĆE BITI REČI ...

- ▶ TXM softver
- ▶ srpELTeC korpus
- ▶ Osnovni koncepti i alati TXM-a za analizu teksta



TEKSTOMETRIJA I TXM SOFTVER

ŠTA JE TEKSTOMETRIJA?

Metodologija koja omogućava nelinearno proučavanje korpusa, u kvantitativnom i kvalitativnom smislu, kombinujući leksikometrijska i statistička istraživanja sa razvijenim korpusnim tehnologijama (Unicode, XML, TEI, NLP, CQP, R).

SVRHA TEKSTOMETRIJSKE ANALIZE

Opisivanje leksičkih, morfoloških, sintaksičkih i drugih karakteristika teksta na lokalnom nivou.

Proučavanje međuodnosa određenih elemenata teksta (zajedničko pojavljivanje), leksičkih odnosa između tekstova (specifičnosti), lociranje elemenata teksta u korpusu (progresija).

Omogućava poređenje lingvističkih hipoteza ili kreiranje novih pomoću posmatranih podataka.

Omogućava pristup tekstu ne samo linearnim čitanjem, već i ciljanim posmatranjem određenih elemenata i karakteristika, izabranih u skladu sa postavljenim pitanjima.

Kvantifikovanje elemenata nije nov pristup, ali je pojednostavljen primenom postojećih alata.

TXM SOFTVER

Grafičko korisničko okruženje zasnovano na korišćenju CQP pretraživača i R statističkog paketa, koje omogućava analizu digitalnih korpusa bilo kog jezika

Korišćenje je besplatno

Softver otvorenog koda

Kontinuirano ga razvija tim istraživača IHRIM (Institut d'Histoire des Représentaions et des Idées dans les Modernités) laboratorije, ENS de Lyon od 2010. godine

Dostupne su desktop instalacije za Windows, Linux i Mac OS X, kao i veb platforma

► NEKE KARAKTERISTIKE TXM SOFTVERA

- Širok spektar ulaznih formata
- Korpusi pisanih tekstova: TXT (Unicode) / XML /TEI
 - Jedinice teksta: tekstovi korpusa (knjige, članci, intervjuji, ...) koji mogu imati i metapodatke (autor, naslov, datum, žanr,)
 - Opcione strukture teksta: svaki tekst može da sadrži određene unutrašnje strukturne jedinice (odeljci, pasusi, upravni govor, ...) koje mogu imati određena svojstva (naslov, broj,)
 - Leksičke jedinice: svaki tekst je sastavljen od niza reči koje mogu imati određena svojstva (grafički oblik, lema, gramatička kategorija, ...)
- Korpusi transkriptata (sinhronizovani sa izvornim audio ili video snimkom)
- Paraelni korpusi

► NEKE KARAKTERISTIKE TXM SOFTVERA

- Tekstovi su kodirani u skladu sa određenim konvencijama, poput XML kodiranja prema TEI smernicama, koje se mogu koristiti prilikom analize korpusa.
- Tekst se automatski segmentira prilikom uvoza u TXM okruženje, dok se reči označavaju svojstvima, kao što su lema ili gramatička kategorija (TXM koristi TreeTagger).
- Bilo koja kombinacija ovih svojstava reči može se upotrebiti za navođenje njihovih obrazaca i prikaz konteksta u kojima se te reči pojavljuju (liste frekvencija, konkordance, vizuelni prikaz izdanja).
- Statistički modeli primjenjeni na prebrojavanje svojstava reči omogućavaju analizu njihove distribucije po korpusima (faktorska analiza, klaster analiza), njihovu izuzetno visoku ili nisku zastupljenost u pojedinim potkorpusima (analiza specifičnosti), ili analizu privlačnosti koja postoji između pojedinih reči (analiza zajedničkog pojavljivanja).
- Rezultat svake analize može da se izveze u tabelarnom ili grafičkom obliku radi dalje analize i uređivanja pomoću nekog drugog alata.

KORPUS

srpELTeC

<https://distantreading.github.io/ELTeC/srp/index.html>

- ▶ Andra Gavrilović. Prve žrtve (kratka proza), 1893.
- ▶ Tadija Kostić. Gospoda seljaci (kratka proza), 1896.
- ▶ Čedomilj Mijatović. Rajko od Rasine (kratka proza), 1892.
Ikonija, vezirova majka (pričevica), 1891.
- ▶ Milan Milićević. Deset para (pričevica), 1881.
- ▶ Bora Stanković. Uvela ruža (pričevica), 1899.
- ▶ Milutin Uskoković. Potrošene reči (pričevica), 1911.
Došljaci (roman), 1910.
Čedomir Ilić (roman), 1914.
- ▶ Jelena Dimitrijević. Nove (roman), 1912.
- ▶ Dragutin Ilić. Hadži Đera (roman), 1904.
- ▶ Milica Janković. Kaluđer iz Rusije (roman), 1919.
- ▶ Lazar Komarčić. Dragocena ogrlica (roman), 1880.
Jedna ugašena zvezda (roman), 1902.
Prosioci (roman), 1905.
- ▶ Branislav Nušić. Opštinsko dete (roman), 1902.
- ▶ Isidora Sekulić. Đakon Bogorodičine crkve (roman), 1919.

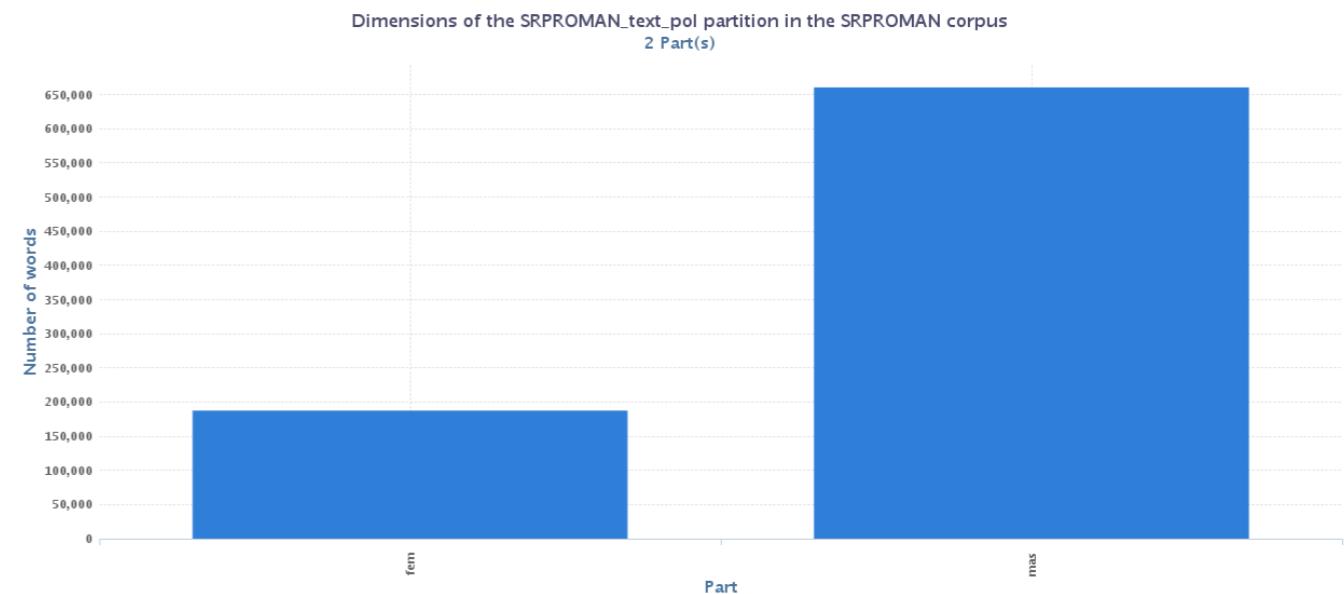
KORPUS

Modul XML - TEI Zero + CSV (uvoz u TXM)

Properties of SRPROMAN (id=SRPROMAN)

Summary Statistics

- Number of words 855412
- Number of word properties 4
- Number of structural units 15



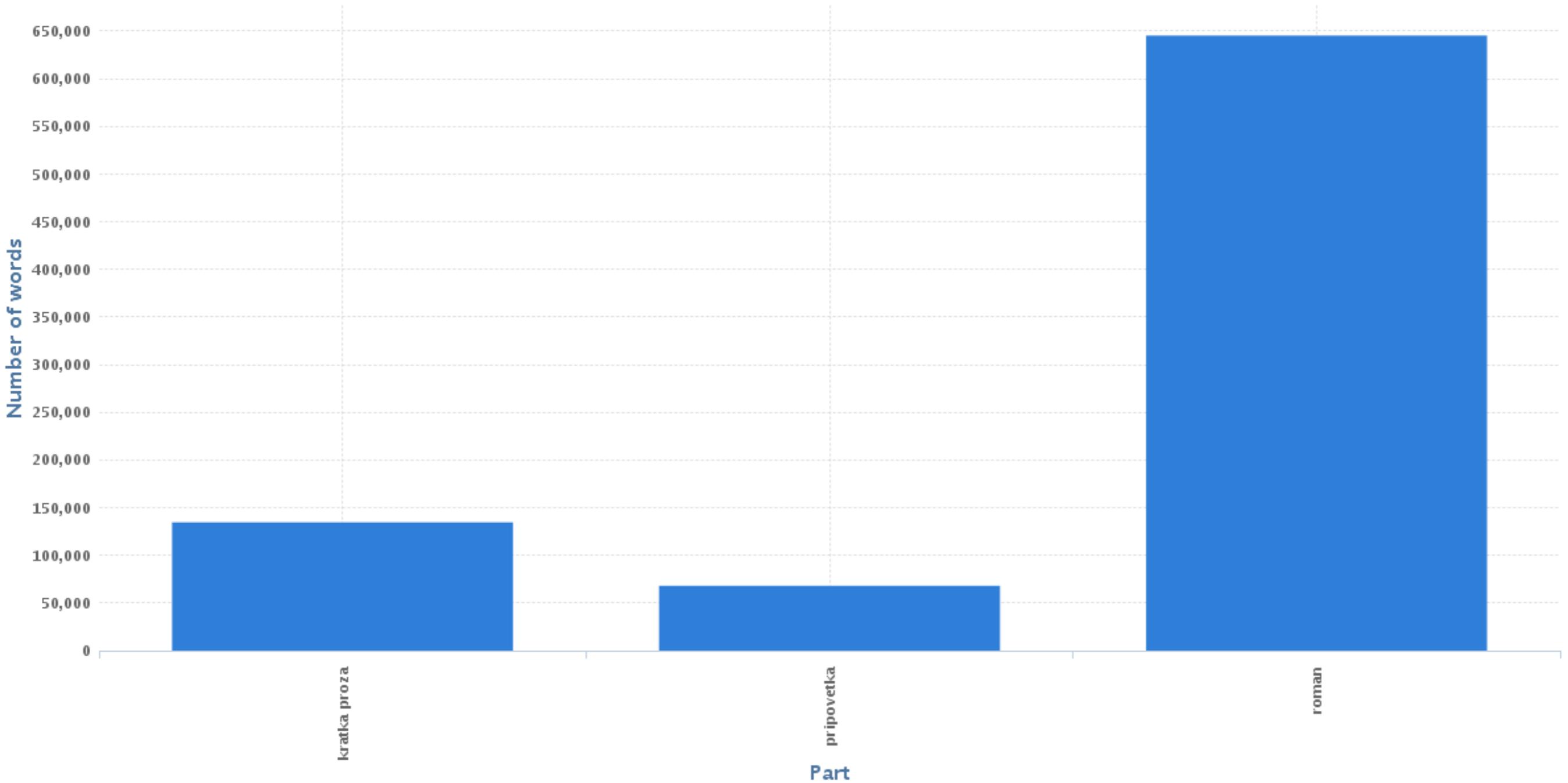
Lexical Units properties

- n : 1, 2, 3, 4, ...
- srlemma : драгоцен, огрица, :, прича, у, свой, време, ...
- srpos : A, N, SENT, PREP, PRO, PUNCT, NUM, CONJ, ABB, ?, V, ADV, ...
- word : Драгоценна, огрица, :, прича, у, своје, време, ...

Structural Units properties

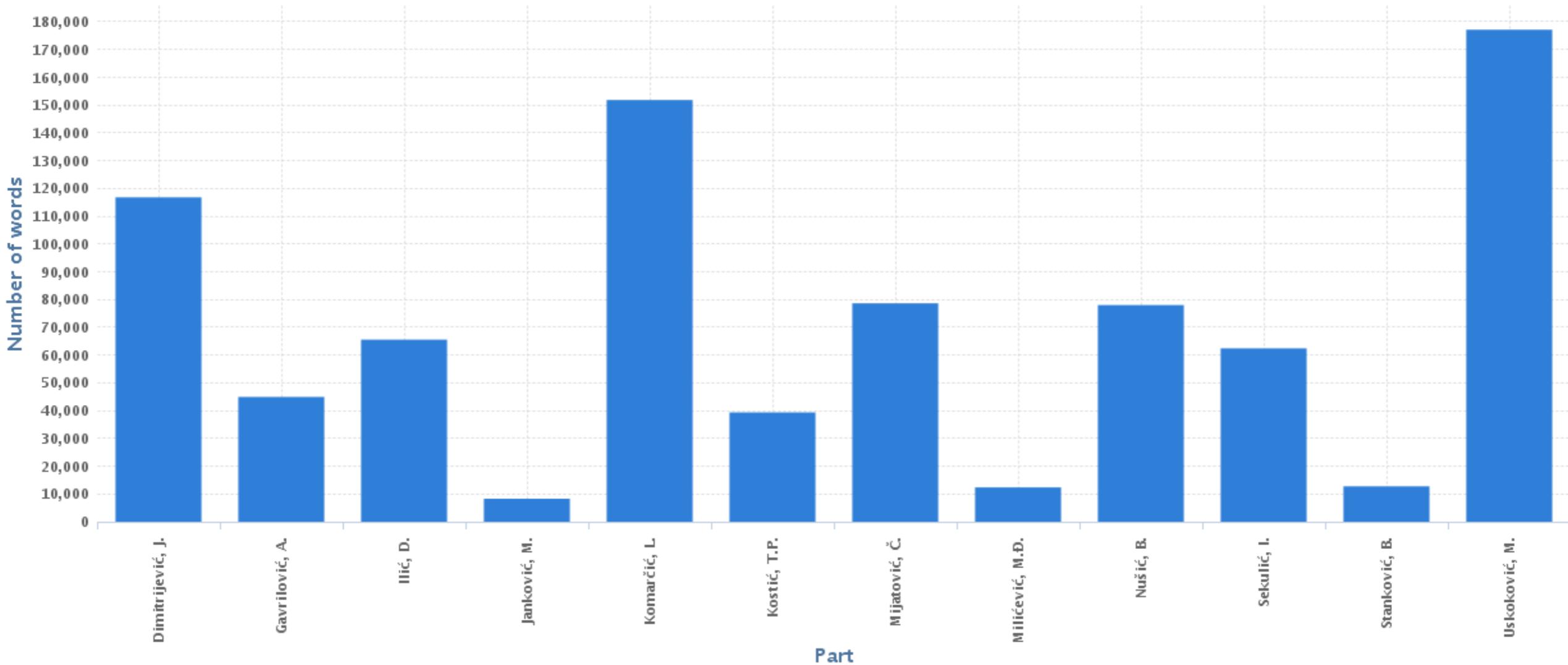
- Back, body, div, foreign, front, head, hi, l, note, p, quote, text (author, date, id, lang, pol, tip, title), title

Dimensions of the SRPROMAN_text_tip partition in the SRPROMAN corpus
3 Part(s)

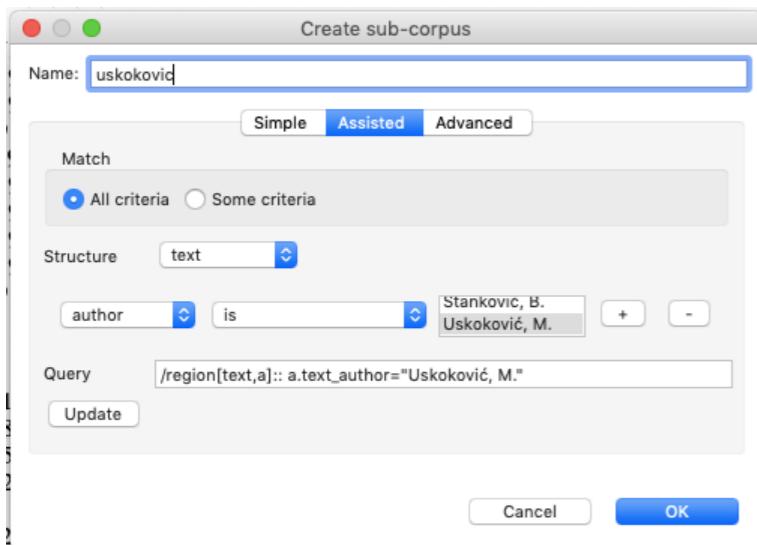


KORPUS

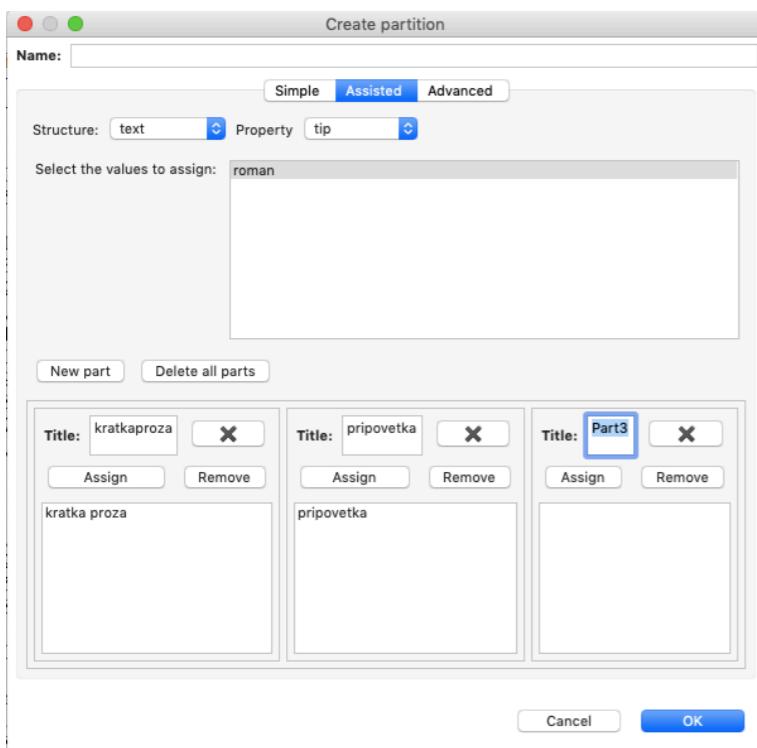
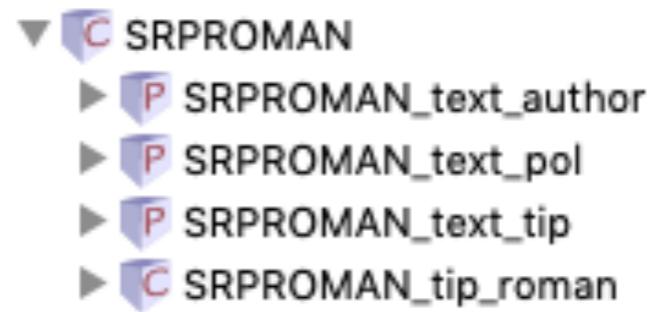
Dimensions of the SRPROMAN_text_author partition in the SRPROMAN corpus
12 Part(s)



▶ POTKORPUSI I PARTICIJE



POTKORPUSI



PARTICIJE

VIZUELNI PRIKAZ TEKSTOVA KORPUSA

SRP18800_DRAGOCENAOGRLICA

author Komarčić, L.
pol m
title Dragocena ogrlica
date 1880
tip roman

ДРАГОЦЕНА ОГРЛИЦА

ПРИЧА У СВОЈЕ ВРЕМЕ

написао

Л. Комарчић

БЕОГРАД

Штампарија Н. Стефановића и Друга

1880

- I -

ДРАГОЦЕНА ОГРЛИЦА
(ПРИЧА У СВОЈЕ ВРЕМЕ)

I

— Није нужно да вам именујем земљу, у којој су се развијали догађаји што иду, започе г. учитељ после кратког ћутања — ви ћете је сами погодити. Доста је да вам напоменем, да је она била поцепана на два три завађена тabora. Ови су, чим се који дочекавао власти, један другог гонили и прогонили. Кад су међусобни раздори и сваковрсна подметања прешла сваку меру, кад су таласи злоупотребе и насиља почели запљускивати у питома села и у мирне вароши, — онда букну револуција. У тој земљи потече братска крв. Тешки ударци револуције заљуђаше, из темеља потресоше и саме друштвене установе. Шта је невиних живота у овој борби пропало?! ... Гром је ударао у стогодишњи дуб. Овај је падао и око себе хиљадама живота смрти предавао!!

Из ове крваве борбе прве богаташке породице излазиле су с просјачким штапом.

Прича се ова односи на једног богатог племића — Артура маркиза де-Ривијера. Он беше најбогатији

1 / 242

2 / 242

12 / 242

8 / 81



▶ LEKSIKON

Lista svojstava reči koje se pojavljuju u korpusu ili potkorpusu po rastućem ili opadajućem redosledu učestalosti

Najfrekventniji oblici su separatori, oblici koji nose malo semantičkih informacija

Srednja frekvencija: centralni rečnik korpusa

Niska frekvencija i hapaks
legomenon: periferni vokabular,
hapaks čini gotovo polovinu rečnika korpusa

t 855412 , v 74371 , fmin 1 , fmax 69224

word	Frequency
,	69224
.	40343
и	26068
је	23158
се	19368
да	16974
у	13821
—	13396
на	8539
!	7947
...	6266
не	5998
а	5092
од	4800
:	4698
су	4561
што	4552
као	4488
?	4262
за	4028
па	3936
то	3751
му	3154
*	3035
:	3027
га	3003
није	2948
кад	2917

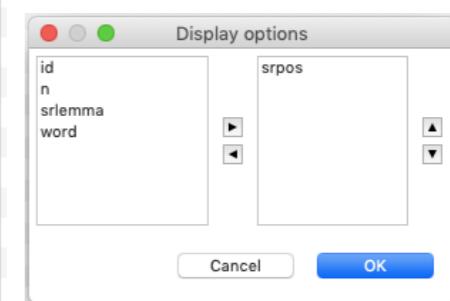
(H) (H) 1 - 100 / 74371 (H) (H)



t 855412 , v 30884 , fmin 1 , fmax 69224

srlemma	Frequency
,	69224
.	40343
јесам	35997
и	28315
се	19368
да	17750
у	15050
—	13396
тј	12453
он	11251
бити	10164
она	9188
на	9176
ја	8285
!	7947
не	6841
а	6723
...	6266
који	5681
од	5053
па	4819
хтети	4815
што	4791
:	4698
као	4691
један	4634
за	4371
оват	4277

(H) (H) 1 - 100 / 30884 (H) (H)



srpos	Frequency
N	177889
V	157514
PUNCT	94602
PRO	88762
CONJ	81837
A	59956
PREP	57451
ADV	45545
SENT	44824
PAR	32967
NUM	8707
	2272
?	1917
INT	629
ABB	486
RN	34
PREF	20

KONKORDANCE

Query	[srlemma="јесам"]	Left context	Pivot	Right context
Gospoda seljaci	канцеларији и среском мутваку. А кад га	је	, доцне, у неким пословима, који су се само ноћу	
Deset para	; а кад ми није вольја, изостајао	сам	, доцнио се, тукао сам се с ђацима, губио своје	
Čedomir Ilić	, удружењиви и безлични, ови млади људи	су	, доцније, мањом главе парламентарних режима. Други тип је био	
Nove	стара Јерменка воли га!... " Ја	сам	, драга госпођице, сто пута пожелела да свој уништим, да	
Dragocena ogrlica	развесели и сигурним гласом додаде: «Сретан	сам	, драга моја Андронито, што вам у напред казати могу,	
Dakon Bogorodičine crkve	нешто у себи. — У тешким часовима	смо	, драги Иринеје, сви без разлике само нејаки људи. Шта	
Opštinsko dete	шнајдер — ево ову. — Дајте ми	је	, драги г. Јово. — Како? Да вам је	
Gospoda seljaci	то су точкови на тој машинерији. Ми	смо	, драги мој, точкови; ми капетани и ви што сте	
Dragocena ogrlica	повикао је тронуто.... «И ти	си	, драги мој Даженоа, сигуран, да се на том тужном	
Dakon Bogorodičine crkve	зле, крвожедне и бестидне, и она	је	, држећи се у ужасу за главу, бежала испред њих с	
Čedomir Ilić	додаде госпођа после неколико тренутака. — Празна	је	, држим само неколико венаца лука. — А шта ћеш му	
Opštinsko dete	Сима — према стању ствари, оваквом каквом	је	, другим речима према фактима, ми моментално немамо деце. —	
Uvela ruža	, и очајног. Беше усталла. Зверала	си	, дрхтала, нијала си се и као хтела где да побегнеш	
Nove	!... Дај да те закитим... Кити	је	, дршћу јој руке. — Слатка суто, шта ти је	
Nove	па од радости заплака. — А где	су	, душо, моја писма? Она поцрвене, збуни се,	
Nove	горд, јунак, мушко! Али ти	си	, душо, сад слободна... Сад си распуштеница Осман-бејова, и	
Nove	истрча се са узвиком: — Шта ти	је	, душо? — Е, шта, турских жена болест —	
Prosioci	само једио просто сеоско уdomљење... И зато	смо	, ево, дошли, да, с пристанком моје властодавке и	
Hadži Đera	и да Милоја нема у кући. Ја	сам	, ево, и дошао, да Милоја уклонимо још данас,	
Ikonija, vezirova majka	жено, згрешио више него Богу! Али	јесте	, ево видим, Бог ми се свети. Његова освета није	
Rajko od Rasine	—, Није, није оче! Чула	сам	, ево и сад чујем лавеж и скамукање њихово! Еј тешко	
Kaluder iz Rusije	су говорили да је Петроград Париз, то	јест	, европска, а не руска варош. Она им је казала	
Dakon Bogorodičine crkve	ми је ова црква толико омилела, да	сам	, ето, остало у њој, иако сте Ви далеко од	
Hadži Đera	виђена и заклоњена, а сада? Сада	је	, ето, сама; сада нема коме да се потужи,	
Hadži Đera	у траг Рамову убици. А сад му	је	, ето, ушао. Ова добра вест, што је Рупић	
Ikonija, vezirova majka	, толико се изумила нисам. А да	јесам	, ех, право да ти кажем, има других соколова за	
Opštinsko dete	иде у општину. — Ама што ти	је	, жено, — ишчуђава се Роска — па сутра је петак	
Ikonija, vezirova majka	Иконија озбильно и тужно..» Теби	сам	, жено, згрешио више него Богу! Али јесте, ево	
Prve žrtve	за ловачки поход буду што свечаније. Соколови	су	, жељни плена, кликтали на дворишту, кад су осетили час	
Opštinsko dete	, ти поклони Паји писару. Перса:	Нисам	, жива ми мајка! Позови га и питај. Ја:	
Jedna ugašena zvezda	син Клеве и његова прекрасна кћер Аруџа-Дара били	су	, живели су бујним и величанственим животом, били су алем камен	
Opštinsko dete	и цео остали збор кликну: „ Тако	је	, живео! " Затим је узео реч пречасни г. архимандрит	
Čedomir Ilić	пустити се у њега. Јер какав да	је	, живот постоји, он се не може мењати у својој суштини	
Prosioci	цвете мој миљени?... „ Срећни	смо	, животе мој! " . Али се на небу наше среће	

KONKORDANCE

Query [srlemma="јесам"]

text_title	word	text: title	Pivot	Right context
Gospoda seljaci	id		је	, доцне, у неким пословима, који су се само ноћу
Deset para	srlemma		сам	, доцнио се, тукао сам се с ђацима, губио своје
Čedomir Ilić	srpos		су	, доцније, махом главе парламентарних режима. Други тип је био
Nove	back: n		сам	, драга госпођице, сто пута пожелела да свој уништим, да
Dragocena ogrlica	body: n		сам	, драга моја Андронито, што вам у напред казати могу,
Đakon Bogorodičine crkve	div: id1		смо	, драги Иринеје, сви без разлике само нејаки људи. Шта
Opštinsko dete	div: n1		је	, драги г. Јово. — Како? Да вам је
Gospoda seljaci	div: type		смо	, драги мој, точкови; ми капетани и ви што сте
Dragocena ogrlica	div: type1		си	, драги мој Даженоа, сигуран, да се на том тужном
Đakon Bogorodičine crkve	div: id		је	, држећи се у ужасу за главу, бежала испред њих с
Čedomir Ilić	div: n		је	, држим само неколико венаца лука. — А шта ћеш му
Opštinsko dete	foreign: lang		је	, другим речима према фактима, ми моментално немамо деце. —
Uvela ruža	foreign: n		си	, дрхтала, нијала си се и као хтела где да побегнеш
Nove	front: n		је	, дршћу јој руке. — Слатка суто, шта ти је
Nove	head: n		су	, душо, моја писма? Она поцрвене, збуни се,
Nove	hi: n		си	, душо, сад слободна... Сад си распуштеница Осман-бејова, и
Nove	l: n		је	, душо? — Е, шта, турских жена болест —
Prosioci	note: id		смо	, ево, дошли, да, с пристанком моје властодавке и
Hadži Đera	note: n		сам	, ево, и дошао, да Милоја уклонимо још данас,
Ikonija, vezirova majka	p: n		јесте	, ево видим, Бог ми се свети. Његова освета није
Rajko od Rasine	quote: n		сам	, ево и сад чујем лавеж и скамукање њихово! Еј тешко
Kaluder iz Rusije	text: author		јест	, европска, а не руска варош. Она им је казала
Đakon Bogorodičine crkve	text: pol		сам	, ето, остало у њој, иако сте Ви далеко од
Hadži Đera	text: project		је	, ето, сама; сада нема коме да се потужи,
Hadži Đera	text: id		је	, ето, ушао. Ова добра вест, што је Рупић
Ikonija, vezirova majka	text: date		јесам	, ех, право да ти кажем, има других соколова за
Opštinsko dete	text: lang		је	, жено, — ишчуђава се Роска — па сутра је петак
Ikonija, vezirova majka	text: tip		сам	, жено, згрешио више него Богу! Али јесте, ево
Prve žrtve	text: base		су	, жељни плена, кликтали на дворишту, кад су осетили час
Opštinsko dete	title: n		Нисам	, жива ми мајка! Позови га и питај. Ja:
Jedna ugašena zvezda	txmcorpus: lang		су	, живели су бујним и величанственим животом, били су алем камен
Opštinsko dete			је	, живео! " Затим је узео реч пречасни г. архимандрит
Čedomir Ilić			је	, живот постоји, он се не може мењати у својој суштини
Prosioci			смо	, животе мој! ". Али се на небу наше среће

Cancel OK

401 - 500 / 35997

KONKORDANCE

Query

text_title	Left context	Pivot	Right context
Gospoda seljaci	канцеларији и среском мутваку. А кад га	је	, доцне, у неким
Deset para	; а кад ми није вольја, изостајао	сам	, доцнио се, тук
Čedomir Ilić	, удружељиви и безлични, ови млади људи	су	, доцније, махом
Nove	стара Јерменка воли га!... " Ja	сам	, драга госпођиц
Dragocena ogrlica	развесели и сигурним гласом додаде: «Сретан	сам	, драга моја Анд
Dakon Bogorodičine crkve	нешто у себи. — У тешким часовима	смо	, драги Иринеје,
Opštinsko dete	шнајдер — ево ову. — Дајте ми	је	, драги г. Јово. —
Gospoda seljaci	то су точкови на тој машинерији. Ми	смо	, драги мој, точк
Dragocena ogrlica	повикао је тронуто.... «И ти	си	, драги мој Даже
Dakon Bogorodičine crkve	зле, крвожедне и бестидне, и она	је	, држећи се у ј
Čedomir Ilić	додаде госпођа после неколико тренутака. — Празна	је	, држим само не
Opštinsko dete	Сима — према стању ствари, оваквом каквом	је	, другим речима према фактима, ми моментално немамо деце. —
Uvela ruža	, и очајног. Беше усталла. Зверала	си	, дрхтала, нијала си се и као хтела где да побегнеш
Nove	!... Дај да те закитим... Кити	је	, дршћу јој руке. — Слатка суто, шта ти је
Nove	па од радости заплака. — А где	су	, душо, моја писма? Она поцрвене, збуни се,
Nove	горд, јунак, мушко! Али ти	си	, душо, сад слободна... Сад си распуштеница Осман-бејова, и
Nove	истрча се са узвиком: — Шта ти	је	, душо? — Е, шта, турских жена болест —
Prosioci	само једио просто сеоско удомљење... И зато	смо	, ево, дошли, да, с пристанком моје властодавке и
Hadži Đera	и да Милоја нема у кући. Ја	сам	, ево, и дошао, да Милоја уклонимо још данас,
Ikonija, vezirova majka	жено, згрешио више него Богу! Али	јесте	, ево видим, Бог ми се свети. Његова освета није
Rajko od Rasine	—, Није, није оче! Чула	сам	, ево и сад чујем лавеж и скамукање њихово! Еј тешко
Kaluder iz Rusije	су говорили да је Петроград Париз, то	јест	, европска, а не руска варош. Она им је казала
Dakon Bogorodičine crkve	ми је ова црква толико омилела, да	сам	, ето, остало у њој, иако сте Ви далеко од
Hadži Đera	виђена и заклоњена, а сада? Сада	је	, ето, сама; сада нема коме да се потужи,
Hadži Đera	у траг Рамову убици. А сад му	је	, ето, ушао. Ова добра вест, што је Рупић
Ikonija, vezirova majka	, толико се изумила нисам. А да	јесам	, ех, право да ти кажем, има других соколова за
Opštinsko dete	иде у општину. — Ама што ти	је	, жено, — ишчуђава се Роска — па сутра је петак
Ikonija, vezirova majka	Иконија озбиљно и тужно..» Теби	сам	, жено, згрешио више него Богу! Али јесте, ево
Prve žrtve	за ловачки поход буду што свечаније. Соколови	су	, жељни плена, кликтали на дворишту, кад су осетили час
Opštinsko dete	, ти поклони Паји писару. Перса:	Нисам	, жива ми мајка! Позови га и питај. Ја:
Jedna ugašena zvezda	син Клеве и његова прекрасна кћер Аруџа-Дара били	су	, живели су бујним и величанственим животом, били су алем камен
Opštinsko dete	и цео остали збор кликну: „ Тако	је	, живео! " Затим је узео реч пречасни г. архимандрит
Čedomir Ilić	пустити се у њега. Јер какав да	је	, живот постоји, он се не може мењати у својој суштини
Prosioci	цвете мој миљени?... „ Срећни	смо	, животе мој! " . Али се на небу наше среће

Left context size 8

Right context size 12

[navigation icons]
401 - 500 / 35997
[navigation icons]

KONKORDANCE

Query

text_title	Left context	Pivot	Right context
Gospoda seljaci	канцеларији и среском мутваку. А кад га	је	, доцне, у неким пословима, који су се само ноћу
Deset para	; а кад ми није вольја, изостајао	сам	, доцнио се, тукао сам се с ђацима, губио своје
Čedomir Ilić	, удружењиви и безлични, ови млади људи	су	, доцније, мањом главе парламентарних режима. Други тип је био
Nove	стара Јерменка воли га!... " Ја	сам	, драга госпођице, сто пута пожелела да свој уништим, да
Dragocena ogrlica	развесели и сигурним гласом додаде: «Сретан	сам	, драга моја Андронито, што вам у напред казати могу,
Dakon Bogorodičine crkve	нешто у себи. — У тешким часовима	смо	, драги Иринеје, сви без разлике само нејаки људи. Шта
Opštinsko dete	шнајдер — ево ову. — Дајте ми	је	, драги г. Јово. — Како? Да вам је
Gospoda seljaci	то су точкови на тој машинерији. Ми	смо	, драги мој, точкови; ми капетани и ви што сте
Dragocena ogrlica	повикао је тронуто.... «И ти	си	, драги мој Даженоа, сигуран, да се на том тужном
Dakon Bogorodičine crkve	зле, крвожедне и бестидне, и она	је	, држећи се у ужасу за главу, бежала испред њих с
Čedomir Ilić	додаде госпођа после неколико тренутака. — Празна	је	, држим само неколико венаца лука. — А шта ћеш му
Opštinsko dete	Сима — према стању ствари, оваквом каквом	је	, другим речима према фактима, ми моментално немамо деце. —
Uvela ruža	, и очајног. Беше усталла. Зверала	си	, дрхтала, нијала си се и као хтела где да побегнеш
Nove	!... Дај да те закитим... Кити	је	
Nove	па од радости заплака. — А где	су	
Nove	горд, јунак, мушко! Али ти	си	
Nove	истрча се са узвиком: — Шта ти	је	
Prosioci	само једио просто сеоско удомљење... И зато	смо	
Hadži Đera	и да Милоја нема у кући. Ја	сам	
Ikonija, vezirova majka	жено, згрешио више него Богу! Али	јесте	
Rajko od Rasine	— , Није, није оче! Чула	сам	
Kaluđer iz Rusije	су говорили да је Петроград Париз, то	јест	
Dakon Bogorodičine crkve	ми је ова црква толико омилела, да	сам	
Hadži Đera	виђена и заклоњена, а сада? Сада	је	
Hadži Đera	у траг Рамову убици. А сад му	је	
Ikonija, vezirova majka	, толико се изумила нисам. А да	јесам	
Opštinsko dete	иде у општину. — Ама што ти	је	, жено, — ишчуђава се Роска — па сутра је петак
Ikonija, vezirova majka	Иконија озбиљно и тужно..» Теби	сам	, жено, згрешио више него Богу! Али јесте, ево
Prve žrtve	за ловачки поход буду што свечаније. Соколови	су	, жељни плена, кликтали на дворишту, кад су осетили час
Opštinsko dete	, ти поклони Паји писару. Перса:	Нисам	, жива ми мајка! Позови га и питај. Ја:
Jedna ugašena zvezda	син Клеве и његова прекрасна кћер Аруџа-Дара били	су	, живели су бујним и величанственим животом, били су алем камен
Opštinsko dete	и цео остали збор кликну: „ Тако	је	, живео! " Затим је узео реч пречасни г. архимандрит
Čedomir Ilić	пустити се у њега. Јер какав да	је	, живот постоји, он се не може мењати у својој суштини
Prosioci	цвете мој миљени?... „ Срећни	смо	, животе мој! " . Али се на небу наше среће

● ● ● Sort options

id	word
n	i, и
srlemma	
srpos	

Cancel OK

401 - 500 / 35997

► REGULARNI IZRAZI

Regуларни израз је израз који описује скуп ниски, у складу са одређеним синтаксним правилима.

Primer

сұнародници, међународни, народске, народности, ...

одговара регуларном изразу [srlemma = ".*народ.*"]

SRPROMAN:[word="међународни"]

Query: [word="међународни"] Keyword: word Edit Search

sort keys: #1 None #2 None #3 None #4 None

|< < 1 - 1 / 1 > >| Hide settings

text_id	Left context	Keyword	Right context
SRP18960...	, господин Проко, на премер, као	међународни	човек, што међу народом живит

word	Frequency
народности	2
народске	2
Народни	1
Народну	1
међународни	1
народнога	1
народност	1
сұнародница	1
сұнароднице	1
сұнародници	1

треба да искажем у овом тренутку свога живота, почаствован од овог одличног скупа да ступим у чувену и одличну партију, којом руководи на далеко чувени Вођа и господин Шеф наш, а коју у нашем срезу и округу тако мудро предводи многопоштовани дугогодишњи посланик, господин Продан Жмурић ... И, заиста, ви, господин Проко, на премер, као међународни човек, што међу народом живите, знate колика потреба захтева да у наполу има изображени људи, који ће n:10941 srpos:A srlemma:међународан и повести путем општег народног напретка и које ја, на премер,

▶ UPITI U TXM OKRUŽENJU

- ▶ Jednostavna pretraga reči: vrednosti

Primeri

данас *народ.*

[word = “златан”] тү[гж].*

мир|рат .+патри.*

- ▶ Pretraga svojstava

[srlemma = “пријатељ”]

[srpos = “A”]

[srpos = “V” & word = игра.*]

- ▶ Pretraga niski leksičkih jedinica

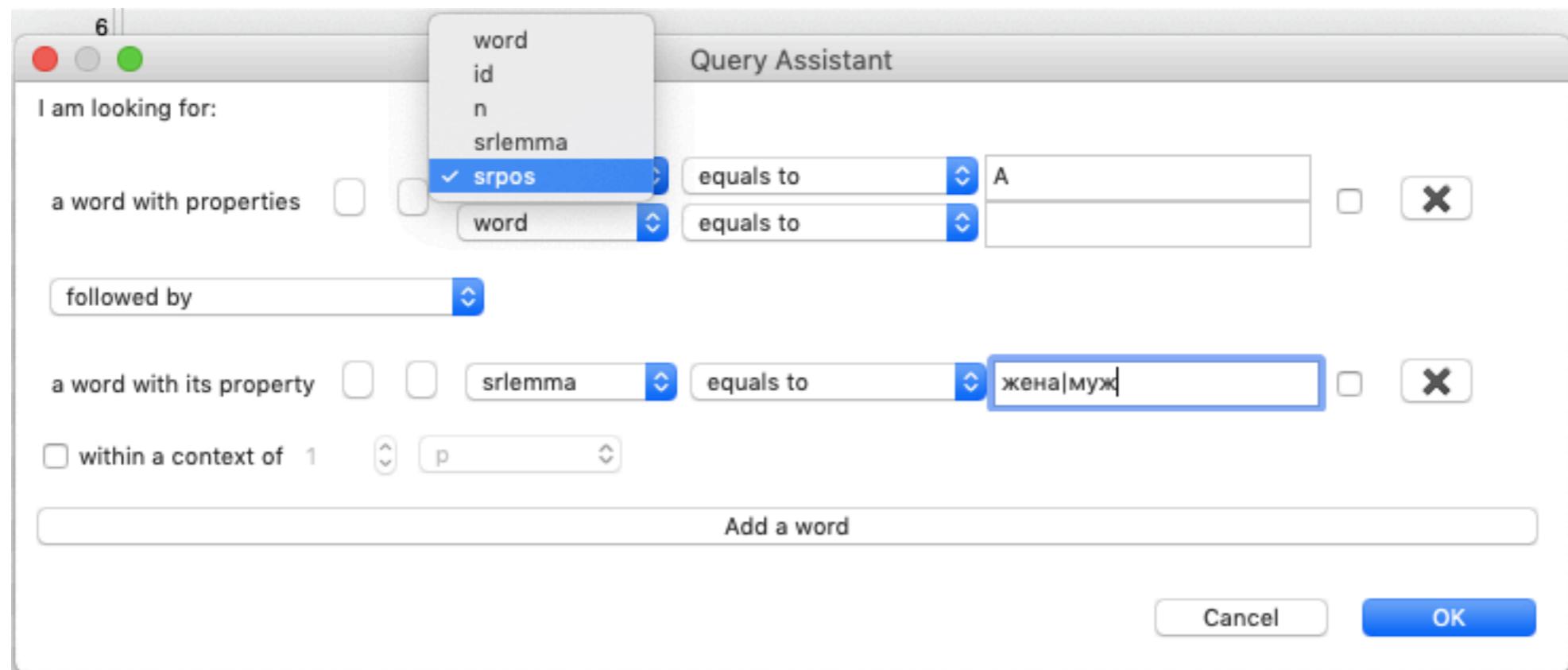
[srlemma = “тежак”] [srpos = “N”]

- ▶ Korišćenje informacija iz konteksta

[word = “срце” & _.text_pol! = “fem”]

▶ INDEKS

Lista frekvencija svojstava određenih pojavljivanja jednog CQL upita izvršenog nad korpusom, potkorpusom ili participijom



► INDEKS

Lista frekvencija svojstava određenih pojavljivanja jednog CQL upita izvršenog nad korpusom, potkorpusom ili participijom

The screenshot shows a software interface for analyzing word frequencies. On the left, there is a main table of words and their frequencies. On the right, there is a smaller table showing a subset of words from the main table, also with their frequencies. A search bar at the top of the right table allows users to filter the results.

I am looking for:

- 6
-
-

a word with property:

- followed by
- a word with its properties
- within a context of

word	Frequency
млада жена	39
турскe жене	24
младе жене	19
Млада жена	13
турских жена	8
сирота жена	7
друга жена	6
друге жене	6
младој жени	6
младих жена	4
покојног мужа	4
турска жена	4
Сирота жена	3
другог мужа	3
лепу жену	3
младу жену	3
паметна жена	3
првог мужа	3
сироту жену	3
стара жена	3
страну жену	3
турским женама	3
турску жену	3
јадна жено	3
љубљена жена	3
веселих жена	2
данашња жена	2
дивна жена	2
добра жено	2
европска жена	2
законита жена	2
лепе жене	2
лепих жена	2

word	Frequency
матора жена	2
млади муж	2
модерној жени	2
несрећна жена	2
нове жене	2
обучених жена	2
официрске жене	2
покојне жене	2
поштену жену	2
прве жене	2
просвећене жене	2
проста жена	2
старе жене	2
трећа жена	2
туђе жене	2
туђу жену	2
чиновничке жене	2
јадна жена	2
јадну жену	2
" Жене	2
» Жена	1
» жена	1
Једне жене	1
Његовој жени	1
Арифин муж	1
Добра жена	1
Многа жена	1
Многе жене	1
Примитивна жена	1
Разјарена жена	1
Стара жена	1
Старинске жене	1
—Та жена	1

Que equ eqi A OK Cancel

► INDEKS

Lista frekvencija svojstava određenih pojavljivanja jednog CQL upita izvršenog nad korpusom, potkorpusom ili participijom

The screenshot shows a software interface with three tables of word frequency data. On the left, there is a sidebar with various search parameters and a status bar indicating '6' results found. The main area contains three tables:

- Table 1 (Left):** Shows words and their frequencies. The first few rows include: млада жена (39), турске жене (24), младе жене (19), Млада жена (13), турских жене (8), сирота жена (7), друга жена (6), друге жене (6), младој жени (6), младих жена (4), покојног мужа (4), турска жена (4), Сирота жена (3), другог мужа (3), лепу жену (3), младу жену (3), паметна жена (3), првог мужа (3), сироту жену (3), стара жена (3), страну жену (3), турским женама (3), турску жену (3), јадна жено (3), љубљена жена (3), веселих жена (2), данашња жена (2), дивна жена (2), добра жено (2), европска жена (2), законита жена (2), лепе жене (2), лепих жена (2).
- Table 2 (Middle):** Shows words and their frequencies. The first few rows include: матора жена (2), млади муж (2), модерној жени (2), несрећна жена (2), нове жене (2), обучених жена (2), официрске жене (2), покојне жене (2), поштену жену (2), прве жене (2), просвећене жене (2), проста жена (2), старе жене (2), трећа жена (2), туђе жене (2), туђу жену (2), чиновничке жене (2), јадна жена (2), јадну жену (2), " Жене (2), » Жена (1), » жена (1), Једне жене (1), Његовој жени (1), Арифин муж (1), Добра жена (1), Многа жена (1), Многе жене (1), Примитивна жена (1), Разјарена жена (1), Стара жена (1), Стариинске жене (1), —Та жена (1).
- Table 3 (Right):** Shows srlemma and their frequencies. The first few rows include: млад жена (86), турски жена (43), други жена (13), сирот жена (13), леп жена (11), стар жена (11), јадан жена (8), туђ жена (6), љубљен жена (6), покојни муж (5), први муж (5), стран жена (5), добар жена (4), модеран жена (4), први жена (4), вољен жена (3), данашњи жена (3), други муж (3), европски жена (3), млад муж (3), паметан жена (3), поштен жена (3), прост жена (3), слаб жена (3), " жена (3), » жена (2), весео жена (2), диван жена (2), добар муж (2), жив жена (2), законит жена (2), заслужан муж (2), заљубљен жена (2).

► ZAJEDNIČKA POJAVLJIVANJA

Query [srlemma = "град.*"]

Parameters

Cooccurrences properties: srlemma [Edit](#) Thresholds: Fmin ≥ 2 Cmin ≥ 2 Score ≥ 2.0

Context: word structure back Use the left context Use the right context include the structure containing the pivot in the count

from - 9 ⌂ to - 0 ⌂ and from 0 ⌂ to 9 ⌂

Cooccurrent	Frequency	CoFrequency	Score	Mean distance
зидина	54	13	16	1.9
бедем	30	11	16	2.3
подграђе	16	9	15	3.4
у	15050	179	14	3.2
овај	4277	75	13	3.4
заповедник	21	7	10	1.6
властела	11	5	8	4.4
од	5053	69	8	3.2
платно	43	7	8	1.0
стража	22	5	6	.0
Љубоја	69	7	6	3.4
главни	139	9	6	4.3
кула	81	7	6	4.1
негда	36	5	5	3.8
царски	104	7	5	1.3
војска	147	8	5	5.0
град	196	9	5	4.7
рушевина	19	4	5	3.2
старешина	21	4	5	3.0
градски	47	5	4	5.6
деспотовина	25	4	4	5.2
царевина	28	4	4	6.5
муслиман	10	3	4	2.3
ударити	97	6	4	4.5

t pivot 280, v соос 132, t соос 0, Т corpus 855412

ZAJEDNIČKA POJAVLJIVANJA

Query [srlemma = "град.*"]

Parameters

Query ([srlemma = "град.*"] []* [srlemma="зидина"]) | ([srlemma="зидина"] []* [srlemma = "град.*"]) within 10

Parameters

Sort keys #1 None #2 None #3 None #4 None Sort

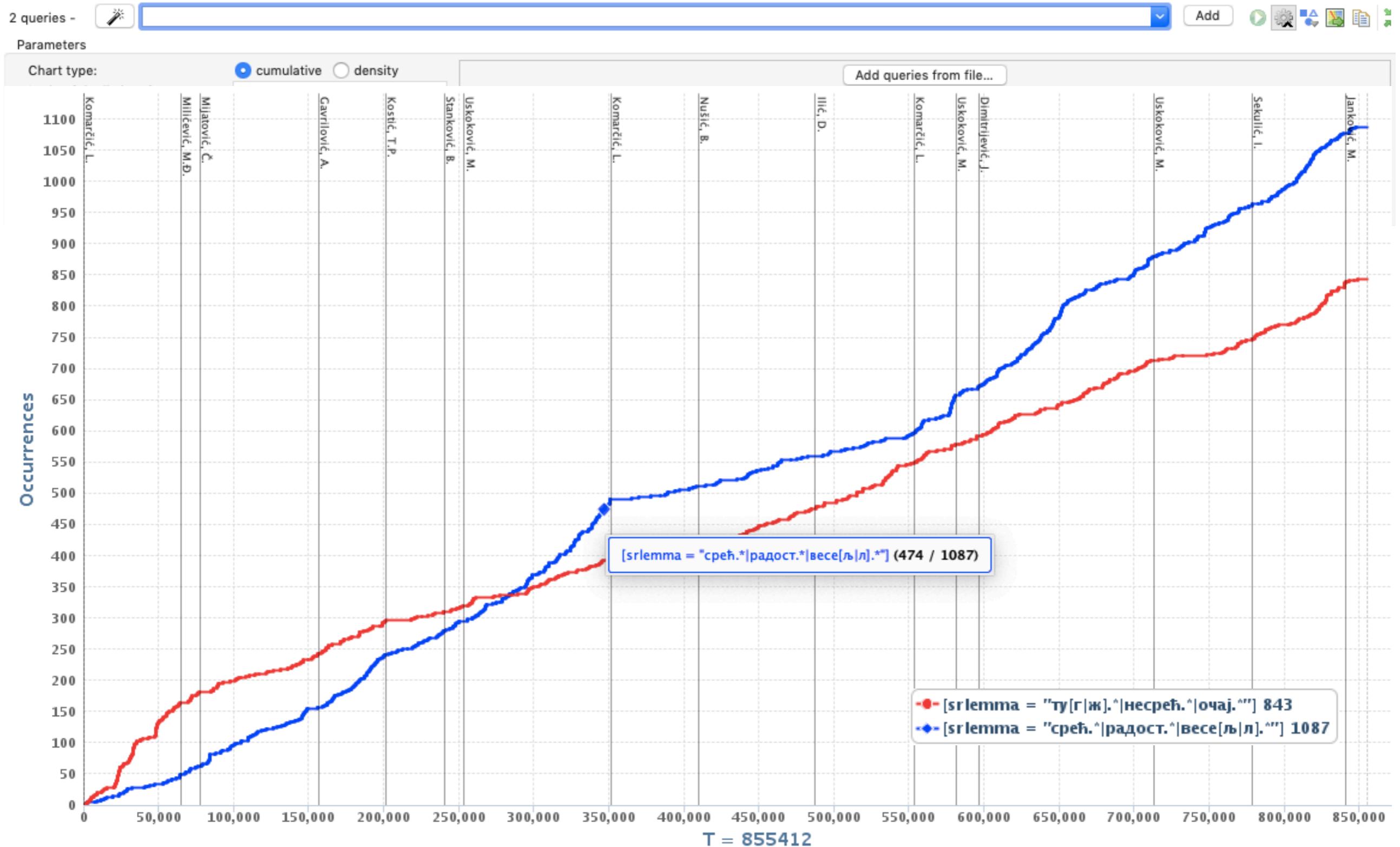
	References	Left context	Pivot	Right context
View	text_id ...	word ...	word ...	word ...
Sort	text_id ...	word ...	word ...	word ...
Size		8 ↑		12 ↑
X				
text_id	Left context	Pivot	Right context	
SRP18920	платна можете најбоље да видите, како је . 2 Вальало им је проћи поред западне	град некада имао двојаких зидина	, као оно што кажу да Цариград има! Па онда баш	
SRP18920	Вальало им је проћи поред западне зидине од	зидине од града	. Уз ту зидину, која се прилично одржала, прислоњена је	
SRP18920	опреми своју сиву ластавицу. 5 Под самим	града. Уз ту зидину	, која се прилично одржала, прислоњена је била овде и онде	
SRP18920	циновски стражари у панцирима од сребра стражарили по	градом кнеза Лазара, а прислоњена уз обрушену једну зидину	, која по свој прилици бејаше остатак од некадашњега спољњег зида градског	
SRP18920	будио успаване људе на посао. Озго испод	зидинама градским	и кроз све јунаштво своје кришом погледали к прозорима дворским од бильура	
SRP18920	један вис, на коме из далека видеше	зидина градских	зачу се куцкање чекића о наковањ. „Ево је иде!	
SRP18920	бих, окупана у сузама, заспала под	зидине од некаквог старог порушеног града	. Али се већ мрак ухватио био кад изађоше на тај вис	
SRP18920	неколике цеви, што их видех горе у	зидином оног старог града	у Крушевцу, ето ње к мени, на сну праћена још	
SRP18920	! Па онда појуре напоље, па на	зидинама од града	!“ рече Маргита; па онда настави пуна збиље и сетности	
SRP18920	сувог грања горела је на највишо и најширој	зидине од града	; нањушише на цеви оловне што воде ваздух доле, па да	
SRP18930	ловачки су се рогови пробе ради разлегали преко	зидини од разваљеног града	. То је стари пустиник давао знак најближем селу под брдом,	
SRP19000	црна као та ноћ, и у хладу	градских зидина	. Био је већ крај месеца октобра, кад и издашна јужна	
		градских зидина	које су виделе толика столећа. Гледајући се очи у очи,	

1 - 13 / 13

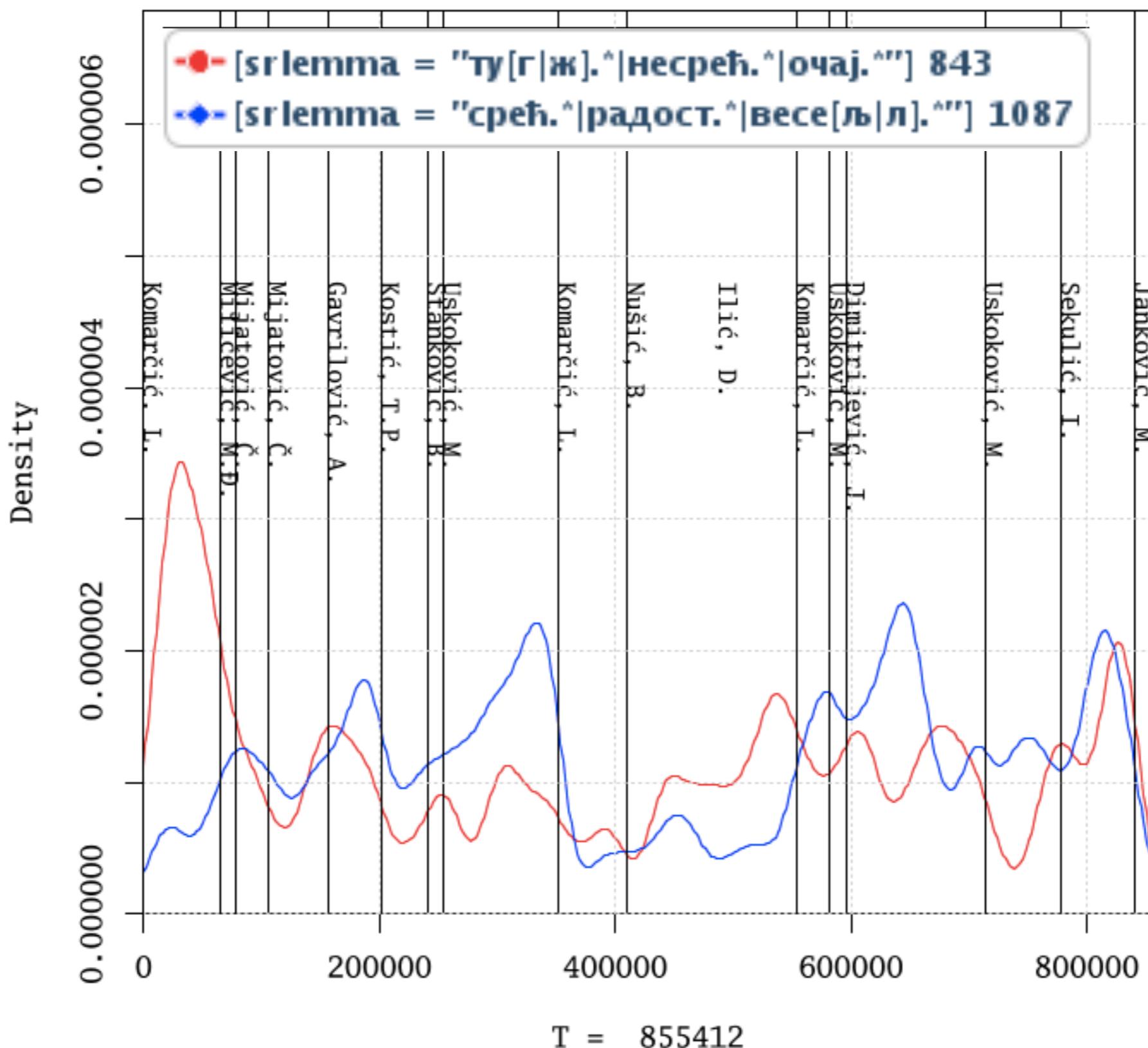
муслиман	10	3	4	2.3	
ударити	97	6	4	4.5	

t pivot 280, v cooc 132, t cooc 0, T corpus 855412

▶ PROGRESIJA



▶ PROGRESIJA



▶ LEKSIČKE TABELE

Property srlemma ⚖ ⏪ ⏩

Parameters

Thresholds

Fmin 2

Number of lines 200

Merge or Delete columns

Merge or Delete rows

srlemma	Frequency	Dimitrijević, J. t=80571	Gavrilović, A. t=27979	Ilić, D. t=42176	Janković, M. t=5751	Komarčić, L. t=101388	Koštić, T.P. t=25515	Mijatović, Č. t=51471
«	1104	0	0	0	51	575	0	355
»	1040	0	0	1	49	555	0	380
Г	771	0	0	1	0	640	23	0
Милош	482	0	0	0	0	0	0	6
Зорка	423	0	0	0	0	0	6	1
"	469	70	5	0	2	71	4	126
:	3004	750	22	129	10	977	39	409
-	1464	14	6	5	5	22	21	19
'	695	46	85	20	0	27	214	69
госпођа	449	61	9	0	0	21	28	81
господин	528	2	0	0	3	58	116	19
одговорити	627	46	53	148	13	103	24	24
овде	490	69	34	38	2	95	13	94
седети	501	134	13	52	1	26	43	48
твој	742	46	53	27	0	120	15	137
узети	553	127	17	15	2	63	17	92
година	691	109	24	26	8	221	28	50
ти	592	74	39	35	2	54	18	110
"	2390	730	2	45	0	719	127	704
те	783	98	22	103	0	96	46	104
преко	466	63	28	38	4	85	42	40
радити	438	59	13	35	4	74	40	20
код	593	118	28	19	2	91	33	12
колико	499	59	32	47	3	100	28	26
немати	430	77	21	31	4	79	25	35
наћи	581	44	33	47	9	127	15	56
соб	478	127	15	13	1	54	35	14
"	3034	802	12	44	0	871	136	806
право	535	47	27	11	3	148	11	73
мајка	785	268	14	39	3	119	10	43
три	461	109	16	21	1	79	22	56
отићи	623	173	39	52	7	88	20	41
тек	531	33	14	64	2	61	21	81
сунце	524	28	10	8	4	274	6	27
наш	1213	163	85	11	2	429	92	83
страна	543	34	65	58	1	130	14	48
	616	57	30	65	0	80	44	82

► SPECIFIČNOSTI

Model verovatnoće (Lafon, 1980) zasnovan na hipergeomatrijskoj raspodeli omogućava:

- ▶ Proučavanje distribucije učestalosti reči/svojstava u (pot)korpusu podeljenom na nekoliko delova
- ▶ Poređenje delova u smislu specifične (višak/deficit) ili osnovne upotrebe reči/svojstava

Indeks specifičnosti

znak (+/-) ako je posmatrana frekvencija veća ili manja nego u “normalnoj” distribuciji (uzimajući u obzir veličinu dela u odnosu na celinu)

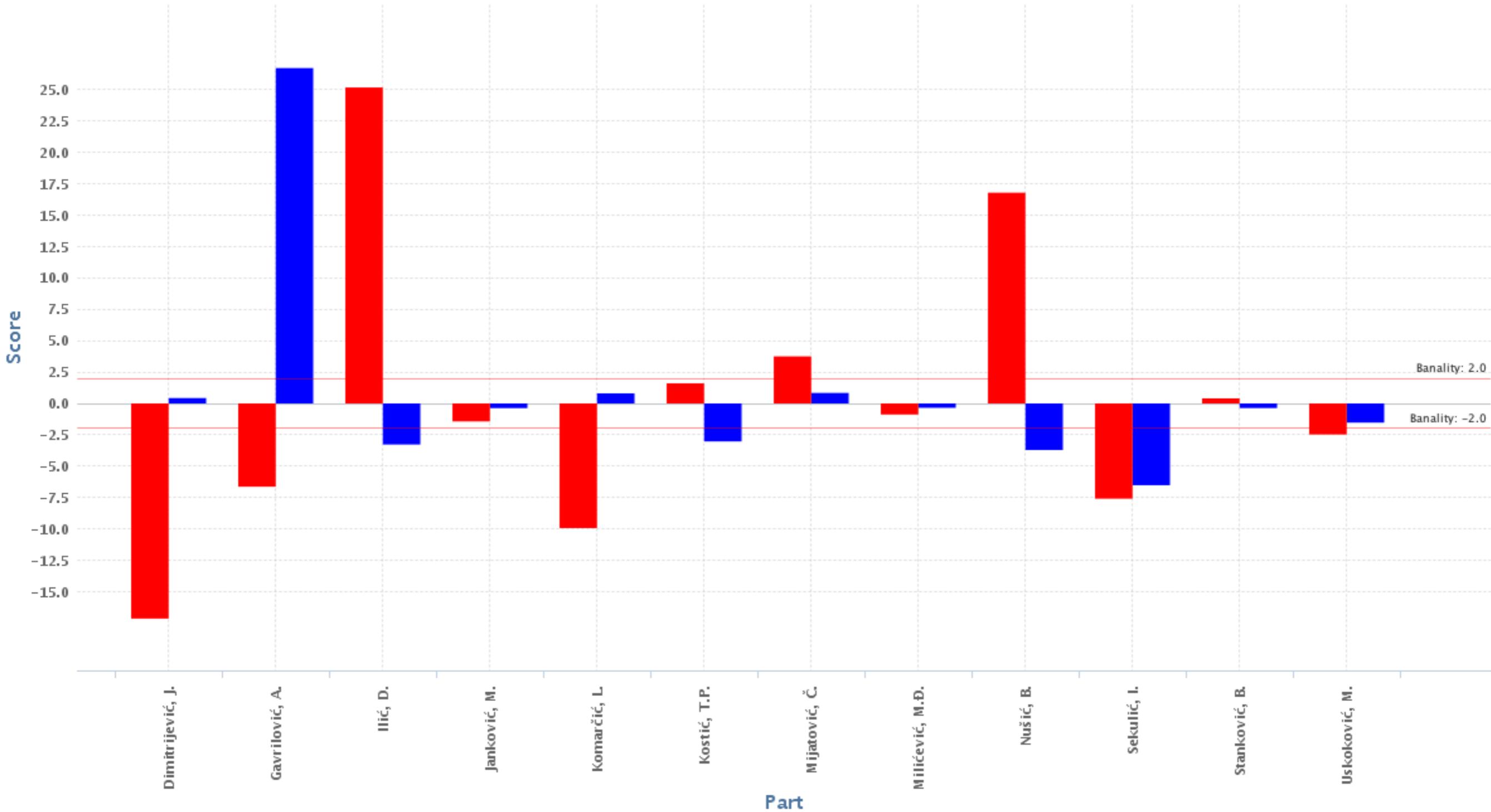
Query  [srlemma = "град|село"]

srlemma	Frequency T=847992	Dimitrijević, J. t=116782	Gavrilović, A. t=44929	Ilić, D. t=65554
село	335	3	1	90
град	196	29	58	4

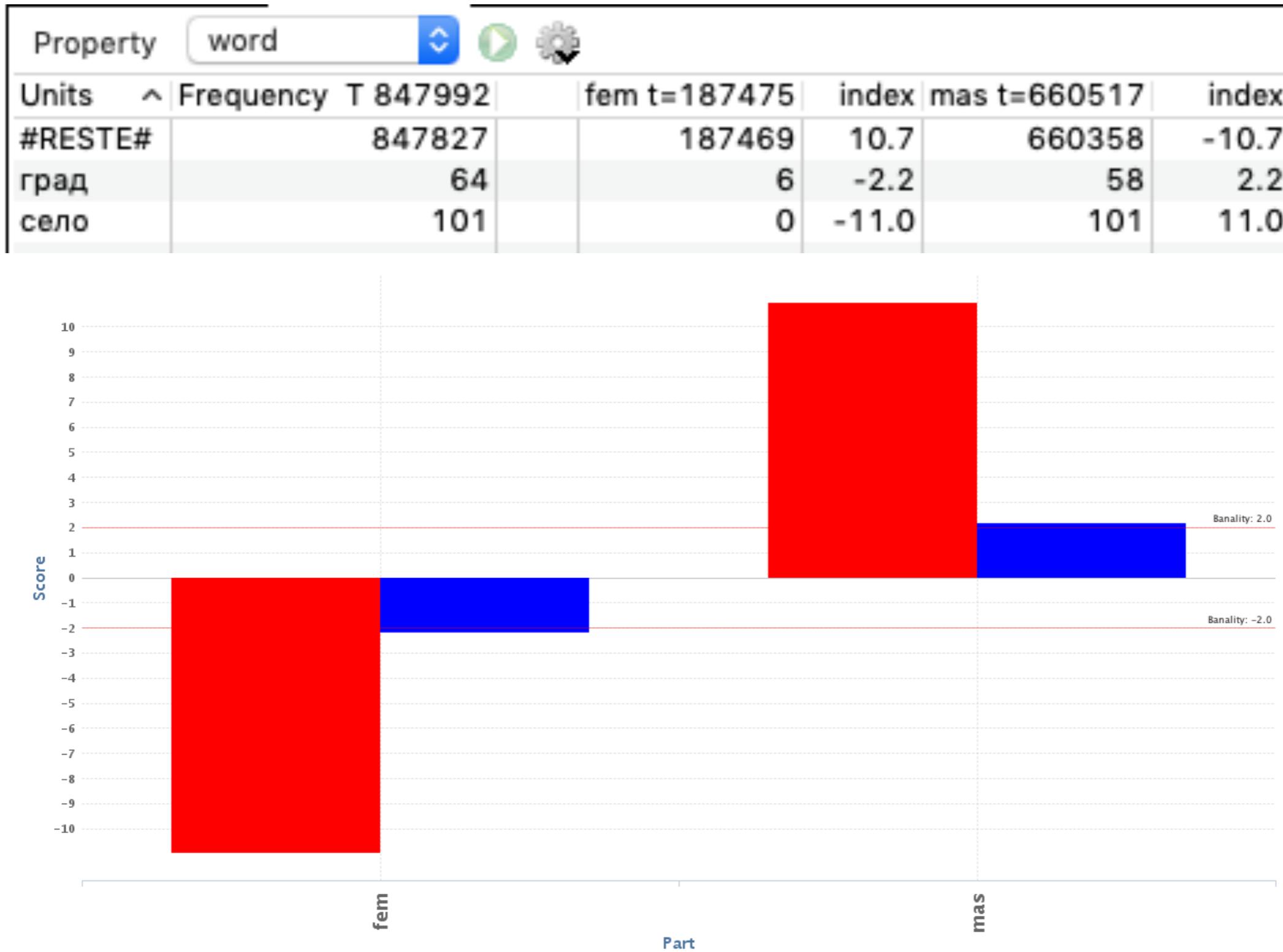
Property **srlemma**   

Units	Frequency T	847992	Dimitrijević, J. t=116782	index	Gavrilović, A. t=44929	index	Ilić, D. t=65554	index	
#RESTE#		847461		116750	8.1	44870	-7.0	65460	-13.3
град		196		29	0.4	58	26.7	4	-3.3
село		335		3	-17.1	1	-6.6	90	25.2

► GRAFIČKI PRIKAZ SPECIFIČNOSTI

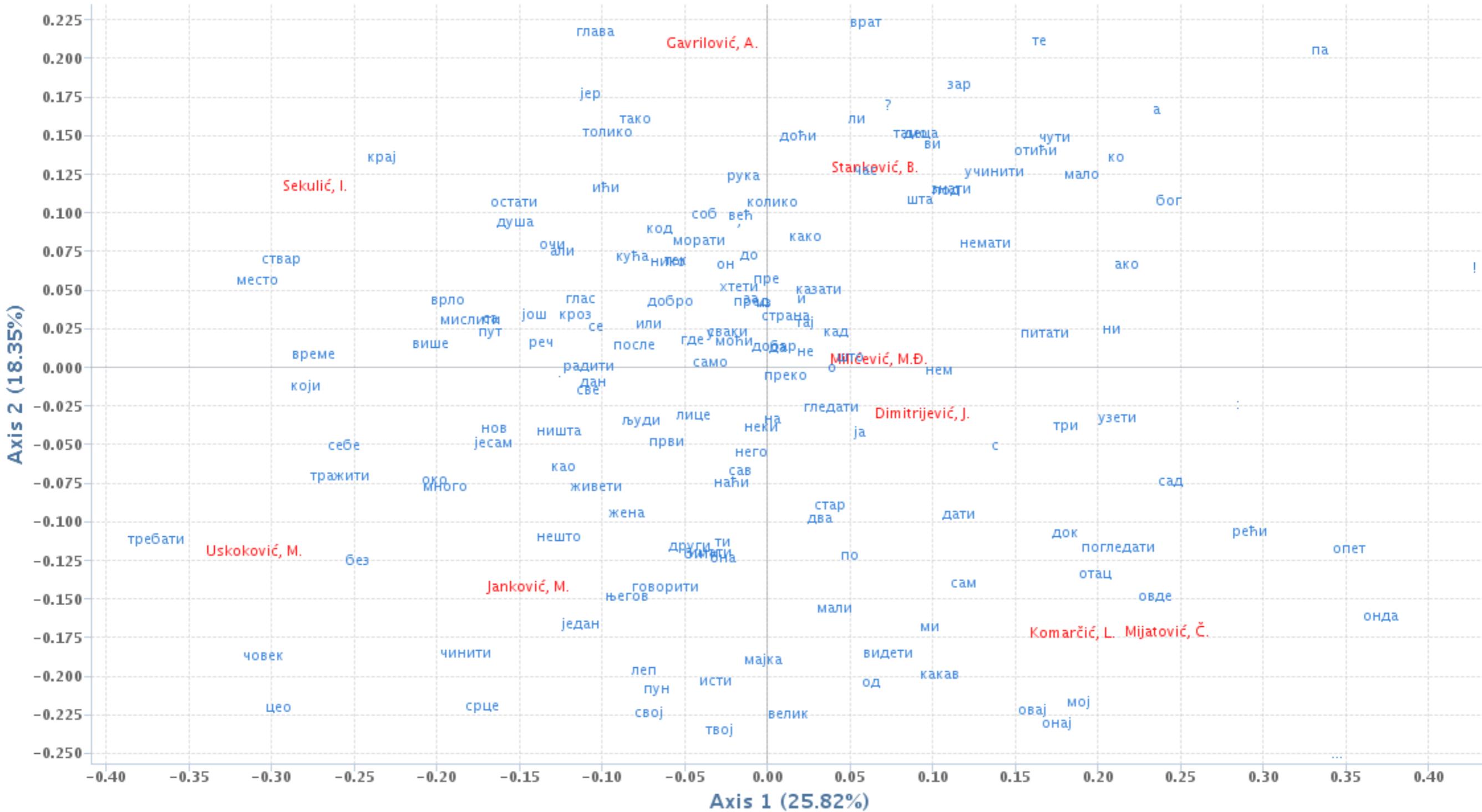


SPECIFIČNOSTI



► FAKTORSKA ANALIZA

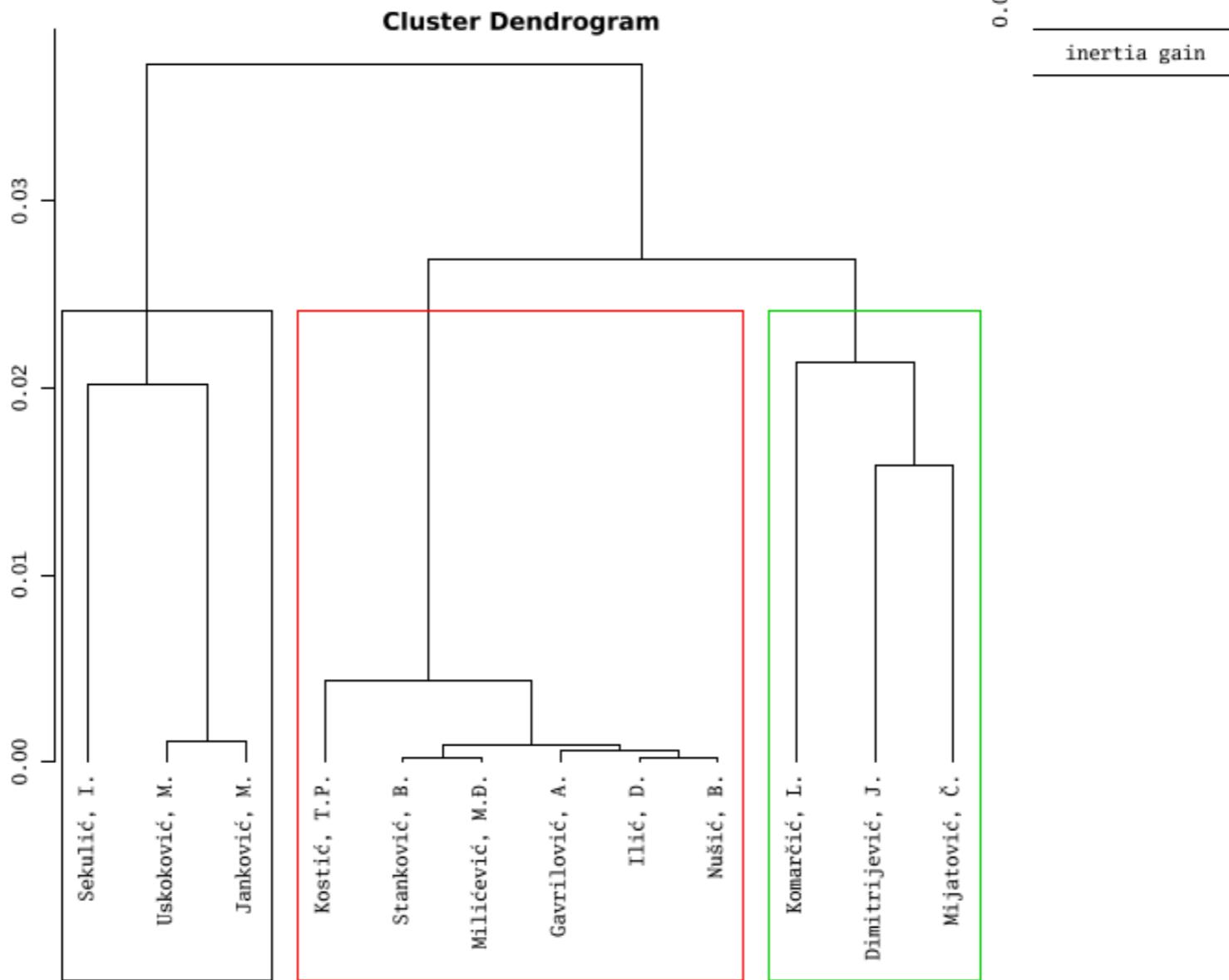
Correspondence analysis factorial plane of SRPROMAN_text_author:srlemma 2 / 200





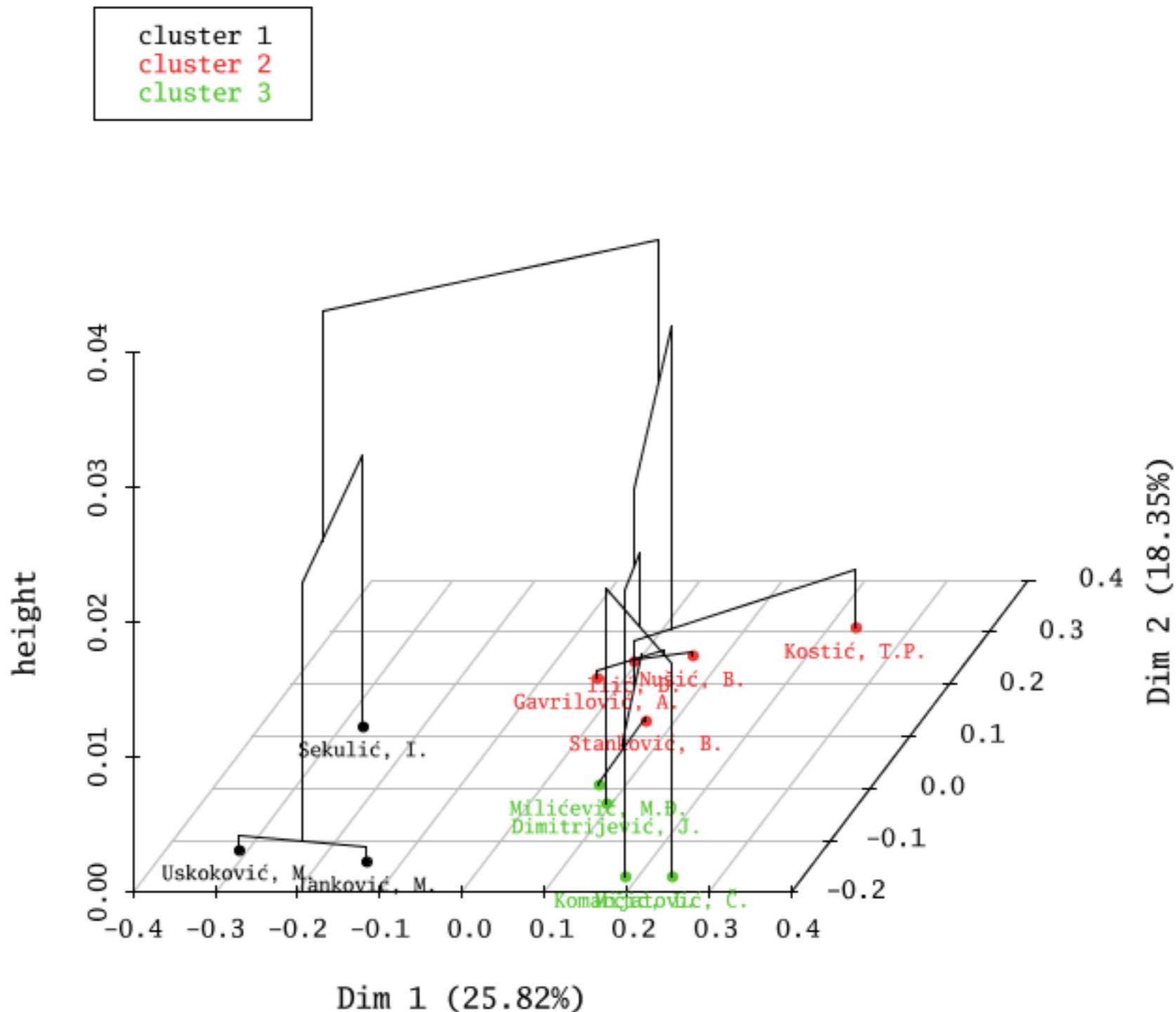
KLASTER ANALIZA

Hierarchical clustering



► KLASTER ANALIZA

Hierarchical clustering on the factor map



► ANOTACIJA KORPUSA

Nekoliko alata za anotaciju korpusa razvijeno je u okviru TXM okruženja:

- ▶ Na osnovu CQP modela korpusa
 - ▶ Anotiranje svojstava reči, koje čine konkordance, kreiranjem novih ili ispravljanjem postojećih svojstava reči
 - ▶ Jednostavno ili napredno anotiranje reči ili niski reči, koje su predstavljene konkordancama, novim strukturama tekstova i njihovim vrednostima
- ▶ Na osnovu integrisanog URS (Unit-Relation-Scheme) modela za anotiranje reči ili niski reči

▶ ŠTA BISMO JOŠ MOGLI

Usavršiti TreeTgger model za srpski jezik

Anotacija korpusa

Radionice za lingviste i istraživače iz oblasti humanističkih nauka Srbije



HVALA NA PAŽNJI!

jelena.jacimovic@stomf.bg.ac.rs

Acknowledgements

**COST Action 16204 – Distant Reading for European Literary History
STSM-CA16204-42562**

**Host institution: Institut d'Histoire des Représentaions et des Idées dans les
Modernités UMR 5317, École Normale Supérieure de Lyon, Lyon, France**