



Slovníčne analize ročno označčenega korpusa ssj500k z orodjem Q-CAT

Kaja Dobrovoljc

- 1 Center za jezikovne vire in tehnologije, Filozofska fakulteta, Univerza v Ljubljani
- 2 Laboratorij za umetno inteligenco, Inštitut Jožef Stefan

Dogodek "Kvantitativne in kvalitativne korpusne analize z novimi orodji iz nacionalnega projekta Nova slovnica sodobne standardne slovenščine: viri in metode (J6-8256)", Ljubljana, 19. 11. 2019



O (kvalitativnem) korpusnem jezikoslovju

- korpusno jezikoslovje
 - preučevanje jezika na podlagi realnih jezikovnih podatkov
 - ni samo kvantitativno
- analiza **manjšega** števila **konkretnih** primerov rabe
 - obvladljivejše delo
 - **poglobljenejša analiza**
 - upoštevanje konteksta
 - identifikacija izjem in mejnih pojavov
- idealna kombinacija obeh pristopov



Načrt predavanja

1. Predstavitev korpusa ssj500k

- 1.1. Nastanek in vsebina
- 1.2. Ravni označenosti
- 1.3. Dostopnost korpusa

2. Predstavitev orodja Q-CAT

- 2.1. Namestitev
- 2.2. Pregledovanje
- 2.3. Iskanje
- 2.4. Označevanje

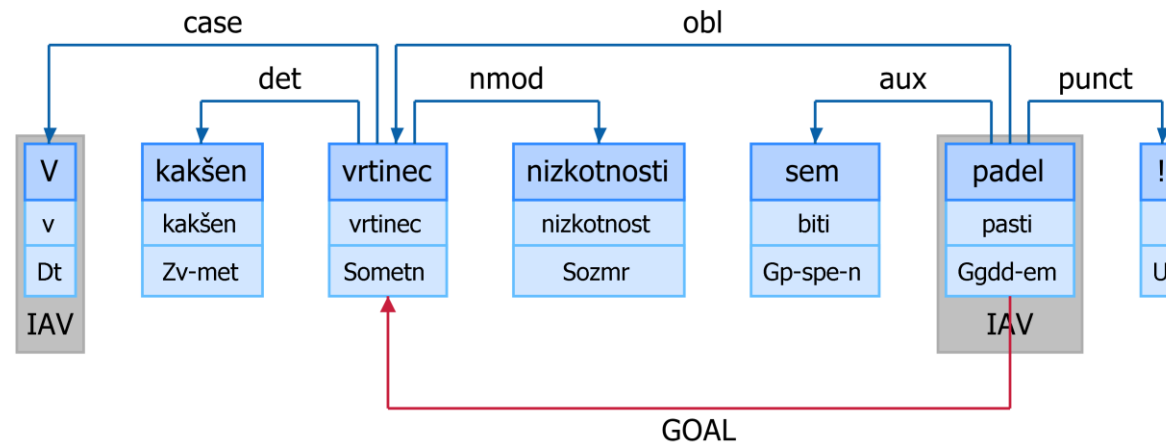
3. Zaključek

Učni korpus **ssj500k**



Kaj je korpus ssj500k

- korpus slovenskih besedil z **(ročno) pripisanimi slovničnimi oznakami** na različnih ravneh



- prosto dostopen
- približno 500.000 besed



Za kaj je uporaben

- temeljna podatkovna zbirka za **razvoj jezikovnih tehnologij**
 - osnova za nadzorovano strojno učenje
 - orodja za slovnično analizo (npr. samodejno oblikosloskladenjsko označevanje)
 - druge aplikacije (npr. samodejno prevajanje, povzemanje ...)
- pomembna tudi za **razvoj jezikoslovnega opisa** slovenščine
 - preverjanje obstoječih teorij in tipologij na avtentičnem gradivu
 - sistematična detekcija mejnih pojavov
 - razvoj novih tipologij
 - jezikoslovne analize



Zgodovina nastajanja

- več kot dve desetletji nenehnega razvoja
- 1998-: **MULTEXT-East** – razvoj prvih jezikovnih virov za oblikoskladenjsko označevanje slovenščine (specifikacije, prvi korpus – roman *1984*)
- 2007-2009: **Jezikoslovno označevanje slovenščine (JOS)** – nadgradnja smernic za lematizacijo, oblikoskladenjsko in (novo) skladenjsko označevanje, prenos na novi korpus jos100k (vzorec korpusa FidaPLUS)
- 2008-2013: **Sporazumevanje v slovenskem jeziku (SSJ)** – označenih še dodatnih 400.000 besed (ssj500k), dodane tudi imenske entitete
- 2013-2018: niz manjših projektov
 - **Označevanje udeleženskih vlog** v slovenščini in hrvaščini
 - **PARSEME**: glagolske večbesedne enote
 - **UD**: mednarodno usklajeno (obliko)skladenjsko označevanje
 - **JANES**: imenske entitete
- itd.

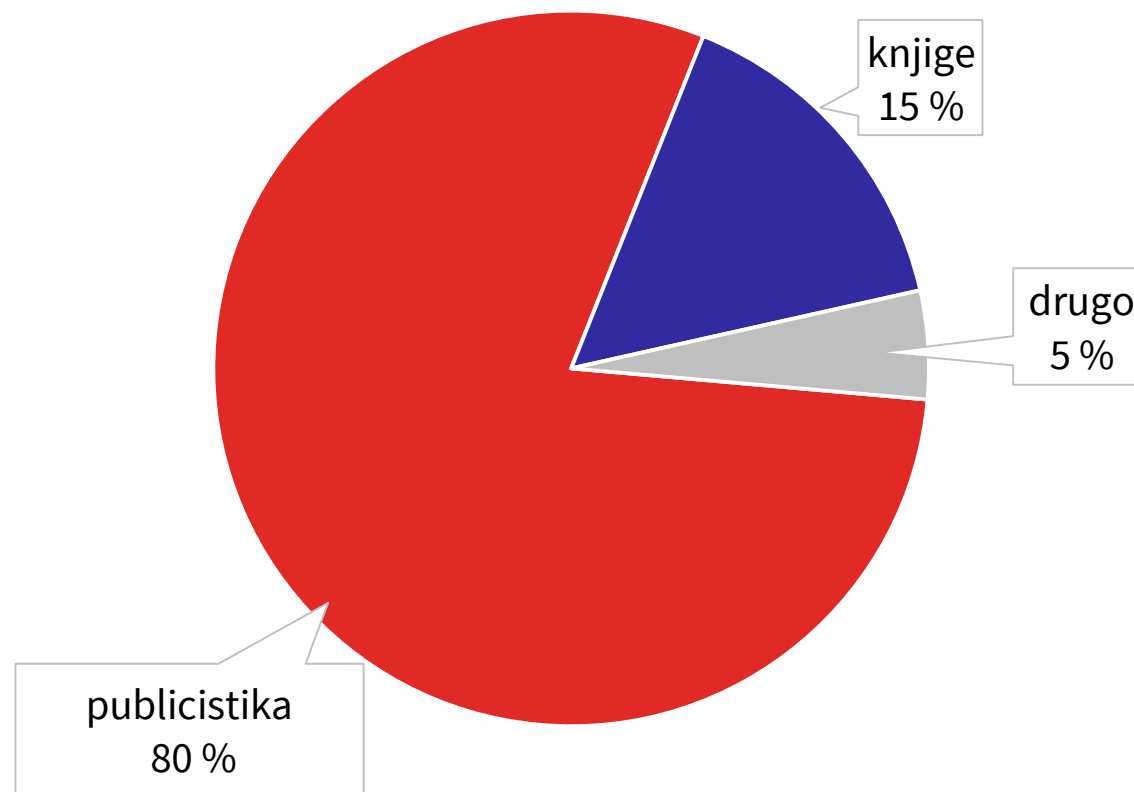
→ trenutna različica: **2.2** (Krek et al. 2019)



Besedilna zgradba

- vzorec referenčnega korpusa sodobne pisne slovenščine **FidaPLUS** (Arhar in Gorjanc 2007)

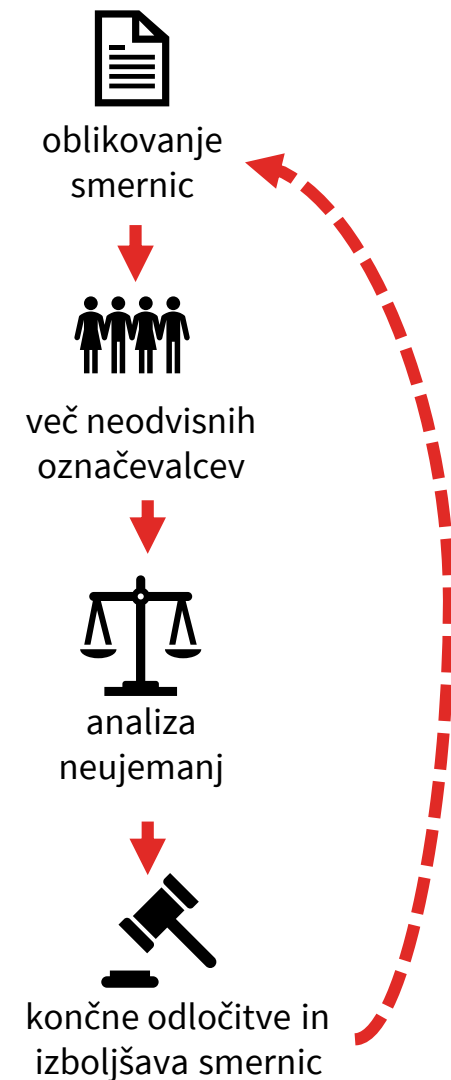
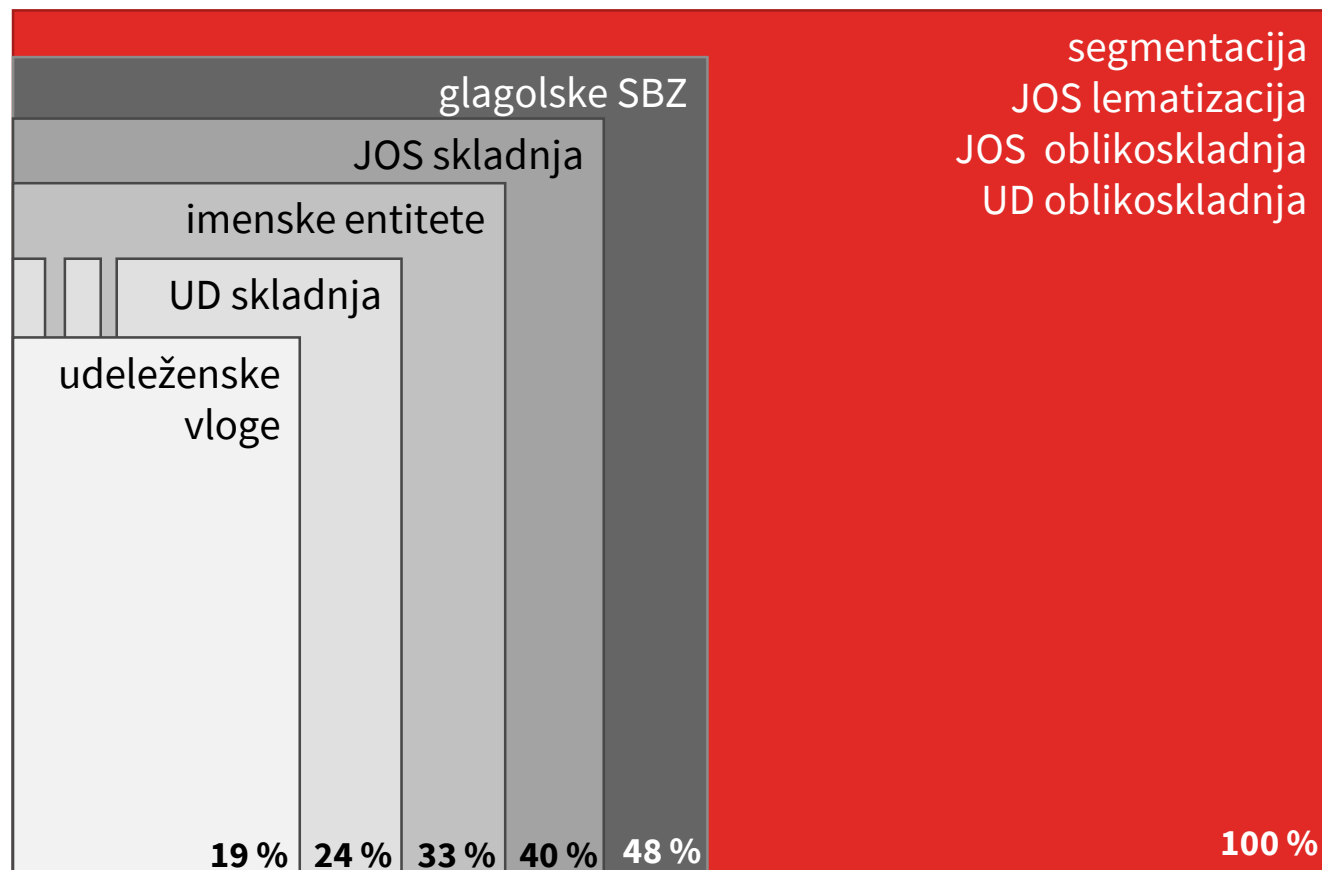
Enota	Pojavitev
Besedil	1.655
Odstavkov	8.137
Povedi	27.829
Pojavnic	586.248
Besed	500.295





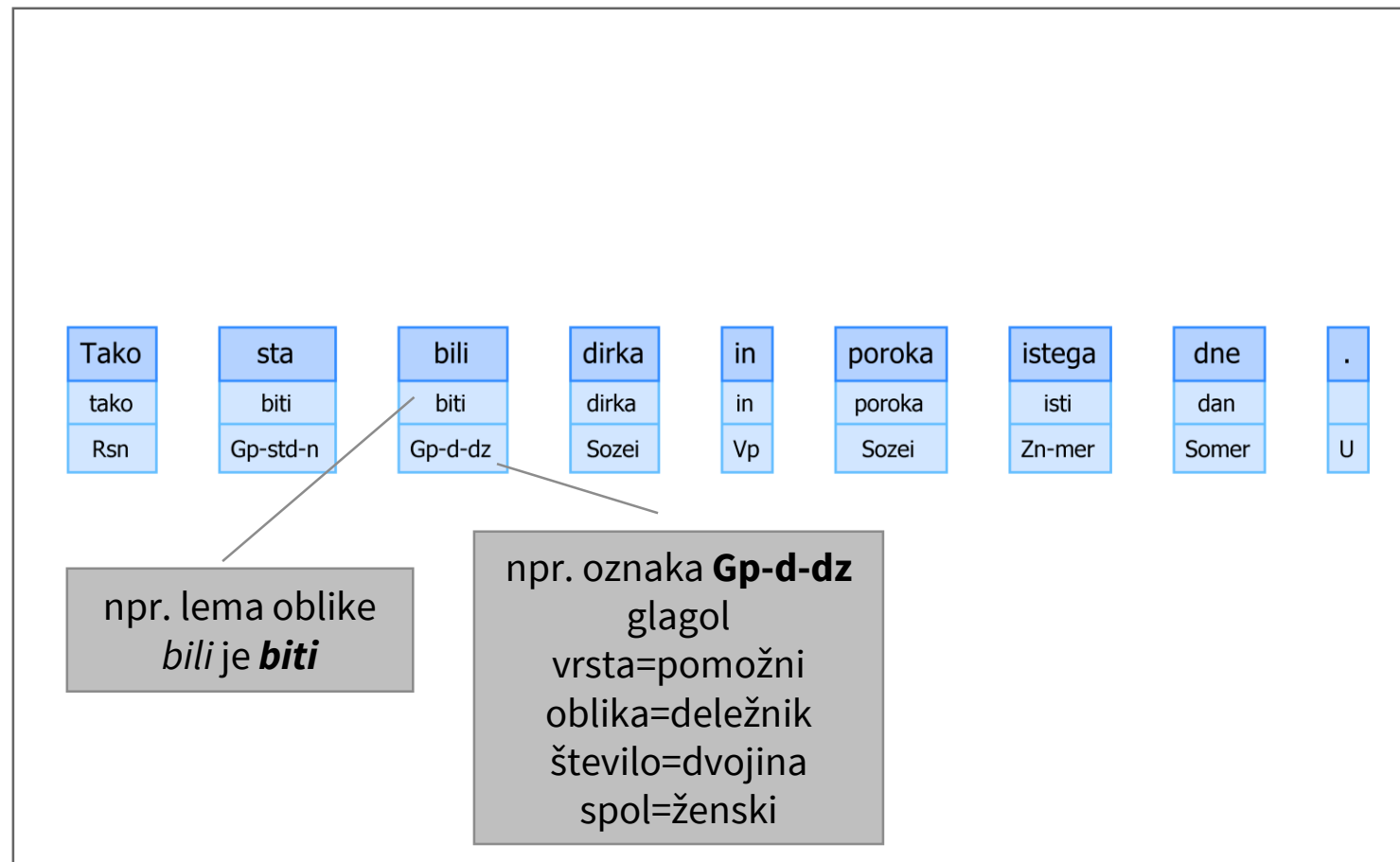
Označenost korpusa ssj500k

- več ravni označenosti različnega obsega





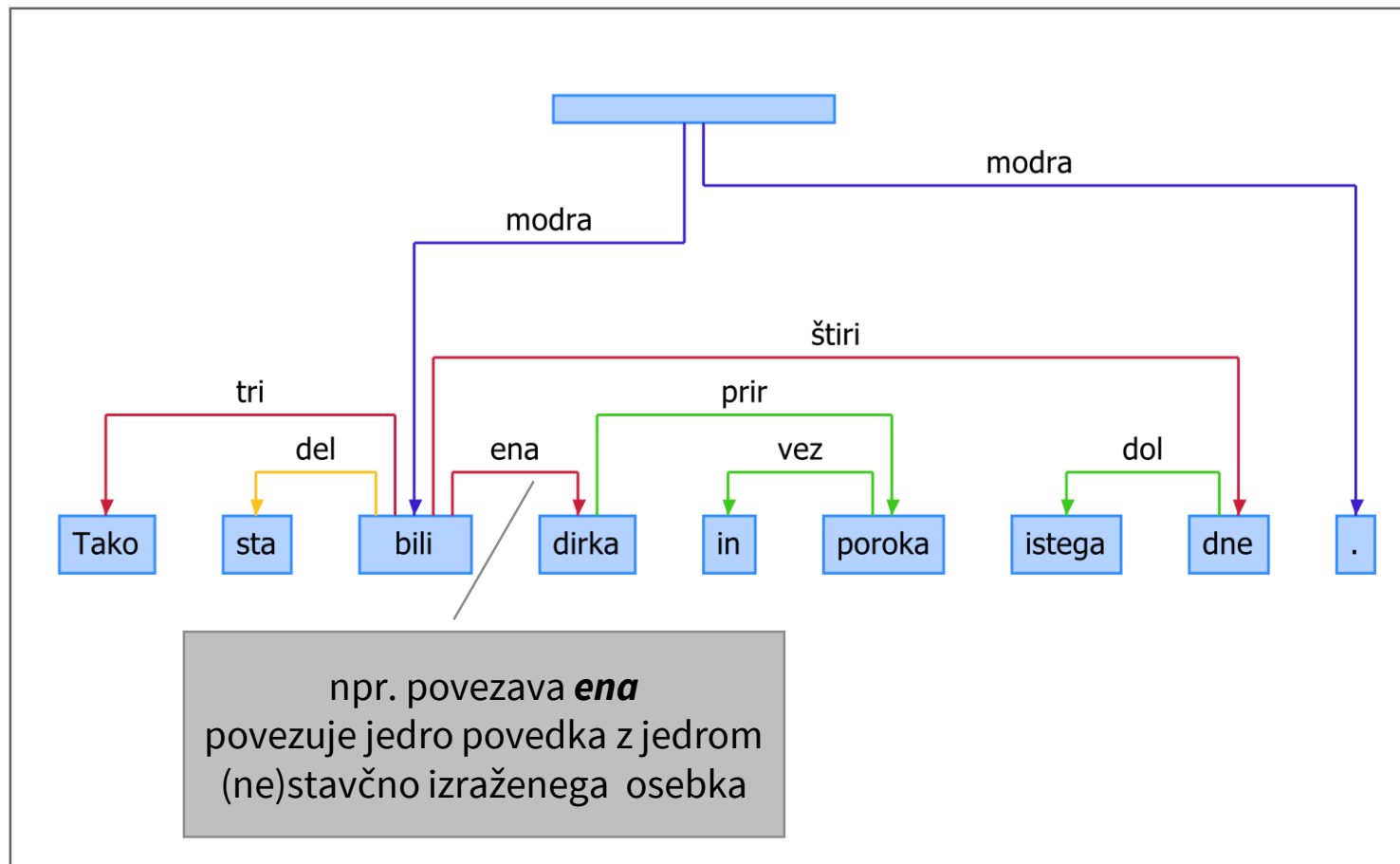
Oblikoskladnja JOS



- pripisovanje **osnovnih oblik** besed (lem) in njihovih **oblikoslovnih lastnosti** (oblikoskladenjskih oznak)
- označevalni sistem **JOS** (Erjavec in Krek 2008; Holozan et al. 2008) – opredelitev nabora besednih vrst, oblikoskladenjskih lastnosti in vrednosti za slovenščino, pravila lematizacije
- <http://nl.ijs.si/jos/>
- označen celoten korpus



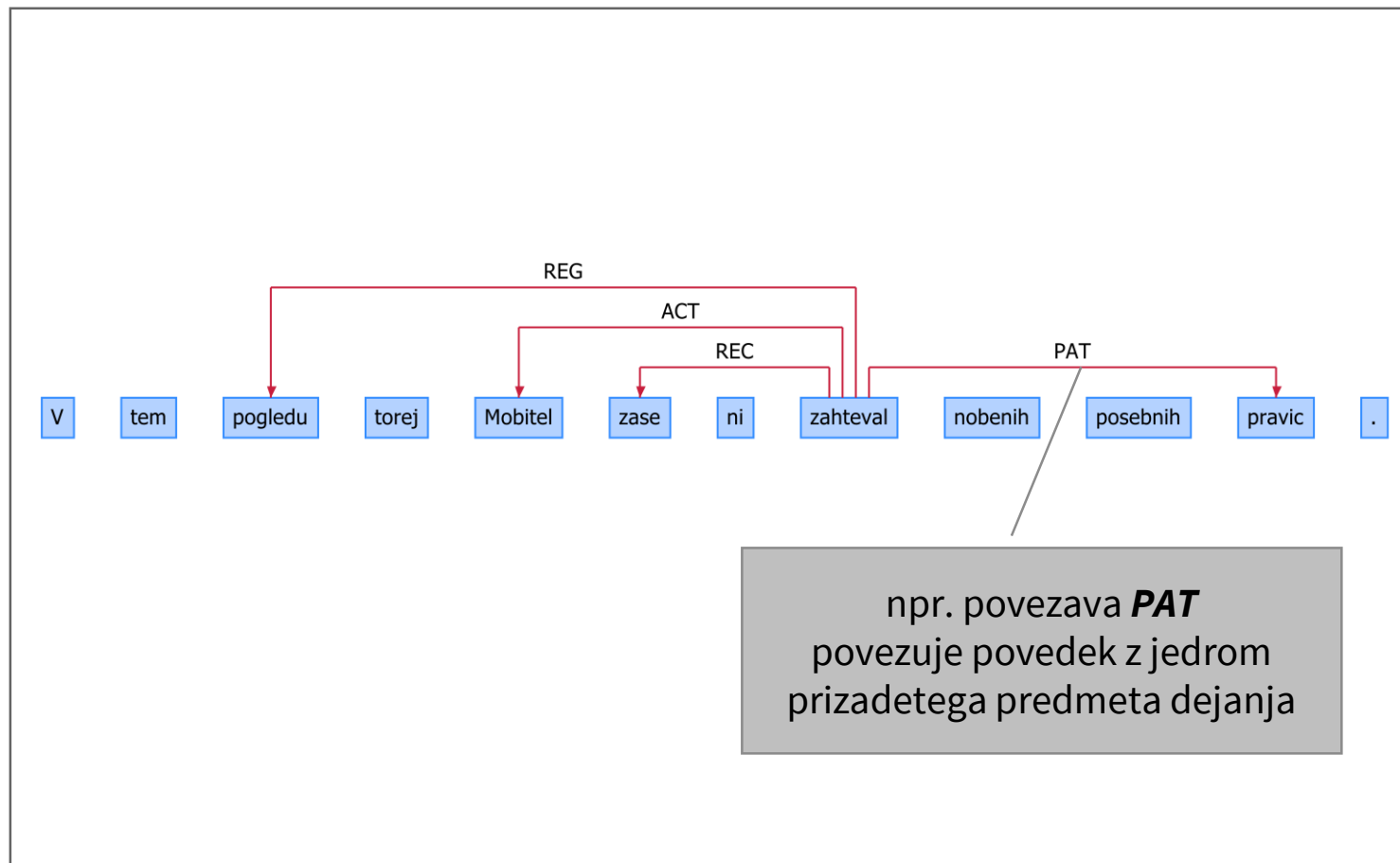
Skladnja JOS



- razčlenjevanje površinske **skladenjske zgradbe povedi** (teorija odvisnostne skladnje)
- označevalni sistem **JOS** (Erjavec et al. 2010; Holozan et al. 2008) – opredelitev **deset jedrnih tipov** skladenjskih povezav
- osredotočenost na **besednozvezno skladnjo** in **glagolsko vezljivost**
- označenih pribl. 40 % korpusa (11.400 povedi)



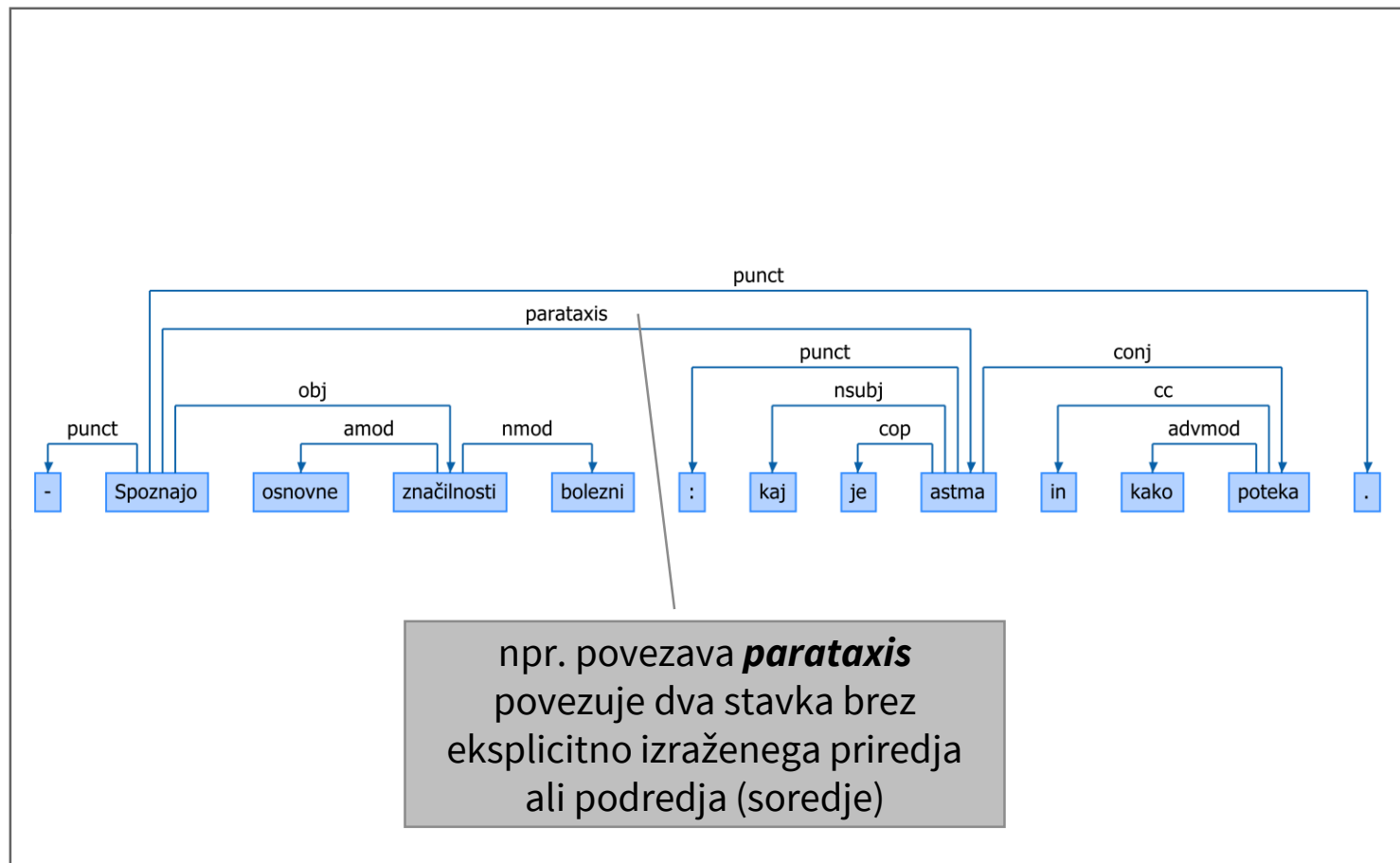
Udeleženske vloge



- razčlenjevanje **pomenske zgradbe povedi** – povedje in njegova določila
- slovensko-hrvaška **označevalna shema SRL** (Gantar et al. 2018) na podlagi češke in drugih sorodnih tipologij
- **25 udeleženskih vlog:** 5 delovalnikov, 17 prislovnih dopolnil, 3 povezave za stalne zveze
- označenih pribl. 19 % korpusa (5.501 povedi)



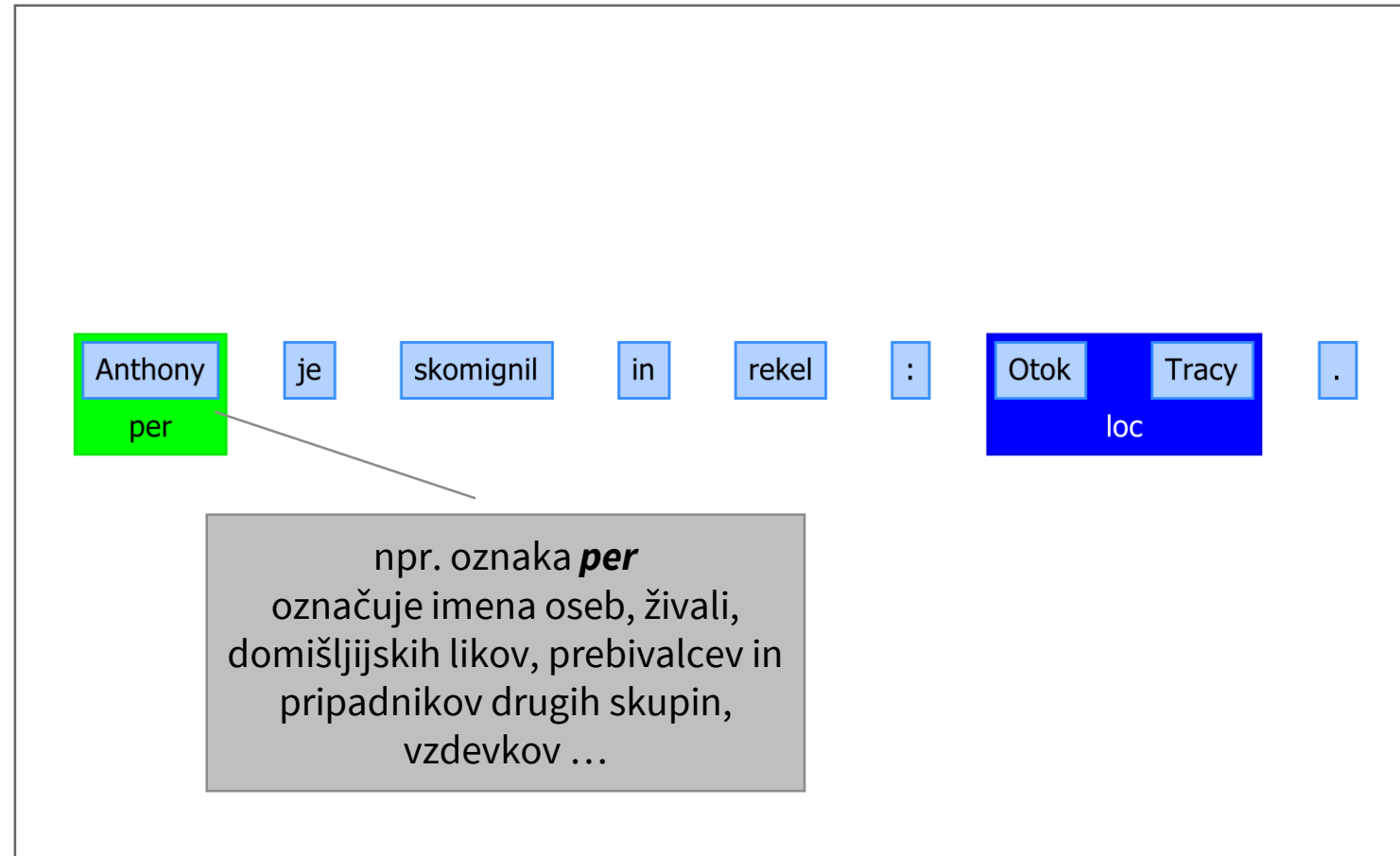
Universal Dependencies



- **mednarodno poenoten** sistem oblikoslovnega in skladskega označevanja korpusov
- 37 tipov skladske povezav, ki pokrivajo **celotno strukturo povedi**
- **strojna pretvorba** na podlagi pravil preslikave med JOS in UD (Dobrovoljc et al. 2017)
- označenih 100 % (oblikoslovje) oz. 24 % (8.000 razčlenjenih povedi) korpusa



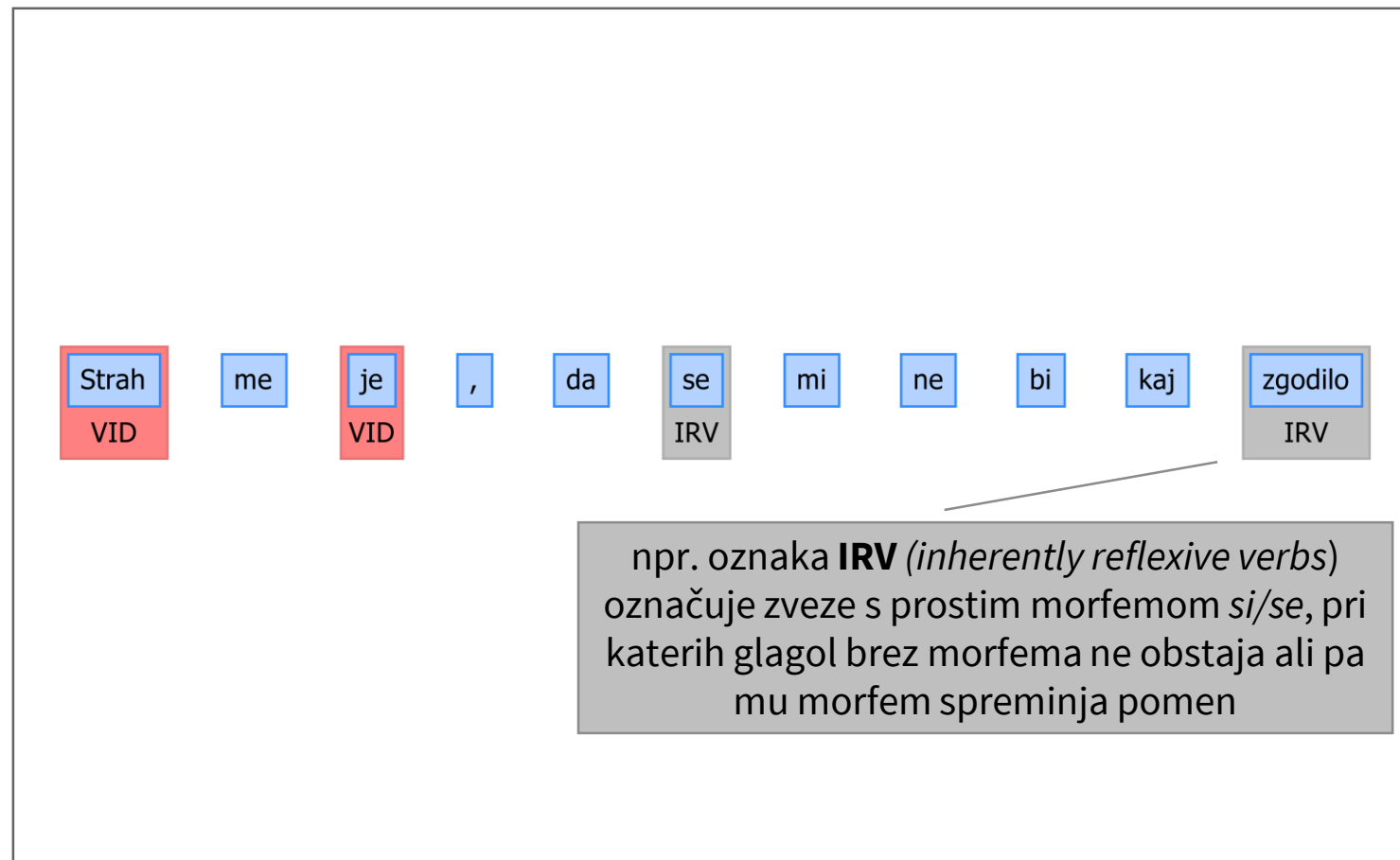
Imenske entitete



- označevanje **lastnoimenskih** besed in besednih zvez
- smernice, razvite znotraj projekta **JANES** (Fišer et al. 2018):
<http://nl.ijs.si/janes/wp-content/uploads/2017/09/SlovenianNER-eng-v1.1.pdf>
- **5 tipov entitet** (osebe in izpeljani svojilni pridevniki, kraji, organizacije, drugo)
- označenih pribl. 33 % korpusa (9.488 povedi, 7.016 entitet)




Glagolske večbesedne enote



- označevanje **stalnih besednih zvez** z glagolsko sestavino
- mednarodno usklajen sistem označevanja, razvit znotraj akcije COST **PARSEME** (Candito et al. 2016; Gantar et al. 2018)
- **4 tipi zvez**: idiomi, zveze s pomensko oslabljenimi glagoli, predložnimi zvezami in povratnimi zaimki
- označenih pribl. 48 % korpusa (13.511 povedi, 3.364 enot)



Dostopnost korpusa ssj500k 2.2

- **podatkovna baza** na repozitoriju CLARIN.SI
 - <http://hdl.handle.net/11356/1210> 
 - različni formati (XML TEI, .vert, CONLL-U)
- konkordančnika **noSketchEngine** in **Kontext**
 - https://www.clarin.si/noske/run.cgi/corp_info?corpname=ssj500k22
 - https://www.clarin.si/kontext/first_form?corpname=ssj500k22
 - brskanje po besednih oblikah in njihovih oznakah
 - ogled konkordanc, izdelava frekvenčnih seznamov
- označevalno orodje **Q-CAT**
 - <http://hdl.handle.net/11356/1262>
 - naprednejše brskanje po kombinacijah oznak
 - spreminjanje ali dodajanje oznak

Orodje **Q-CAT**



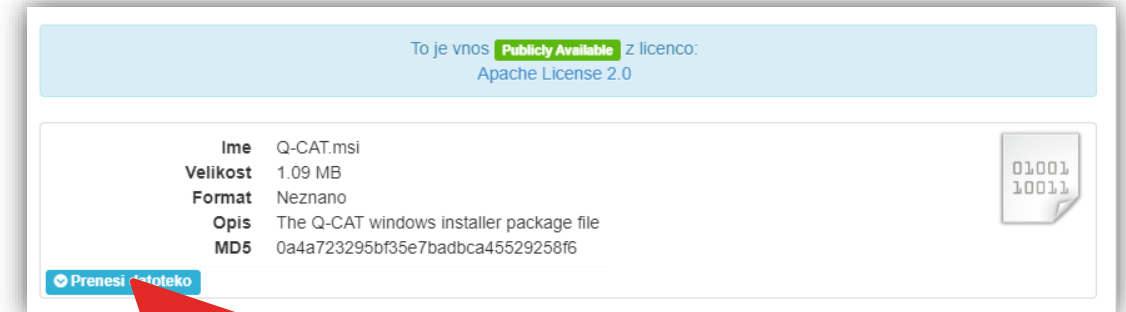
Kaj je orodje Q-CAT

- Q-CAT: (*Querying-Supported*) *Corpus Annotation Tool*
- računalniški **program za izdelavo in analizo** ročno označenih besedilnih korpusov
- avtor: **Janez Brank** (IJS)
- prvotno razvit v okviru projekta [Sporazumevanje v slovenskem jeziku](#) (2008-2013)
 - pod imenom **SentenceMarkup**
 - skladiščno razčlenjevanje in pripisovanje imenskih entitet
- nadgrajen v okviru projekta [Nova slovnica sodobne standardne slovenščine: viri in metode](#) (2017–2020)
 - uvoz poljubnega korpusa s poljubnimi ravnmi označevanja
 - dinamično spreminjanje nastavitev
 - omogočena kompleksnejša iskanja
 - izboljšana uporabniška izkušnja



Dostopnost in namestitvev

- prosto dostopen na repozitoriju CLARIN.SI
 - <http://hdl.handle.net/11356/1262>
- deluje na operacijskem sistemu **Windows**
- preprost za namestitvev
 - prenos in zagon datoteke Q-CAT.msi
- dobro dokumentiran



Priročnik za uporabo:
<https://bit.ly/32ypbhy>



Prikaz uporabe na primeru korpusa ssj500k

- okvirni scenarij
 - 1. Uvoz** korpusa in **brskanje**
 - 2. Iskanje** po označenem korpusu
 - iskanje po posameznih ravneh
 - kombiniranje iskanj po več ravneh
 - pregledovanje in shranjevanje zadetkov
 - 3. Označevanje**
 - popravljanje ali dodajanje že opredeljenih oznak
 - oblikovanje nove ravni označevanja

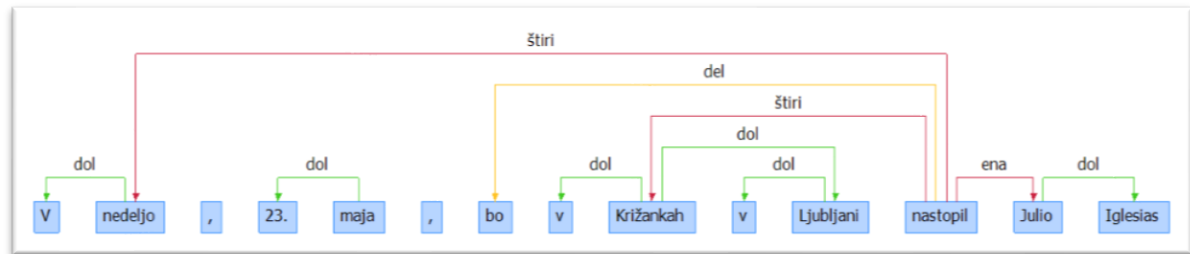
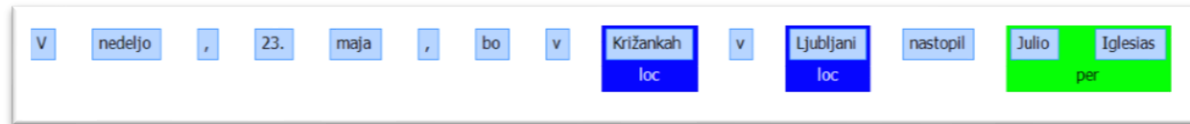
Priročnik za uporabo:
<https://bit.ly/32ypbhy>



Demo

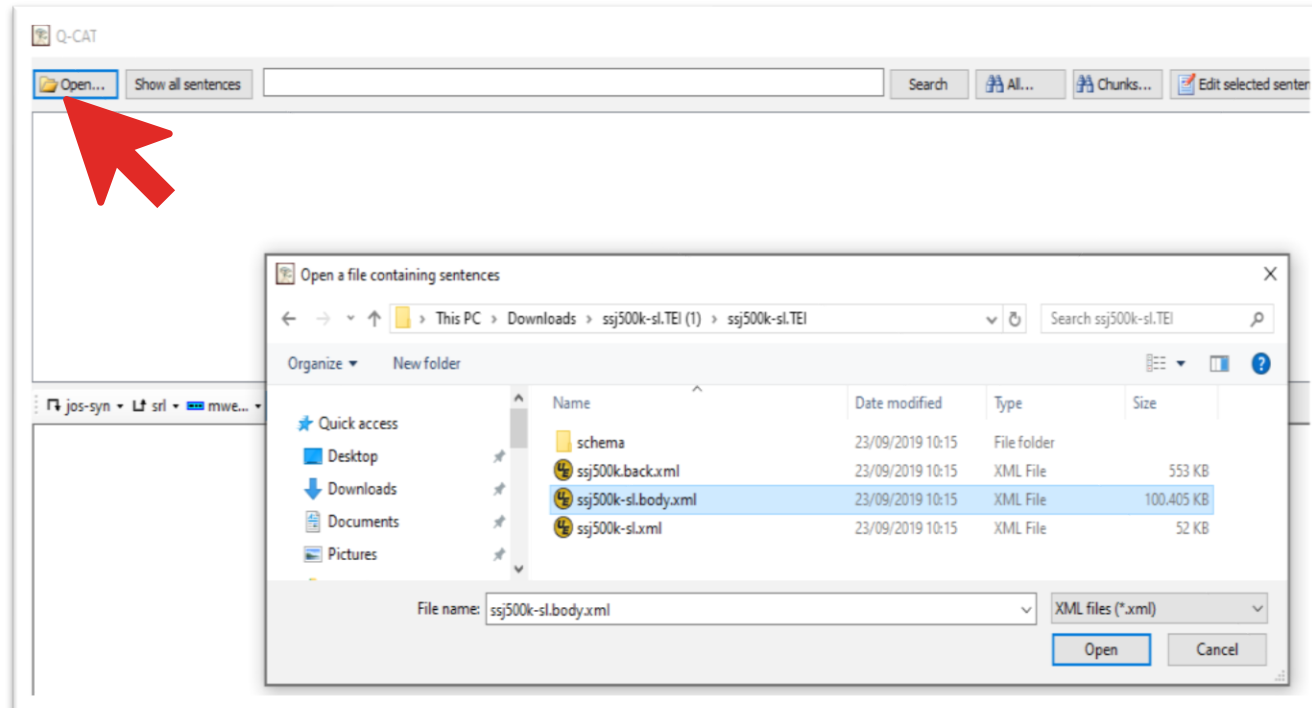
Vrste oznak

- **oznake oblik**
 - npr. oblika, lema, msd
 - modra polja
- **oznake nizov** (angl. *chunks*)
 - npr. stalne besedne zveze
 - barvne ploščice
- **oznake povezav** (angl. *links*)
 - npr. skladienjske povezave
 - barvne puščice



Demo

Uvoz korpusa



- zagon datoteke **QCat.exe**
- izbira korpusa v format XML (obvezno **segmentiran na besede in povedi**)
- brez dodatnih nastavitev lahko uvozimo
 - korpus **brez slovničnih oznak**
 - korpus **z eno ali več ravnmi označenosti ssj500k**
- v primeru uvoza korpusa z nepoznanimi oznakami, potrebne predhodne spremembe nastavitev

Odpiranje korpusa

Prikaz vseh povedi

Enostavno iskanje povedi

Iskanje po označenem korpusu

Označevanje povedi

Urejanje nastavitv

Okno s prikazom vseh povedi

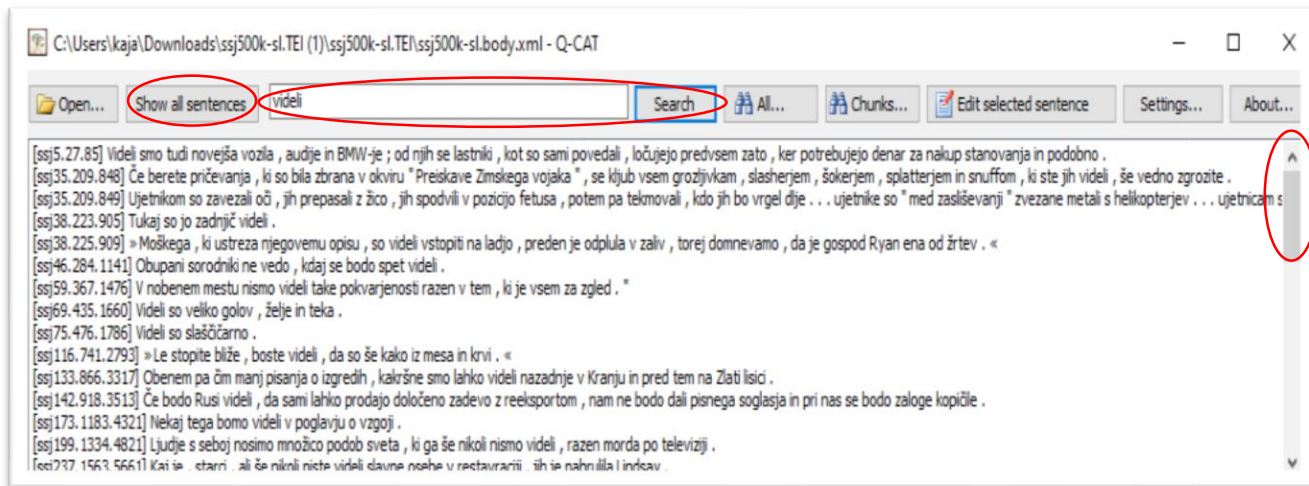
Nastavitve prikaza označene povedi

Okno za prikaz označene povedi

The screenshot shows the CAT software interface. At the top, there's a menu bar with options like 'Open...', 'Show all sentences', 'Search', 'All...', 'Chunks...', 'Edit selected sentence', 'Settings', and 'About...'. Below the menu bar, a list of sentences is displayed, each with a unique ID and a snippet of text. A red callout points to this list, labeled 'Okno s prikazom vseh povedi'. Below the list, there's a toolbar with various icons and a dropdown menu. A red callout points to this area, labeled 'Nastavitve prikaza označene povedi'. The main part of the interface shows a detailed view of a sentence with its constituent words and their grammatical roles. A red callout points to this view, labeled 'Okno za prikaz označene povedi'. The words are arranged in a tree-like structure, with arrows indicating relationships between them. The words are: " Tistega večera sem preveč popil , zgodilo se je mesec dni po tem , ko sem izvedel , da me žena vara .

Demo

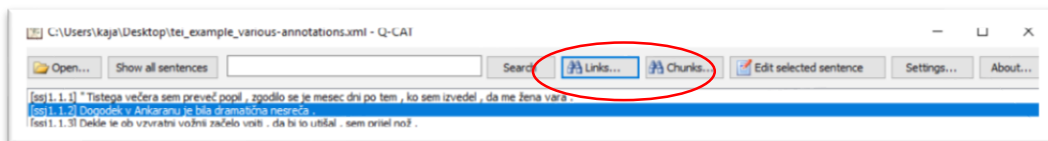
Brskanje po povedih



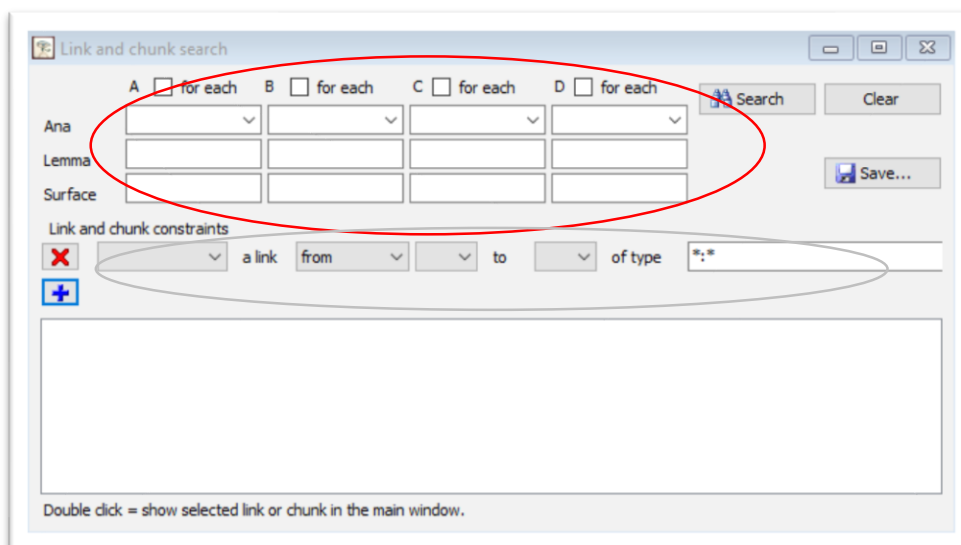
- po seznamu povedu se premikamo s puščicami ali drsnikom na desni
- iskalno okno za enostavno iskanje povedi išče po **besednih oblikah** (npr. *videli*)
- v primeru vnosa več besednih oblik (npr. *so videli*) program poišče povedi, ki jih vsebujejo, ne glede na dejanski vrstni red

Demo

Naprednejša iskanja po korpusu



- orodje omogoča dve vrsti naprednejših iskanj po korpusu: **splošno iskanje** (*All*) in **iskanje zgolj po nizih** (*Chunks*)
- vmesnik za iskanje je sestavljen iz območja za opredelitev **lastnosti besed** (zgoraj) in območja za opredelitev **oznak nizov/povezav** (spodaj)



Demo

Rezultati iskanja

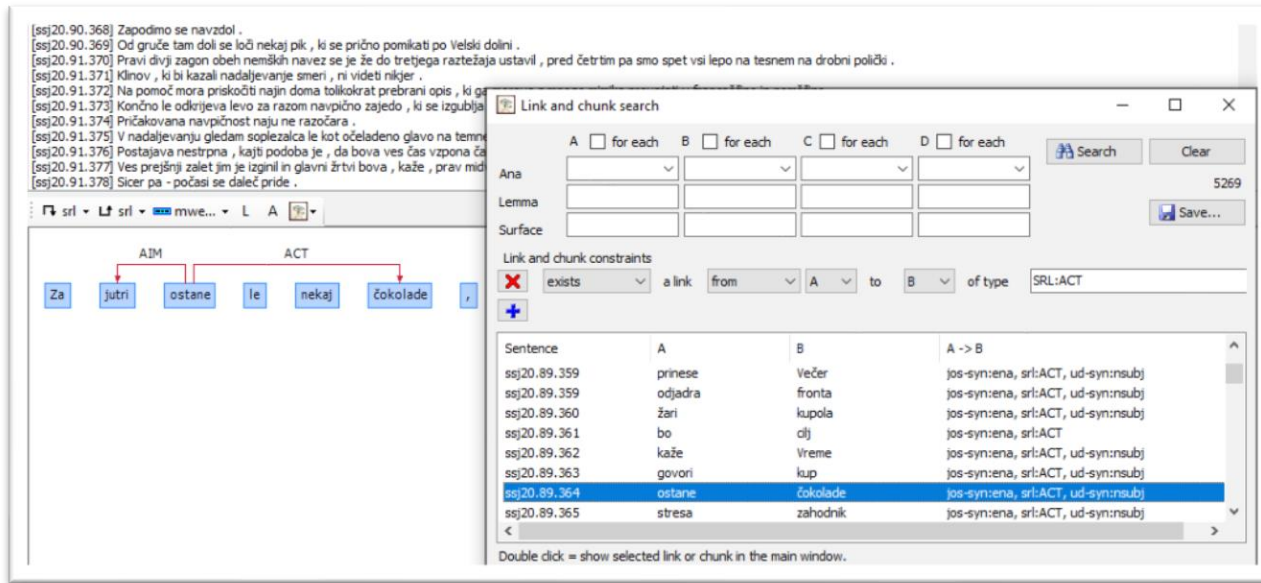
The screenshot shows a search interface with a list of search results on the left and a 'Link and chunk search' dialog box on the right. The dialog box has a search input field containing 'Gg*' and a table of search results. A red arrow points from the 'biti' button in the dialog box to the search results table.

Sentence	A
ssj225.1468.5340	bijejo
ssj643.3303.13760	bije
ssj758.3795.1484	biti
ssj1022.4979.17341	bil
ssj1308.6347.22013	bito
ssj1330.6408.22245	biti
ssj1345.6479.22504	biti
ssj1353.6527.22663	bito
ssj1353.6529.22679	bita
ssj1378.6706.23278	bito

- iskanje vrne **število zadetkov** (desno zgoraj) in seznam vseh **ustreznih povedi** in opredeljenih pojavnih
- z dvojnim klikom na zadetek se v glavnem oknu v ozadju prikaže ustrezna poved

Demo

(Splošno) iskanje po eni vrsti povezave



- S klikom na **modri znak plus** se odpre območje za omejevanje oznak povezav ali nizov, v katerem lahko opredelimo:
 - obstoj (*exists*) ali neobstoj (*doesn't exist*) določene povezave med dvema pojavnicama
 - konkretno (*from*) ali poljubno (*between*) smer povezave med dvema pojavnicama
 - označevalno raven in oznako povezave (polje *:*), npr. SRL:ACT
- iskanje lahko kombiniramo tudi z opredelitvijo dodatnih lastnosti pojavnic (zgoraj), npr. Gg*

Demo

(Splošno) iskanje po več povezavah hkrati

The screenshot shows the 'Link and chunk search' window. At the top, there are search criteria for 'A', 'B', 'C', and 'D' with checkboxes for 'for each'. The search term 'Gg*' is entered in the 'A' field. Below this, there are fields for 'Lemma' and 'Surface'. Under 'Link and chunk constraints', three rules are defined: 'exists' a link from 'A' to 'B' of type 'SRL:ACT', 'exists' a link from 'A' to 'C' of type 'SRL:PAT', and 'exists' a link from 'A' to 'D' of type 'SRL:REC'. A table of results is shown below, with columns for 'Sentence', 'A', 'B', 'C', 'D', and 'A -> B'. The first row is highlighted in blue.

Sentence	A	B	C	D	A -> B
ssj93.619.2350	ustvari	Krepost	spodobnost	nas	jos-syn:ena, srl:
ssj95.621.2369	pojasnjeval	akademik	resnico	nam	jos-syn:ena, srl:
ssj95.625.2392	predstavljajo	Popoldnevi	užitek	družini	jos-syn:ena, srl:
ssj98.629.2436	nimajo	Poraženci	časa	kaj	jos-syn:ena, srl:
ssj98.633.2441	pripravile	sile	akcijo	njega	jos-syn:ena, srl:
ssj100.642.2468	očita	Barberis	domišljije	ji	jos-syn:ena, srl:
ssj100.647.2478	očitajo	Slovenci	požrešnost	državi	jos-syn:ena, srl:
ssj104.670.2549	ukradel	Ta	ogenj	bogovom	jos-syn:ena, srl:
ssj109.686.2595	pomagala	sorodnika	zločinu	mu	jos-syn:ena, srl:
ssj110.688.2598	prinašajo	zakoni	revolucije	področju	jos-syn:ena, srl:

- S klikom na **modri znak plus** odpiramo nova polja za vnos dodatnih iskalnih pogojev
- npr. SRL:ACT + SRL:PAT + SRL:REC, če nas zanimajo povedi s povedki z vzorcem *kdo komu kaj*

Demo

(Splošno) iskanje po različnih ravneh hkrati

The screenshot shows a software interface for multi-level search. The main window displays a sentence: "V bolnišnici bodo uvedli... štirje... tri... šolo... za... starše" with a dependency graph. A "Link and chunk search" dialog is open, showing search criteria and a table of results.

Link and chunk search dialog:

- Search criteria: Ana: Gg*, Lemma: , Surface:
- Link and chunk constraints:
 - exists: a link from A to B of type SRL:ACT
 - doesn't exist: a link from A to B of type JOS-SYN:ena

Table of results:

Sentence	A	B	A -> B
ssj1.1.8	uspelo	morilcu	jos-syn:dve, srl:ACT, ud-syn:obj
ssj2.2.9	uvedli	bolnišnici	jos-syn:štiri, srl:ACT, ud-syn:obl
ssj2.2.10	pripravili	bolnišnici	jos-syn:štiri, srl:ACT, ud-syn:obl
ssj3.5.18	poteka	astma	srl:ACT
ssj3.7.24	gre	zdravljenje	jos-syn:dve, srl:ACT, ud-syn:obl
ssj3.7.25	razkrila	bioenergetičarka	srl:ACT
ssj3.8.27	okrepijo	nihanja	srl:ACT
ssj3.8.27	postavijo	nihanja	srl:ACT
ssj3.8.27	invertrajajo	nihanja	srl:ACT
ssj3.9.32	popije	ženska	srl:ACT
ssj3.9.32	sprosti	ženska	srl:ACT
ssj3.9.32	zaspi	ženska	srl:ACT
ssj3.10.39	prihaja	koželjnica	srl:ACT

- hkratno iskanje lahko izvedemo tudi s povezavami na **več različnih ravneh označevanja**
- npr. iskanje vršilcev (SRL:ACT), ki niso osebkki (JOS-SYN:ena)
- npr. iskanje vršilcev (SRL:ACT) v stalnih besednih zvezah (MWE:*)

Demo

Kompleksnejše iskanje z negacijo

The screenshot shows the Q-CAT software interface. The main window displays a list of sentences with various annotations. A search window titled "Link and chunk search" is open, showing search criteria and constraints. The search criteria include "A" (Gg*), "B" (checked "for each"), "C" (unchecked "for each"), and "D" (unchecked "for each"). The search results are displayed in a table below the constraints.

Sentence	A	C	D	A -> C
ssj3.9.38	začutil	sebi	življenje	jos-syn:štiri, srl:REC, ud-syn:obl
ssj5.20.73	zganila	meni	ga	jos-syn:štiri, srl:REC
ssj5.29.91	posvečati	tehnik	pozornosti	jos-syn:dve, srl:REC, ud-syn:obl
ssj5.29.92	paziti	slednjih	varnost	jos-syn:štiri, srl:REC, ud-syn:obl
ssj5.30.95	priznajo	nam	strategijo	jos-syn:dve, srl:REC, ud-syn:obl
ssj6.34.132	imel	večera	odnos	jos-syn:dve, srl:REC, ud-syn:obl
ssj12.46.210	odpreti	Noki	pokrovček	jos-syn:štiri, srl:REC
ssj12.47.215	poda	nam	vrednosti	jos-syn:dve, srl:REC, ud-syn:obl
ssj12.49.226	prodajali	strankam	glasbo	jos-syn:dve, srl:REC, ud-syn:obl
ssj20.93.395	pokloni	mi	kredenco	jos-syn:dve, srl:REC, ud-syn:obl
ssj21.102.431	ukradli	mi	denarnico	jos-syn:dve, srl:REC, ud-syn:obl
ssj21.102.433	pokazala	mu	ročaja	jos-syn:dve, srl:REC, ud-syn:obl
ssj21.102.434	izmaknil	mi	denarnico	jos-syn:dve, srl:REC

- kadar funkcijo negacije uporabimo za povezave med več različnimi pari pojavnic (npr. A in B ter A in C), moramo opredeliti tudi, ali negacija velja za vse možne kombinacije pojavnic ali ne.
- Temu je namenjeno polje **for each** pri vsaki pojavnici, ki ga izberemo, če mora izbrani pogoj veljati za vsako pojavnico v povedi.

Demo

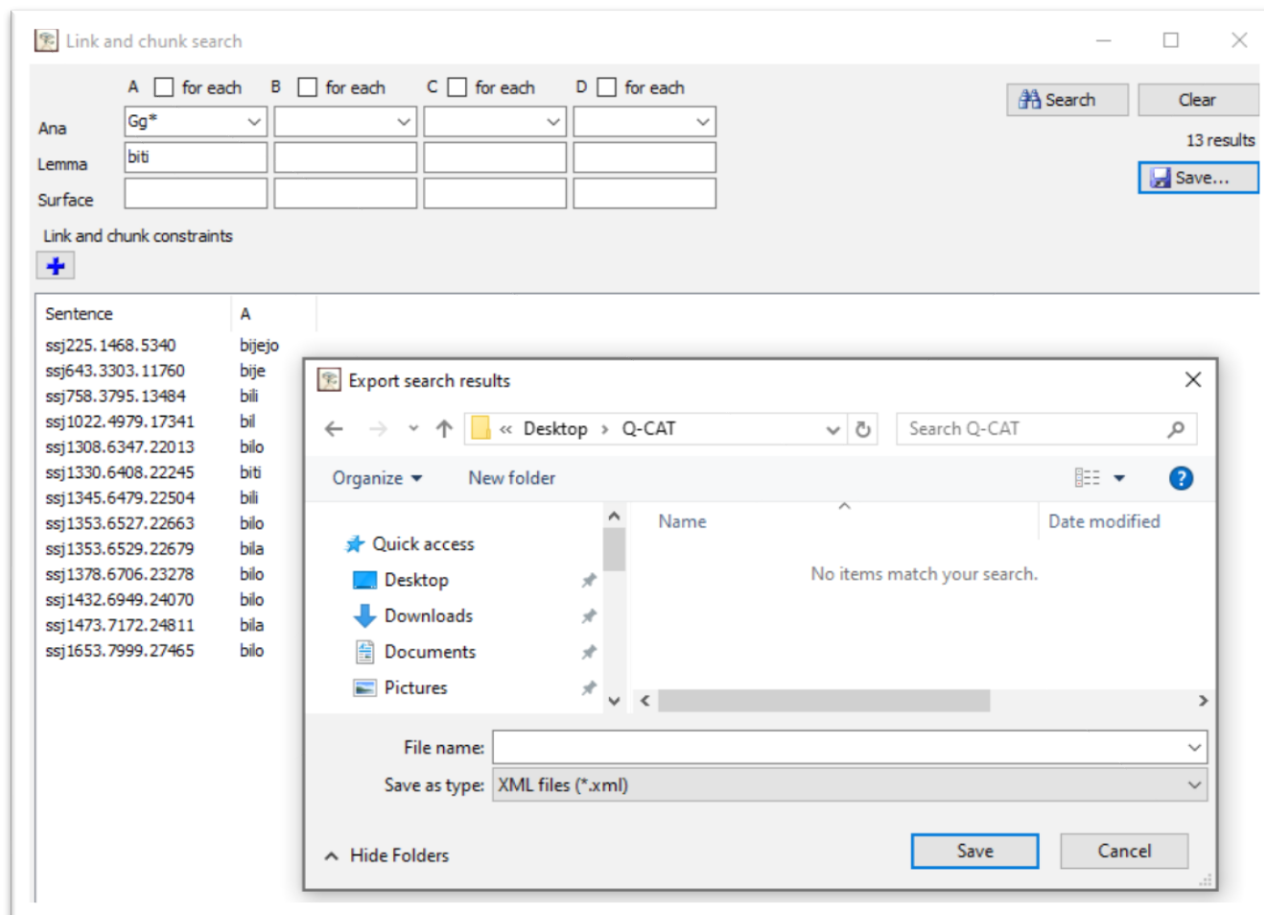
Iskanje po nizih

Sentence	Type	Surface	A
ssj5.21.74	MWE:VID	Na voljo imeli	imeli
ssj38.219.898	MWE:VID	imaš glavo na ramenih	imaš
ssj53.331.1356	MWE:VID	imava rada	imava
ssj69.436.1663	MWE:VID	Glavno besedo imata	imata
ssj81.527.1994	MWE:VID	na voljo imamo	imamo
ssj81.532.2022	MWE:VID	imamo prav	imamo
ssj84.549.2091	MWE:VID	imajo na voljo	imajo
ssj84.556.2117	MWE:VID	nima nič z	nima
ssj91.608.2291	MWE:VID	imeli v delu	imeli
ssj95.625.2398	MWE:VID	nima smisla	nima
ssj98.629.2435	MWE:VID	ima pred nosom	ima
ssj111.698.2632	MWE:VID	imajo na voljo	imajo
ssj114.712.2672	MWE:VID	rad imam	imam
ssj116.725.2715	MWE:VID	ima na skrbi	ima
ssj121.776.2933	MWE:VID	imajo prav	imajo
ssj126.805.3042	MWE:VID	imel prav	imel
ssj131.825.3121	MWE:VID	imel na vesti	imel
ssj133.862.3299	MWE:VID	Za seboj ima	ima
ssj133.874.3350	MWE:VID	nimaš popravnega izpita	nimaš
ssj134.878.3367	MWE:VID	imela zadosti	imela
ssj135.881.3380	MWE:VID	s imeti opravka	imeti
ssj141.907.3477	MWE:VID	imate v pesti	imate
ssj141.912.3492	MWE:VID	imejte glavno besedo	imejte
ssj142.920.3520	MWE:VID	pomena nima	nima
ssj144.984.3725	MWE:VID	ima rada	ima
ssj149.1023.3836	MWE:VID	Radi imamo	imamo
ssj160.1080.4036	MWE:VID	imeti radi	imeti
ssj177.1176.4293	MWE:VID	nimata nima	nimata

- deluje podobno kot splošno iskanje, razen:
 - vsa **iskanja se izvajajo izključno po pojavnicah z oznakami nizov** (npr. imenske entitete, MWE)
 - pri izpisu rezultatov se v stolpcu *Surface* **izpiše celotna zveza**
- npr. iskanje leme *imeti* vrne vse stalne zveze s to lemo
- npr. iskanje leme *imeti* s povezavo MWE:VID vrne vse frazeme s to lemo

Demo

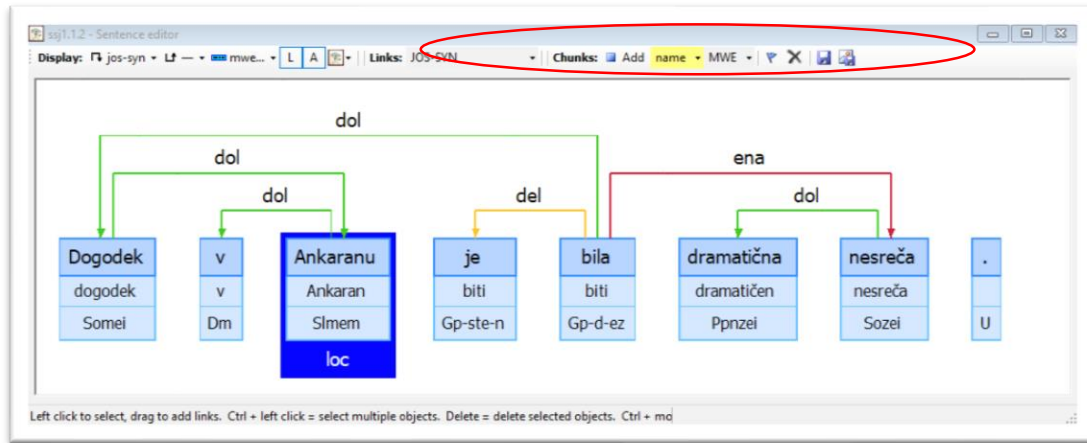
Shranjevanje rezultatov iskanj



- rezultate iskanj lahko shranimo s klikom na gumb *Save*
- Z izbiro shranjevanja v obliki **XML** shranimo vse povedi v enaki obliki, kot je poved zapisana v izhodiščnem korpusu, z vsemi oznakami vred. S tem načinom shranjevanja torej **ustvarimo podkorpus** izhodiščnega korpusa.
- Z izbiro shranjevanja v obliki **.txt** shranimo vse rezultate v obliki podobnega **seznama**, ločenega s tabulatorji, kot se prikaže pri rezultatih iskanja v orodju. Uporabno za nadaljni pregled v Excelu.

Demo

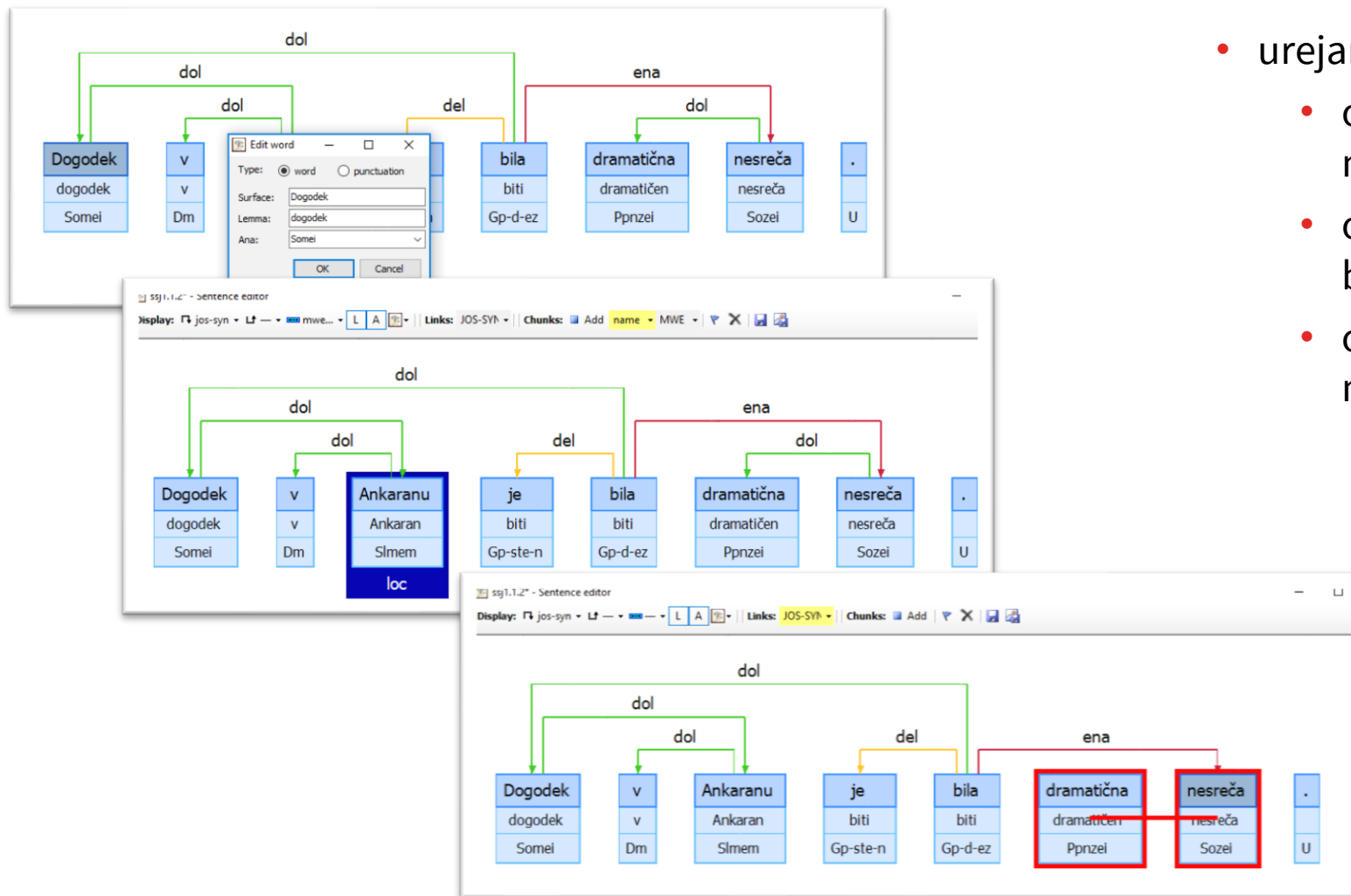
Označevanje korpusa



- s klikom na *Edit selected sentence* v zgornji vrstici vmesnika se odpre okno za urejanje povedi
- v njem **vklopimo vsaj prikaz ravni, ki jo želimo urejati**, lahko pa tudi druge

Demo

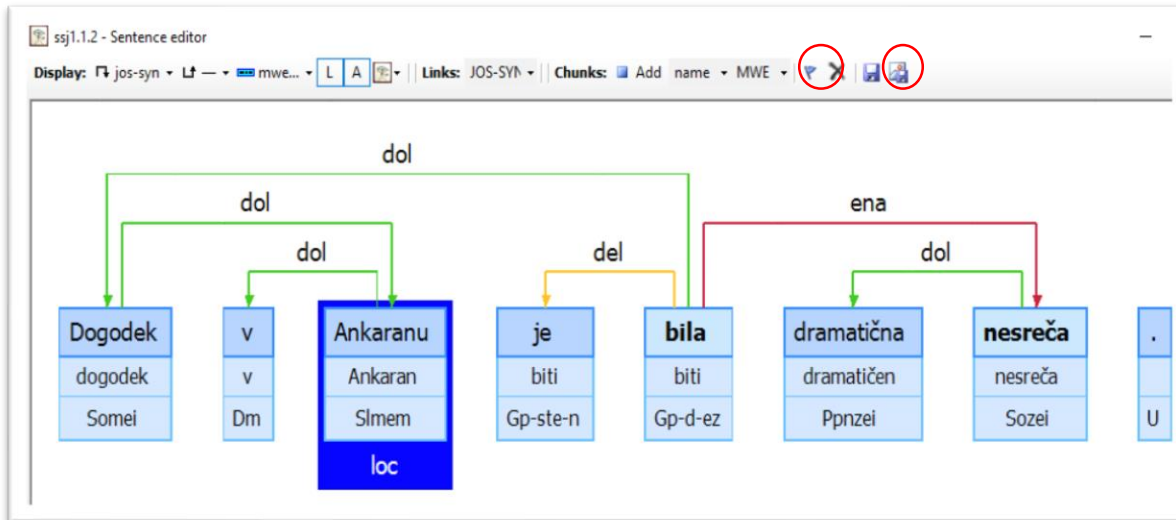
Označevanje korpusa



- urejamo lahko:
 - oblike, leme, msd-je (dvoklik na modra polja pod besedilom)
 - oznake nizov (klik na ploščice pod besedami ali besednimi zvezami)
 - oznake povezav (klik na puščice med besedama)

Demo

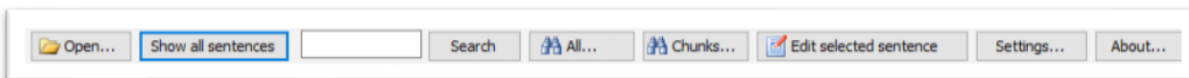
Shranjevanje in izvoz slike



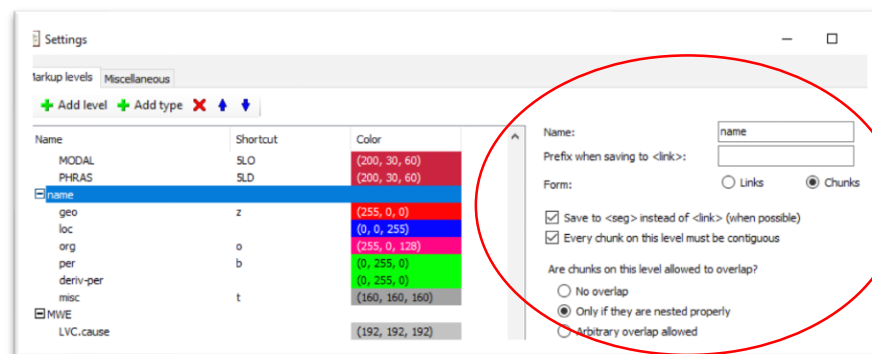
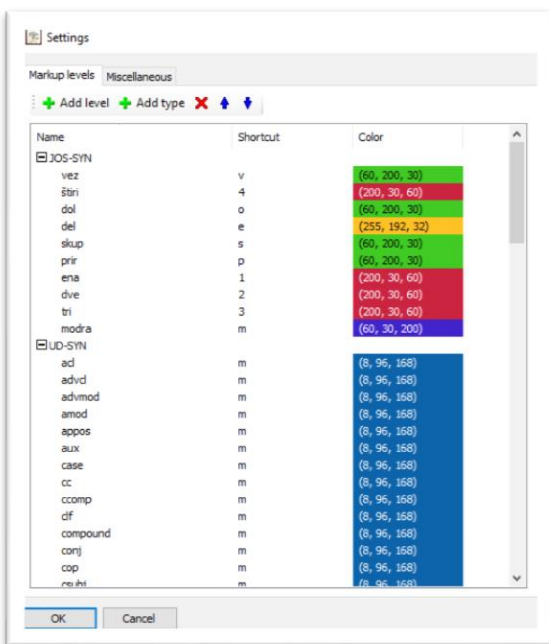
- spremembe shranimo s klikom na ikono za shranjevanje oz. nas na to opozori program ob izhodu iz okna za urejanje. Če ne določimo drugače, se spremembe **shranijo v izvorno datoteko**.
- označeno poved lahko tudi **shranimo kot sliko** visoke ločljivosti

Demo

Nastavitve označevanje



- v vmesniku za nastavitve označevanja (gumb *Settings*) lahko **spremenimo nastavitve obstoječih ravni** označevanja, npr.
 - dodajanje/brisanje/omejevanje oznak
 - urejanje bližnjic
 - urejanje barv
- ali **ustvarimo nove ravni** označevanja



Format

- XML TEI (pomoč na info@clarin.si)
- korpus mora biti vsaj segmentiran na povedi/besede, lahko pa vsebuje tudi:
 - oznake oblik (npr. lema, msd)
 - ena ali več označevalnih ravni iz ssj500k
 - poljubne druge oznake nizov/povezav – spremembe nastavitev pred začetkom dela
- primeri v mapi Programi/IJS/Q-CAT/Samples

Zaključek

- predstavili **korpus ssj500k**, največjo prostodostopno zbirko slovenskih besedil z ročno pripisanimi slovničnimi oznakami
- poleg uporabnosti za razvoj jezikovnih tehnologij ima korpus velik potencial za **kvalitativne korpusne analize** slovničnih pojavov v slovenščini
- ta analiza bistveno olajšana z **orodjem Q-CAT**, ki je prenosljivo tudi na druge korpuse
- poleg obstoječih ravni označenosti ssj500k **v pripravi** še:
 - stalne besedne zveze po smernicah LBS (Gantar 2015; Gantar et al. 2019)
 - oznake koreferenčnih povezav (Žitnik in Bajec 2018)
 - odstranjevanje nestandardnih besedil (projekt NSSS)
 - ...
- ssj500k kot preizkusni kamen za druge jezikoslovne teorije in tipologije v našem prostoru?

Bibliografija

- Krek, S. et al. (2019). **The ssj500k Training Corpus for Slovenian Language Processing**. *V pripravi*.
- Dobrovoljc, K. (2019). **Q-CAT: Orodje za ročno označevanje in analizo besedilnih korpusov. Priročnik za uporabo**. <https://bit.ly/32ypbhy>
- Arhar, S. and Gorjanc, V. (2007). Korpus FidaPLUS: Nova generacija slovenskega referencnega korpusa (The FidaPLUS Corpus: A New Generation of the Slovene Reference Corpus). *Jezik in slovstvo*, 52(2):95–110.
- Dobrovoljc, K., Erjavec, T., and Krek, S. (2017). The Universal Dependencies Treebank for Slovenian. In *Proceedings of the 6th Workshop on Balto-Slavic Natural Language Processing, BSNLP@EACL 2017*, pages 33–38.
- Erjavec, T. and Krek, S. (2008). The JOS morphosyntactically tagged corpus of Slovene. In *LREC 2008*.
- Erjavec, T., Fišer, D., Krek, S., and Ledinek, N. (2010). The JOS Linguistically Tagged Corpus of Slovene. In *Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10)*, Valletta, Malta, May. European Language Resources Association (ELRA).
- Erjavec, T. (2004). MULTTEXT-east version 3: Multilingual morphosyntactic specifications, lexicons and corpora. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal, May. European Language Resources Association (ELRA).
- Fišer, D., Ljubešič, N., and Erjavec, T. (2018). The Janes project: language resources and tools for Slovene user generated content. *Language Resources and Evaluation*, Sep.
- Gantar, P., Arhar Holdt, Š., Čibej, J., Kuzman, T., and Kavčič, T. (2018a). Glagolske večbesedne enote v učnem korpusu ssj500k 2.1. In *Proceedings of the conference on Language Technologies Digital Humanities*, Ljubljana, Slovenia. Ljubljana University Press, Faculty of Arts.
- Gantar, P., Despot, K. S., Krek, S., and Ljubešič, N. (2018b). Towards semantic role labeling in slovene and croatian. In *Proceedings of the conference on Language Technologies Digital Humanities*, Ljubljana, Slovenia. Ljubljana University Press, Faculty of Arts.
- Holozan, P., Krek, S., Pivec, M., Rigač, S., Rozman, S., and Velušček, A. (2008). Specifikacije za učni korpus. Projekt "Sporazumevanje v slovenskem jeziku". Technical report.
- Krek, S., Dobrovoljc, K., Erjavec, T., Može, S., Ledinek, N., Holz, N., Zupan, K., Gantar, P., Kuzman, T., Čibej, J., Arhar Holdt, , Kavčič, T., Škrjanec, I., Marko, D., Jezeršek, L., and Zajc, A. (2019). *Training corpus ssj500k 2.2*. Slovenian language resource repository CLARIN.SI. <http://hdl.handle.net/11356/1210>.

Hvala.

info@cjvt.si

Center za
jezikovne vire
in tehnologije

Večna pot 113
1000 Ljubljana
Slovenija

www.cjvt.si
00386 14798299
info@cjvt.si