

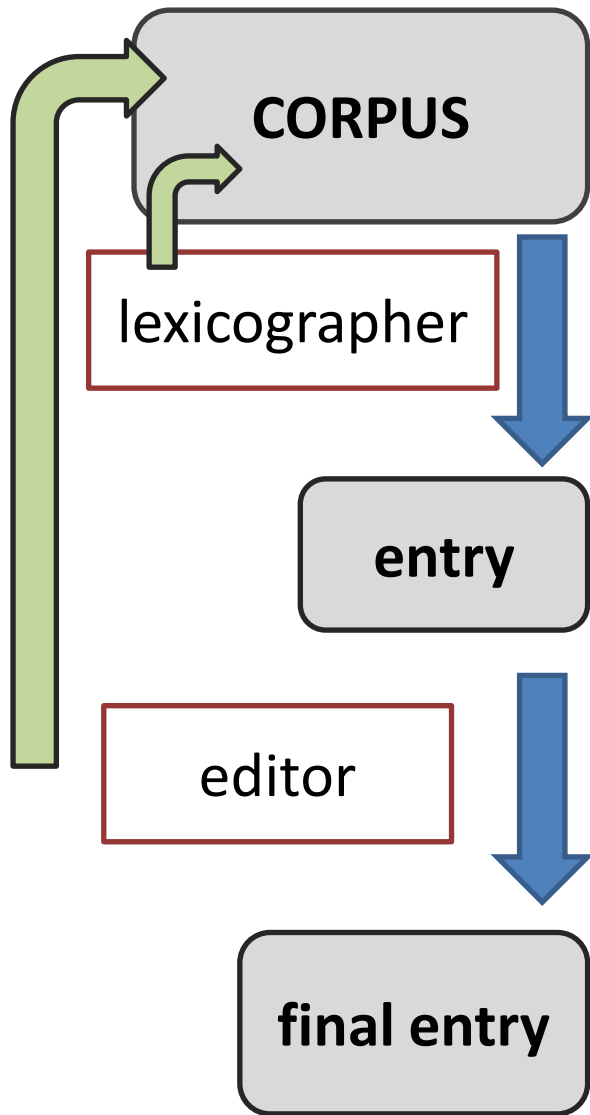
Monolingual Lexicography

**Automatic extraction of data from corpora
for lexicographic purposes**

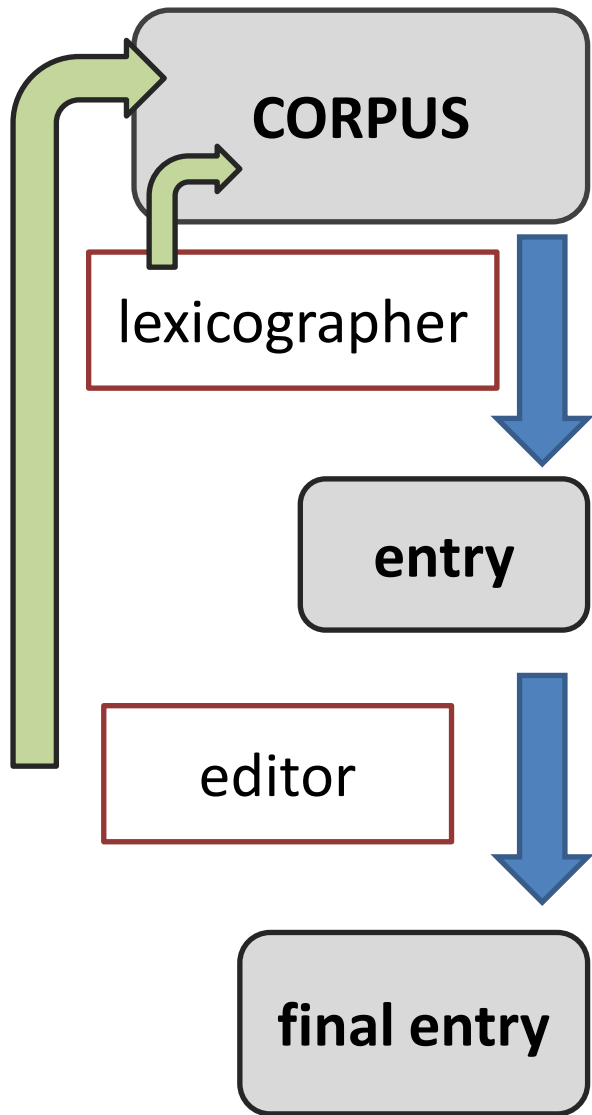
**Proposal for the Compilation of the
Dictionary of Modern Slovene**

Session 3

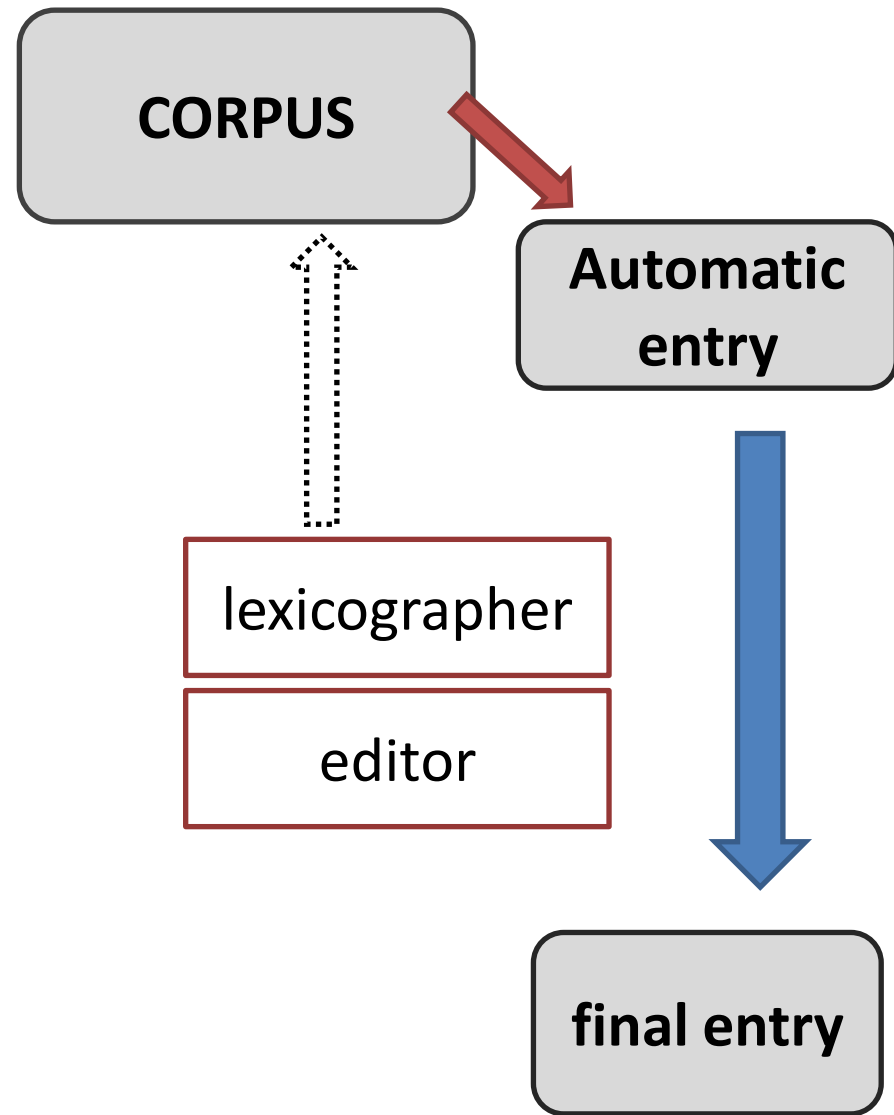
NOW



NOW



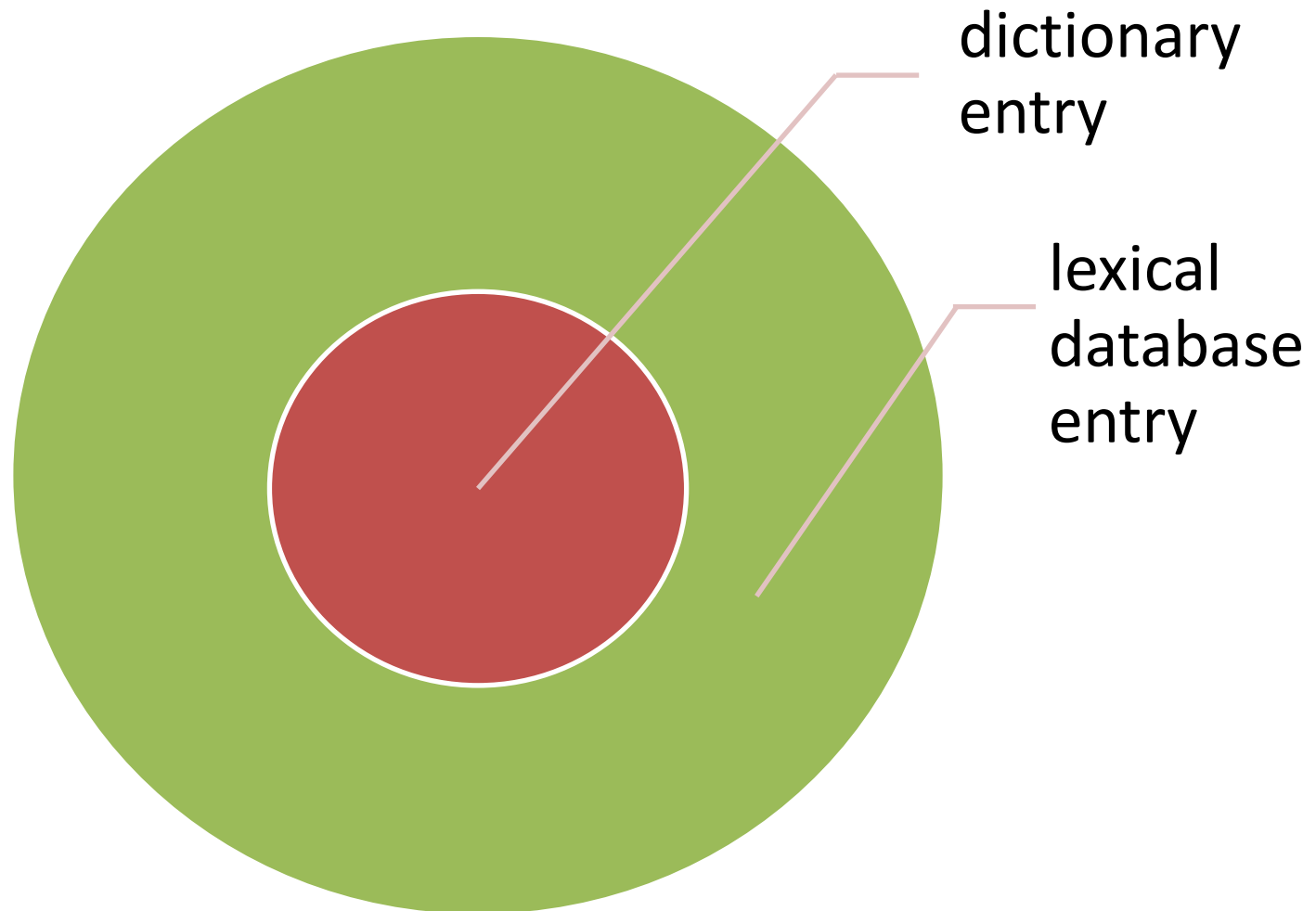
Rundell & Kilgarriff (2011)



Slovenian situation

- The only (scholarly) monolingual dictionary of Slovene published in 1991 (with many entries older, from 70s and 80s)
 - Based on material from late 19th century to 1970s
 - Problems with IPR (many authors / owners of copyright)
- New edition announced for October 2014
 - will build on the old version + the supplement from 2012
 - the same IPR problem as the first edition

Slovene Lexical Database



Towards automation

- Sketch grammar
 - Version 2 (Krek, 2010): for manual analysis
 - Version 3 (2012): for automation
 - finely-grained sketch grammar based on the structures identified in the lexical database
- GDEX for Slovene
 - Version 1 (Kosem et al., 2011): for manual analysis
 - Version 2 (2012): for automation
 - different configurations for different word classes

I. LEMMA

- headword
- part-of-speech

svitati se (to dawn)

verb

II. SENSE

- indicator

1. *daniti se (day)*

2. *dojemati (understand)*

- semantic

unary relations & constructions

a DAN.
najati sonce

če se ČLOVEKU začne svitati o nekem DOGAJANJU. začne dojemati. kar prej ni vedel. ali pa je bilo to pred njim skrito

gramrels

III. SYNTAX

- lable

only in 3rd pers.

- structure
- pattern

gbz Inf-GBZ

kaj se svita (sth is dawning)

rbz GBZ

komu se svita o čem (sth is dawning to sb about sth)

word sketches

- synt. combin.

- collocation

[začeti. pričeti] se svitati

[počasi. malo. malce] se svita

GDEX

V. EXAMPLES

- example

Preden se začne zjutraj svitati. je najtemnejša noč.

Počasi se mi je začelo svitati. zakaj Jasni oči tako žarijo.

Na vzhodu se je že svital dan. ko sta se poslovila.

Petru se pričinja svitati o nekdanji zvezi ned Chadom in Heather.

- multi-word unit

VI. PHRASEOLOGY • phraseological units

Sketch grammar

- regular expressions over POS tags
=a_modifier/modifies
2:[tag="P.*"] 1:[tag="S.*"]
- the name of the arguments (order)
- 1: 2: = words to be extracted as the first/second argument
- |, ., (), {} and * - standard metacharacters (RE)

Regular gramrels

Number	SR-14
Type	enodelna
Name	=nedoločnik
English (appr.)	=infinitive
Sketch Grammar	Query
1:[]	Verb, Adjective, Adverb, Noun
SLDB syntactic structure	nedoločnik
GBZ Inf-gbz	uspeti [doseči]
SBZ1 Inf-gbz	(imeti) možnost [pritožiti se]
PBZ Inf-gbz	pripravljen [oditi]; sposoben [izpeljati]
RBZ Inf-gbz	dobro [izkoristiti]

DUAL gramrels

Number	SR-01
Type	recipročna (*DUAL)
Name	=kakšen?/kdo-kaj?
English (appr.)	=what_kind?/who-what?
Sketch Grammar	Query
1: [tag="S.*"]	Noun
2: [tag="P.*"]	Adjective
SLDB structure	kakšen? / what_kind?
pbz SBZ1	[boleč, lep] spomin
SLDB structure	kdo-kaj? / who-what?
PBZ sbz1	rdeča [žoga]

TRINARY gramrels

Number	SR-06
Type	tridelna (*TRINARY)
Name	=%s
English (appr.)	=%s
Sketch Grammar	Query
1:[tag="G.*"]	Verb
SLDB syntactic structure	=%s
sbz1 GBZ sbz4 (s sbz6 / med sbz6)	stisniti s [prsti, palci]
sbz1 GBZ sbz4 (s sbz6 / med sbz6)	stisniti med [prsti]
sbz1 GBZ (o sbz5 glede sbz2 za sbz4)	pogajati se o [vdaji, izpustitvi]

Automation – Sketch grammar

- use of macros – easier to read
- direct relation between SLD elements and gramrels included in the grammar
- *SEPARATEPAGE (very complex)
- *CONSTRUCTION (very useful)
- *COLLOC (for „syntactic combinations“ in SLD)

Macros examples

- `define(`nedolocnik', `[tag="G.n.*"]')`
- `define(`pomoznik', `[tag="Gv.*"]')`
- `define(`deleznik', `[tag="Gpd.*"]')`
- `define(`gl_nebiti', `[tag="G.*" & lemma!="biti"]')`
- `define(`gl_sed_3', `[tag="Gpp.t.*"]')`
- `define(`brez_GSVD', `[tag!="[GSVD].*" & word!="[,;()-]"]')`

Macros used in gremrels

- =predl-pred
 - 2:predlog 1:samostalnik
- =%s_s6
 - 1:samostalnik 3:predlog brez_GSVD{0,5}
 - 2:samost_oro
- =S_V_O3_O2
 - 2:osebek brez_PSVD{0,5} 1:glagol brez_SVD{0,5}
 - predmet_daj{1,4} brez_SVD{0,5} predmet_rod

Example: NOUN_1-noun_2

VERB + prep + NOUN-gen

„dobiti iz česa“ / to get from sth

- `<struktura>GBZ %s sbz2</struktura>`

- *SEPARATEPAGE `koga-česa_g2`

- *TRINARY

`=%s_g2`

1:glagol sise{0,2} 3:predlog brez_GSVDK{0,5}

2:samost_rod

3:predlog brez_GSVDK{0,5} 2:samost_rod sise{0,1}

1:glagol

koga-česa_g2	485	
od-d_g2	206	17.9
iz-d_g2	107	5.0
do-d_g2	22	1.6
brez-d_g2	21	3.7
v-d_g2	21	12.0
poleg-d_g2	21	8.9
zaradi-d_g2	18	2.0
z-d_g2	15	2.1
za-d_g2	14	4.4

NOUN_1-noun_2

dobiti

(glagol)

FidaPlus (20M) freq = 11162 (735.5 per million)

displaying only: koga-česa_g2

whole word sketch

brez-d_g2	21	3.7	iz-d_g2	107	5.0	na-d_g2	5	3.1	za-d_g2	14	4.4	od-d_g2	206	17.9
laganje	1	9.61	Nigerija	2	8.82	priložnost	1	3.18	spletka	1	8.64	dalajlama	2	8.05
recept	6	8.34	onostranstvo	1	8.21	stran	3	3.12	pust	1	8.33	gradbenik	2	7.93
odredba	1	6.69	arest	1	8.14	sredstvo	1	2.41	karta	1	5.09	prednik	2	7.43
natečaj	1	6.06	katran	1	8.13				mleko	1	4.25	Avstrijec	2	7.34
težava	9	4.87	Ligojna	1	8.13	o-d_g2	1	8.0	denar	4	3.8	Pelhan	1	7.3
razpis	1	3.95	Juršinci	1	8.07	informacija	1	2.91	naloga	1	3.56	Adanič	1	7.29
problem	1	3.02	sukcesija	1	8.03	zaradi-d_g2	18	2.0	stanovanje	1	3.29	Izabela	1	7.21
razlog	1	2.89	limfa	1	8.01	kolegialnost	1	10.68	vloga	1	2.64	deklič	1	7.17
			vrečica	1	7.92	črevesje	1	8.06	pravica	1	2.12	Lek	2	7.15
			ZPIZ	1	7.91	panika	1	7.32	cesta	1	1.92	FIBA	1	7.13
			proračun	20	7.74	taktika	1	7.17				NEK	1	7.12
			NT	1	7.71	obnašanje	2	7.15				Uefa	1	7.11
			Montreal	1	7.71	noša	1	7.08				bršljan	1	7.09
			Kremelj	1	7.63	prostornina	1	6.63				Portugalska	1	6.98
			mozeg	1	7.62	pnevmatika	1	6.32				vol	1	6.93
			Carigrad	1	7.6							Nik	1	6.93

*CONSTRUCTION (very useful)

- Element <vzorci> = syntactic patterns
 - who/what does sb sth
 - who/what does sth to sb etc.
- In entries with verbs as headwords
- Under structures + collocations
- Now: examples with binary collocations
- CONSTRUCTION: examples with complete patterns

Example: S_V_O3_O4

=S_V_O3_O4

"subject"

"indirect
object"

"direct
object"

2:osebek brez_PSVD{0,5} 1:glagol brez_SVD{0,5}
predmet_daj{1,4} brez_SVD{0,5} predmet_toz

2:osebek brez_PSVD{0,5} 1:glagol brez_SVD{0,5}
predmet_toz{1,4} brez_SVD{0,5} predmet_daj

2:osebek brez_PSVD{0,5} predmet_daj{1,4}
brez_SVD{0,5} 1:glagol brez_SVD{0,5} predmet_toz

2:osebek brez_PSVD{0,5} predmet_toz{1,4}
brez_SVD{0,5} 1:glagol brez_SVD{0,5} predmet_daj

Example from SkE

<u>z_nikalnim</u> 155 7.3	<u>s_prislovom</u> 112 3.3	<u>kakšen-p?</u> 6 1.2	<u>S_V_O3_O4</u> 47 18.0	<u>S_V_O3_O2</u> 18 5.2
mir 36 9.37	dušek 2 8.62	Hast 1 12.19	Mikulin 1 9.39	Kacin 1 9.64
soglasje 15 8.6	močnik 1 8.13	Jehan 1 12.19	Henigman 1 9.3	gostilničar 1 8.64
veto 2 8.03	spodbuda 4 7.79	Votan 1 12.0	Požun 1 9.19	Istrabenz 1 8.15
gol 9 7.97	tornado 1 7.76	jurjev 1 11.19	razvijalec 1 8.61	dekan 1 7.93
predujem 1 7.5	individualnost 1 7.76	kratek 1 3.36	Tuđman 1 8.54	pod 1 6.57
maksimum 1 7.36	zalet 1 7.57	dober 1 0.63	Kristus 1 8.14	namestnik 1 6.3
frustracija 1 7.32	samozavest 2 7.45	<u>kakšnega-p</u> 1 0.9	Jarc 1 7.94	plod 1 5.89
pokoj 1 7.06	šarm 1 7.33	preprost 1 4.82	iskrica 1 7.9	mora 1 5.83
golaž 1 6.94	drobiž 1 7.31		Pahor 1 7.81	Rusija 1 5.72
kis 1 6.9	kad 1 7.17		Harry 1 7.53	uedba 1 5.45
povod 1 6.83	plastenka 1 7.14		mineral 1 7.45	profesor 1 4.92
malica 1 6.79	optimizem 2 7.14		pečat 1 7.15	Krka 1 4.5
mladič 1 6.68	poudarek 3 7.12		pobudnik 1 7.1	oblast 1 3.14
breza 2 6.64	morala 1 7.09		Bill 1 7.02	večina 1 2.57
rezultat 13 6.58	pianist 1 7.05		prireditelj 1 6.96	območje 1 2.48
odgovor 8 6.47	kovanec 1 6.93		gospoda 1 6.92	človek 2 2.12
plaketa 1 6.32	avtonomija 1 6.89		govornik 1 6.72	država 1 0.84
dar 1 6.29	modrost 1 6.43		Washington 1 6.53	
modrost 1 6.28	odpoved 1 6.3		Sevnica 1 6.27	
prid 1 6.03	prid 1 6.16		testiranje 1 6.23	
odmerek 1 5.98	pooblastilo 1 6.1		Cerkev 1 6.04	
pojasnilo 1 5.87	skrb 3 6.06		volilec 1 5.86	
dovoljenje 5 5.76	žito 1 5.93		potrošnik 1 5.82	
opis 1 5.66	gol 2 5.86		Martin 1 5.8	
moka 1 5.65	inflacija 1 5.83		Hrvaška 1 5.01	

Examples – high precision

poročil z njo. Tako združene **moči** so tovarni **dale** nov zagon in postala je najboljša tovarna predsednik Borut **Pahor** je Drnovšku ponovno **dal** košarico ne sicer za večno, ampak vsaj posebej zadovoljen, ker so neuvrščene **države dale** vso prednost jedrski razorožitvi. Udeleženci Učite se od mojstrov. " Najboljši **govorniki dajo** svojim poslušalcem vselej občutek, kot Vrhniki. Gostilna **Iskrica** je pohodnikom **dala** lonec pasulja, Marko Breclj je uredil na terenu, še preden mednarodna **skupnost da** ZN mandat za ukrepanje, " je izjavil Anan Jelovec, del Sredme). **Sveti Martin je dal** svoje ime cerkvi s prepoznavnim baročnim privoščič še mineralno kopel; **minerali dajo** vodi posebno zeleno barvo. V spremstvu premier Akajeva, je tako kot vrhovno **sodišče dal** prednost staremu parlamentu, ki je Bakijeva 1997 je mestni svetnik Mihael **Jarc** sicer **dal** pobudo mestnemu svetu, da bi po Janezu podaljšati - Državne **ustanove** so jagrom **dale** čepice - Zadovoljstvo v Luksemburgu - Tudi sta) **Kongresno testiranje** varnosti **dalo** porazno sliko **Kot švicarski sir** leto pa napoveduje, da bodo **prireditelji dali** večji poudarek tudi pohodom, ki jih bodo štela za plačano z dnem, ko bo **potrošnik dal** nalog taki organizaciji. Ali pa, če bo V nadaljevanju je trener Miro **Požun** spet **dal** priložnost mladim igralcem in prav vsi, je upravičena domneva, da je **Washington dal** Manili tiho podporo za poskus vojaške rešitve

*COLLOC for „syntactic combinations“

- Element <zveza> = syntactic combinations
 - "v razmerju do" (in relation to)
- Mainly nominal headwords
- Under (sub)sense after syntactic structures as a separate category

COLLOC: d_sam_d

- =d_sam_d
- *COLLOC "%(2.lemma)_%(3.lemma)-p"
- 2:predlog 1:samostalnik 3:predlog

preposition

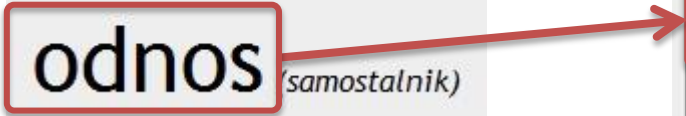
noun

preposition

Example SkE: "in relation to"

corpus: FidaPlus (20M)

odnos (*samostalnik*)



<u>d_sam_d</u>	<u>412</u>	<u>7.0</u>
v_do	<u>92</u>	11.69
za_z	<u>109</u>	10.86
v_med	<u>38</u>	10.4
o_med	<u>18</u>	10.07
o_z	<u>18</u>	9.59
na_med	<u>13</u>	9.14
o_do	<u>8</u>	9.12
z_do	<u>8</u>	8.63
za_med	<u>5</u>	7.96
glede_do	<u>3</u>	7.86
v_z	<u>59</u>	7.84

avtoriteta (samostalnik) Fida PLUS 620m (SLD sketch grammar) frekvenca = 9301 (12.6 na milijon)

gl-pred	5465	0.8
<input type="checkbox"/> spodkopavati	<u>21</u>	27.93
<input type="checkbox"/> spodkopati	<u>29</u>	27.79
<input type="checkbox"/> rušiti	<u>44</u>	27.39
<input type="checkbox"/> spodbijati	<u>26</u>	27.1
<input type="checkbox"/> priznavati	<u>57</u>	26.73
<input type="checkbox"/> zastaviti	<u>61</u>	22.89
<input type="checkbox"/> priznati	<u>100</u>	19.77
<input type="checkbox"/> podrejati	<u>15</u>	19.41
<input type="checkbox"/> izžarevati	<u>17</u>	19.05
<input type="checkbox"/> utrjevati	<u>20</u>	18.87
<input type="checkbox"/> upirati	<u>31</u>	18.81
<input type="checkbox"/> sklicevati	<u>30</u>	18.57
<input type="checkbox"/> uveljaviti	<u>46</u>	17.71
<input type="checkbox"/> izgubljeni	<u>34</u>	17.68
<input type="checkbox"/> spoštovati	<u>39</u>	17.07
<input type="checkbox"/> vzpostaviti	<u>37</u>	16.73
<input type="checkbox"/> omajati	<u>11</u>	15.69
<input type="checkbox"/> imeti	<u>575</u>	15.18
<input type="checkbox"/> uživati	<u>46</u>	14.94
<input type="checkbox"/> uveljavljati	<u>26</u>	14.89
<input type="checkbox"/> utrditi	<u>20</u>	14.57
<input type="checkbox"/> smešiti	<u>6</u>	14.22
<input type="checkbox"/> izrabljati	<u>10</u>	14.06
<input type="checkbox"/> okrepiti	<u>28</u>	13.96
<input type="checkbox"/> zavračati	<u>20</u>	13.87

[>>](#)

kakšen?	4098	3.3
<input type="checkbox"/> moralen	<u>337</u>	67.97
<input type="checkbox"/> nesporen	<u>182</u>	66.8
<input type="checkbox"/> brezpriziven	<u>16</u>	38.76
<input type="checkbox"/> cerkven	<u>70</u>	38.37
<input type="checkbox"/> očetovski	<u>26</u>	38.17
<input type="checkbox"/> starševski	<u>33</u>	38.09
<input type="checkbox"/> učiteljev	<u>23</u>	36.58
<input type="checkbox"/> vrhoven	<u>57</u>	36.21
<input type="checkbox"/> strokoven	<u>126</u>	35.38
<input type="checkbox"/> absoluten	<u>41</u>	35.19
<input type="checkbox"/> visok	<u>184</u>	35.18
<input type="checkbox"/> velik	<u>348</u>	33.94
<input type="checkbox"/> političen	<u>137</u>	33.75
<input type="checkbox"/> nezmotljiv	<u>15</u>	33.4
<input type="checkbox"/> papežev	<u>23</u>	30.93
<input type="checkbox"/> verski	<u>47</u>	30.82
<input type="checkbox"/> nedotakljiv	<u>14</u>	28.19
<input type="checkbox"/> vodilen	<u>35</u>	25.43
<input type="checkbox"/> svetoven	<u>85</u>	25.43
<input type="checkbox"/> lažen	<u>18</u>	25.15
<input type="checkbox"/> nedvomen	<u>10</u>	25.11
<input type="checkbox"/> očetov	<u>18</u>	24.87
<input type="checkbox"/> znanstven	<u>28</u>	24.51
<input type="checkbox"/> duhoven	<u>22</u>	24.21
<input type="checkbox"/> priznan	<u>20</u>	23.91

[>>](#)

gl-za	3478	0.5
<input type="checkbox"/> omajati	<u>9</u>	15.38
<input type="checkbox"/> priznavati	<u>15</u>	14.72
<input type="checkbox"/> spoštovati	<u>23</u>	14.36
<input type="checkbox"/> razvaljati	<u>6</u>	13.1
<input type="checkbox"/> vzpostavljati	<u>9</u>	12.91
<input type="checkbox"/> načeti	<u>15</u>	12.8
<input type="checkbox"/> spodkopati	<u>6</u>	12.74
<input type="checkbox"/> podrejati	<u>6</u>	12.32
<input type="checkbox"/> temeljiti	<u>21</u>	12.25
<input type="checkbox"/> obvladovati	<u>12</u>	11.68
<input type="checkbox"/> močiti	<u>6</u>	11.01
<input type="checkbox"/> vplivati	<u>29</u>	10.99
<input type="checkbox"/> uživati	<u>21</u>	10.51
<input type="checkbox"/> zagovarjati	<u>12</u>	9.96
<input type="checkbox"/> izhajati	<u>16</u>	9.88
<input type="checkbox"/> izkazovati	<u>6</u>	9.28
<input type="checkbox"/> zagotavljati	<u>19</u>	9.11
<input type="checkbox"/> odločati	<u>20</u>	8.95
<input type="checkbox"/> nadomestiti	<u>13</u>	8.73
<input type="checkbox"/> izvirati	<u>9</u>	8.61
<input type="checkbox"/> prepovedovati	<u>6</u>	8.51
<input type="checkbox"/> uveljavljati	<u>9</u>	8.38
<input type="checkbox"/> sklicevati	<u>7</u>	8.23
<input type="checkbox"/> tikati2	<u>8</u>	8.13
<input type="checkbox"/> graditi	<u>14</u>	8.09

[>>](#)

koga-kaj	1759	6.5
<input type="checkbox"/> spodkopati	<u>26</u>	31.8
<input type="checkbox"/> zastaviti	<u>62</u>	29.85
<input type="checkbox"/> rušiti	<u>32</u>	29.27
<input type="checkbox"/> spodkopavati	<u>16</u>	29.12
<input type="checkbox"/> spodbijati	<u>19</u>	28.17
<input type="checkbox"/> izgubljeni	<u>33</u>	23.15
<input type="checkbox"/> izžarevati	<u>15</u>	22.31
<input type="checkbox"/> priznavati	<u>22</u>	21.46
<input checked="" type="checkbox"/> vzpostaviti	<u>31</u>	20.72
<input type="checkbox"/> uveljaviti	<u>31</u>	19.64
<input type="checkbox"/> zavračati	<u>21</u>	19.36
<input type="checkbox"/> utrjevati	<u>13</u>	19.03
<input type="checkbox"/> prezirati	<u>9</u>	18.58
<input type="checkbox"/> utrditi	<u>16</u>	17.28
<input type="checkbox"/> okrepiti	<u>22</u>	17.03
<input type="checkbox"/> zbijati	<u>6</u>	16.36
<input type="checkbox"/> pridobiti	<u>40</u>	16.27
<input type="checkbox"/> uživati	<u>28</u>	16.16
<input type="checkbox"/> spoštovati	<u>20</u>	16.14
<input type="checkbox"/> izrabljati	<u>8</u>	15.77
<input type="checkbox"/> izgubiti	<u>42</u>	15.56
<input type="checkbox"/> omajati	<u>7</u>	15.18
<input type="checkbox"/> razvaljati	<u>6</u>	15.01
<input type="checkbox"/> uveljavljati	<u>15</u>	14.86
<input type="checkbox"/> priznati	<u>30</u>	14.36

[>>](#)

priredje	1236	2.4
<input type="checkbox"/> ugled	<u>85</u>	50.25
<input type="checkbox"/> verodostojnost	<u>23</u>	37.59
<input type="checkbox"/> moč	<u>61</u>	34.21
<input type="checkbox"/> karizma	<u>14</u>	33.61
<input type="checkbox"/> avtoriteta	<u>18</u>	30.03
<input type="checkbox"/> integriteta	<u>12</u>	29.9
<input type="checkbox"/> kredibilnost	<u>11</u>	28.86
<input type="checkbox"/> moč	<u>21</u>	27.57
<input type="checkbox"/> hierarhija	<u>11</u>	25.44
<input type="checkbox"/> legitimnost	<u>9</u>	25.3
<input type="checkbox"/> oblast	<u>25</u>	22.33
<input type="checkbox"/> znanje	<u>18</u>	19.51
<input type="checkbox"/> odgovornost	<u>13</u>	18.85
<input type="checkbox"/> zaupanje	<u>10</u>	18.47
<input type="checkbox"/> samozavest	<u>7</u>	17.71
<input type="checkbox"/> spoštovanje	<u>8</u>	16.85
<input type="checkbox"/> pooblastilo	<u>7</u>	16.68
<input type="checkbox"/> družben	<u>11</u>	15.91
<input type="checkbox"/> vpliv	<u>11</u>	15.22
<input type="checkbox"/> izkušnja	<u>11</u>	14.04
<input type="checkbox"/> položaj	<u>12</u>	13.7
<input type="checkbox"/> vloga	<u>12</u>	13.4
<input type="checkbox"/> tradicija	<u>6</u>	12.14
<input type="checkbox"/> red	<u>7</u>	10.23
<input type="checkbox"/> podoben	<u>6</u>	8.56

[>>](#)

gl-pred	5465	0.8
<input type="checkbox"/> spodkopavati	21	27.93
<input type="checkbox"/> spodkopati	29	27.79
<input type="checkbox"/> rušiti	44	27.39
<input type="checkbox"/> spodbijati	26	27.31
<input type="checkbox"/> priznavati	57	26.73
<input type="checkbox"/> zastaviti	61	22.89
<input type="checkbox"/> priznati	100	19.77
<input type="checkbox"/> podrejati	15	19.41
<input type="checkbox"/> izžarevati	17	19.05
<input type="checkbox"/> utrjevati	20	18.87
<input type="checkbox"/> upirati	31	18.81
<input type="checkbox"/> sklicevati	30	18.57
<input type="checkbox"/> uveljaviti	46	17.71
<input type="checkbox"/> izgubljati	34	17.68
<input type="checkbox"/> spoštovati	39	17.07
<input type="checkbox"/> vzpostaviti	37	16.73
<input type="checkbox"/> omajati	11	15.69
<input type="checkbox"/> imeti	575	15.18
<input type="checkbox"/> uživati	46	14.94
<input type="checkbox"/> uveljavljati	26	14.89
<input type="checkbox"/> utrditi	20	14.57
<input type="checkbox"/> smešiti	6	14.22
<input type="checkbox"/> izrabljati	10	14.06
<input type="checkbox"/> okreptiti	28	13.96
<input type="checkbox"/> zavračati	20	13.87

kakšen?	4098	3.3
<input type="checkbox"/> moralen	337	67.97
<input type="checkbox"/> nesopren	182	66.8
<input type="checkbox"/> brezpriziven	16	38.76
<input type="checkbox"/> cerkven	70	38.37
<input type="checkbox"/> očetovski	26	38.17
<input type="checkbox"/> starševski	33	38.09
<input type="checkbox"/> učiteljev	23	36.58
<input type="checkbox"/> podrejeni	57	36.21
<input type="checkbox"/> strokoven	126	35.38
<input type="checkbox"/> absoluten	41	35.19
<input type="checkbox"/> visok	184	35.18
<input type="checkbox"/> velik	348	33.94
<input type="checkbox"/> političen	137	33.75
<input type="checkbox"/> nezmotljiv	15	33.4
<input type="checkbox"/> papežev	23	30.93
<input type="checkbox"/> verski	47	30.82
<input type="checkbox"/> nedotakljiv	14	28.19
<input type="checkbox"/> vodilen	35	25.43
<input type="checkbox"/> svetoven	85	25.43
<input type="checkbox"/> lažen	18	25.15
<input type="checkbox"/> nedvomen	10	25.11
<input type="checkbox"/> očetov	18	24.87
<input type="checkbox"/> znanstven	28	24.51
<input type="checkbox"/> duhoven	22	24.21
<input type="checkbox"/> priznan	20	23.91

gl-za	3478	0.5
<input type="checkbox"/> omajati	9	15.38
<input type="checkbox"/> priznavati	15	14.72
<input type="checkbox"/> spoštovati	23	14.36
<input type="checkbox"/> razvaljati	6	13.1
<input type="checkbox"/> vzpostavljati	9	12.91
<input type="checkbox"/> načeti	15	12.8
<input type="checkbox"/> spodkopati	6	12.74
<input type="checkbox"/> podrejati	6	12.32
<input type="checkbox"/> temeljiti	21	12.25
<input type="checkbox"/> obvladovati	12	11.68
<input type="checkbox"/> močiti	6	11.01
<input type="checkbox"/> vplivati	29	10.99
<input type="checkbox"/> uživati	21	10.51
<input type="checkbox"/> zagovarjati	12	9.96
<input type="checkbox"/> izhajati	16	9.88
<input type="checkbox"/> izkazovati	6	9.28
<input type="checkbox"/> zagotavljati	19	9.1
<input type="checkbox"/> odločati	20	8.95
<input type="checkbox"/> nadomestiti	13	8.73
<input type="checkbox"/> izvirati	9	8.61
<input type="checkbox"/> prepovedovati	6	8.51
<input type="checkbox"/> uveljavljati	9	8.38
<input type="checkbox"/> sklicevati	7	8.23
<input type="checkbox"/> tikati2	8	8.13
<input type="checkbox"/> graditi	14	8.09

koga-kaj	1759	6.5
<input type="checkbox"/> spodkopati	26	31.8
<input type="checkbox"/> zastaviti	62	29.85
<input type="checkbox"/> rušiti	32	29.27
<input type="checkbox"/> spodkopavati	16	29.12
<input type="checkbox"/> spodbijati	19	28.17
<input type="checkbox"/> izgubljati	33	23.15
<input type="checkbox"/> izžarevati	15	22.31
<input type="checkbox"/> priznavati	22	21.46
<input type="checkbox"/> vzpostaviti	31	20.72
<input type="checkbox"/> uveljaviti	31	19.64
<input type="checkbox"/> zavračati	21	19.36
<input type="checkbox"/> utrjevati	13	19.03
<input type="checkbox"/> prezirati	9	18.58
<input type="checkbox"/> utrditi	16	17.28
<input type="checkbox"/> okreptiti	22	17.03
<input type="checkbox"/> zbijati	6	16.36
<input type="checkbox"/> pridobiti	40	16.27
<input type="checkbox"/> uživati	28	16.16
<input type="checkbox"/> spoštovati	20	16.14
<input type="checkbox"/> uveljavljati	8	15.77
<input type="checkbox"/> omajati	7	15.18
<input type="checkbox"/> razvaljati	6	15.01
<input type="checkbox"/> uveljavljati	15	14.86
<input type="checkbox"/> priznati	30	14.36

privedje	1236	2.4
<input type="checkbox"/> ugled	85	50.25
<input type="checkbox"/> verodostojnost	23	37.59
<input type="checkbox"/> moč	61	34.21
<input type="checkbox"/> karizma	14	33.61
<input type="checkbox"/> integriteta	12	29.9
<input type="checkbox"/> kredibilnost	11	28.86
<input type="checkbox"/> moč	21	27.57
<input type="checkbox"/> hierarhija	11	25.44
<input type="checkbox"/> legitimnost	9	25.3
<input type="checkbox"/> oblast	25	22.33
<input type="checkbox"/> znanje	18	19.51
<input type="checkbox"/> odgovornost	13	18.85
<input type="checkbox"/> zaupanje	10	18.47
<input type="checkbox"/> samozavest	7	17.71
<input type="checkbox"/> spoštovanje	8	16.85
<input type="checkbox"/> pooblastilo	7	16.68
<input type="checkbox"/> družben	11	15.91
<input type="checkbox"/> vpliv	11	15.22
<input type="checkbox"/> izkušnja	11	14.04
<input type="checkbox"/> položaj	12	13.7
<input type="checkbox"/> vloga	12	13.4
<input type="checkbox"/> tradicija	6	12.14
<input type="checkbox"/> red	7	10.23
<input type="checkbox"/> podoben	6	8.56

veznik	1086	2.4
<input type="checkbox"/> ki	483	41.1
<input type="checkbox"/> da	153	24.78
<input type="checkbox"/> saj	46	23.2
<input type="checkbox"/> in	35	22.02
<input type="checkbox"/> zato	27	21.07
<input type="checkbox"/> ampak	21	21.04
<input type="checkbox"/> temveč	16	20.78
<input type="checkbox"/> marveč	8	20.43
<input type="checkbox"/> kar	27	19.69
<input type="checkbox"/> kot	16	19.24
<input type="checkbox"/> kateri	36	18.99
<input type="checkbox"/> če	25	18.44
<input type="checkbox"/> naj	17	16.03
<input type="checkbox"/> kakršen	8	16.0
<input type="checkbox"/> toda	10	15.34
<input type="checkbox"/> ker	17	15.27
<input type="checkbox"/> ter	40	12.94
<input type="checkbox"/> kakor	7	14.74
<input type="checkbox"/> kadar	6	14.54
<input type="checkbox"/> čeprav	10	14.48
<input type="checkbox"/> vendar	13	13.91
<input type="checkbox"/> a	12	13.37
<input type="checkbox"/> ali	10	13.21
<input type="checkbox"/> drug	7	13.0
<input type="checkbox"/> ko	13	11.1

predi-za	1056	0.8
<input type="checkbox"/> nad	35	21.67
<input type="checkbox"/> na	246	21.43
<input type="checkbox"/> v	336	20.91
<input type="checkbox"/> med	49	16.87
<input type="checkbox"/> za	133	15.86
<input type="checkbox"/> pri	54	15.75
<input type="checkbox"/> znotraj	6	12.49
<input type="checkbox"/> kot	33	11.86
<input type="checkbox"/> glede	6	10.23
<input type="checkbox"/> z	56	8.0
<input type="checkbox"/> brez	8	7.21
<input type="checkbox"/> iz	20	7.11
<input type="checkbox"/> zaradi	9	6.11
<input type="checkbox"/> pred	8	4.6
<input type="checkbox"/> do	11	4.21
<input type="checkbox"/> od	9	3.16
<input type="checkbox"/> po	7	1.94
<input type="checkbox"/> po	8	1.3

osebek_od	1001	2.3
<input type="checkbox"/> priznavati	14	18.97
<input type="checkbox"/> omajati	6	15.35
<input type="checkbox"/> temeljiti	16	15.21
<input type="checkbox"/> obotajati	17	12.01
<input type="checkbox"/> izhajati	10	11.05
<input type="checkbox"/> tikati2	7	11.03
<input type="checkbox"/> spoštovati	7	9.56
<input type="checkbox"/> postati	22	8.76
<input type="checkbox"/> predstavljati	14	8.73
<input type="checkbox"/> odločati	9	8.25
<input type="checkbox"/> dati	29	8.22
<input type="checkbox"/> jemati	6	8.05
<input type="checkbox"/> ostajati	7	6.96
<input type="checkbox"/> delovati	9	5.81
<input type="checkbox"/> praviti	14	5.47
<input type="checkbox"/> meniti	8	5.44
<input type="checkbox"/> jesti	8	4.49
<input type="checkbox"/> dejati	9	4.43
<input type="checkbox"/> pomeniti	10	4.43
<input type="checkbox"/> pomagati	8	4.29
<input type="checkbox"/> imeti	39	3.52
<input type="checkbox"/> znati	8	3.43
<input type="checkbox"/> veljati	6	3.12
<input type="checkbox"/> kazati	6	3.02
<input type="checkbox"/> reči	7	2.26

predlog	965	0.6
<input type="checkbox"/> z	321	26.72
<input type="checkbox"/> brez	52	24.26
<input type="checkbox"/> do	87	21.53
<input type="checkbox"/> proti	31	19.74
<input type="checkbox"/> zaradi	29	15.21
<input type="checkbox"/> kot	42	14.45
<input type="checkbox"/> pred	30	14.18
<input type="checkbox"/> na	112	13.63
<input type="checkbox"/> zoper	6	12.33
<input type="checkbox"/> za	82	11.8
<input type="checkbox"/> od	23	9.08
<input type="checkbox"/> nad	8	8.89
<input type="checkbox"/> k	9	7.94
<input type="checkbox"/> pod	8	7.86
<input type="checkbox"/> med	15	7.58
<input type="checkbox"/> o	16	6.39
<input type="checkbox"/> po	17	5.23
<input type="checkbox"/> ob	6	2.97
<input type="checkbox"/> v	41	2.57
<input type="checkbox"/> iz	7	1.98
<input type="checkbox"/> pri	6	1.73

oba-v-rod	923	1.6
<input type="checkbox"/> moralen	79	51.4
<input type="checkbox"/> nesopren	26	40.09
<input type="checkbox"/> cerkven	22	29.73
<input type="checkbox"/> očetovski	9	28.1
<input type="checkbox"/> političen	43	27.73
<input type="checkbox"/> starševski	10	26.93
<input type="checkbox"/> nezmotljiv	6	25.06
<input type="checkbox"/> strokoven	29	24.35
<input type="checkbox"/> znanstven	14	22.82
<input type="checkbox"/> učiteljev	6	22.81
<input type="checkbox"/> lažen	9	22.31
<input type="checkbox"/> svetoven	31	22.19
<input type="checkbox"/> absoluten	9	21.59
<input type="checkbox"/> priznan	9	20.39
<input type="checkbox"/> potreben	17	20.18
<input type="checkbox"/> verski	11	19.92
<input type="checkbox"/> vrhoven	9	19.35
<input type="checkbox"/> religiozen	6	19.21
<input type="checkbox"/> velik	40	17.92
<input type="checkbox"/> priznan	6	17.21
<input type="checkbox"/> zadosten	6	17.92
<input type="checkbox"/> intelektualen	6	17.04
<input type="checkbox"/> pravi	12	16.86
<input type="checkbox"/> praven	11	15.8
<input type="checkbox"/> lasten	10	14.3
<input type="checkbox"/> tradicionalen	6	13.17

z-d_X	697	5.5
<input type="checkbox"/> težava	26	21.44
<input type="checkbox"/> biti	124	19.25
<input type="checkbox"/> zlovek	23	16.71
<input type="checkbox"/> konflikt	6	16.33
<input type="checkbox"/> obvladovati	8	13.82
<input type="checkbox"/> trener	6	12.27
<input type="checkbox"/> oseba	8	12.01
<input type="checkbox"/> problem	7	11.66
<input type="checkbox"/> moč	6	11.3
<input type="checkbox"/> podpreti	8	9.48
<input type="checkbox"/> poskrbeti	8	7.65
<input type="checkbox"/> voditi	11	7.38
<input type="checkbox"/> vplivati	7	7.28
<input type="checkbox"/> govoriti	10	6.43
<input type="checkbox"/> povezati	7	6.41

v-rodil-s	653	1.8
<input type="checkbox"/> učitelj	28	30.97
<input type="checkbox"/> stvarnik	7	26.46
<input type="checkbox"/> stariš	21	24.96
<input type="checkbox"/> Cerkev	10	22.44
<input type="checkbox"/> država	38	22.05
<input type="checkbox"/> sodstvo	7	21.87
<input type="checkbox"/> argument	10	21.32
<input type="checkbox"/> predsednik	20	18.49
<input type="checkbox"/> institucija	9	18.45
<input type="checkbox"/> papež	6	17.29
<input type="checkbox"/> sodišče	15	16.93
<input type="checkbox"/> cerkev	10	16.38
<input type="checkbox"/> bog	6	15.94
<input type="checkbox"/> parlament	8	14.88
<input type="checkbox"/> kota	8	14.64
<input type="checkbox"/> znanje	8	14.01
<input type="checkbox"/> vlada	11	13.36
<input type="checkbox"/> trener	6	12.45
<input type="checkbox"/> oče	6	12.38
<input type="checkbox"/> vodja	6	12.25
<input type="checkbox"/> zdravnik	6	12.0
<input type="checkbox"/> agencija	6	11.78
<input type="checkbox"/> položaj	7	11.4
<input type="checkbox"/> oblast	6	11.12
<input type="checkbox"/> zakon	6	7.88

s-koga-česa	651	1.8
<input type="checkbox"/> priznavanje	13	33.02
<input type="checkbox"/> pomanjkanje	23	30.01
<input type="checkbox"/> spoštovanje	19	29.48
<input type="checkbox"/> spodkopavanje	7	29.4
<input type="checkbox"/> nesoprotovanje	9	27.92
<input type="checkbox"/> krepitev	7	20.75
<input type="checkbox"/> argument	9	20.14
<input type="checkbox"/> kriza	12	20.13
<input type="checkbox"/> hierarhija	6	20.02
<input type="checkbox"/> izguba	11	19.93
<input type="checkbox"/> struktura	10	18.34
<input type="checkbox"/> uveljavljanje	6	18.18
<input type="checkbox"/> zloraba	7	17.93
<input type="checkbox"/> oblika	15	17.61
<input type="checkbox"/> pozicija	6	17.5
<input type="checkbox"/> prava	6	16.49
<input type="checkbox"/> moč	9	14.86
<input type="checkbox"/> vprašanje	12	14.48
<input type="checkbox"/> problem	6	12.96
<input type="checkbox"/> podlaga	6	12.14
<input type="checkbox"/> vpliv	6	12.02

koga-česa	425	7.0
<input type="checkbox"/> priznati	48	26.89
<input type="checkbox"/> priznavati	17	24.63
<input type="checkbox"/> spodbijati	7	20.82
<input type="checkbox"/> imeti	105	16.88
<input type="checkbox"/> rušiti	6	15.58
<input type="checkbox"/> spoštovati	9	14.27
<input type="checkbox"/> prenesti	10	13.55
<input type="checkbox"/> uveljaviti	7	11.58
<input type="checkbox"/> ustvariti	8	9.98
<input type="checkbox"/> izgubiti	6	6.58
<input type="checkbox"/> znati	7	5.41

komu-čemu	239	6.4
<input type="checkbox"/> upirati	24	30.82
<input type="checkbox"/> podrejati	12	28.71
<input type="checkbox"/> upreti	15	23.79
<input type="checkbox"/> ukloniti	8	23.72
<input type="checkbox"/> podrediti	10	21.31
<input type="checkbox"/> škoditi	7	17.09
<input type="checkbox"/> nasprotovati	7	13.63
<input type="checkbox"/> zaupati	8	13.48

za-d_X	179	1.3
<input type="checkbox"/> veljati	55	29.32
<input type="checkbox"/> skrivati	10	16.39
<input type="checkbox"/> biti	19	9.71
<input type="checkbox"/> iti	13	

Transfer of information

- API using data from Sketch Engine
- Gramrels:
 - Element <struktura> = syntactic structures
 - Element <vzorec> = syntactic patterns
 - Element <zveza> = syntactic combinations
 - Element <oznaka> = labels
- Collocations = element <kolokacija>
- Examples = element <zgled> using GDEX

Gramrel to <struktura>

ADJECTIVE + NOUN

<skladenjska struktura>

<struktura>kakšen?</struktura>

<kolokacije>

<kolokacija id="839596"><k>nov</k></kolokacija>

<kolokacija id="839746"><k>deloven</k></kolokacija>

<kolokacija id="840017"><k>spleten</k></kolokacija>

<kolokacija id="839637"><k>glaven</k></kolokacija>

<kolokacija id="839725"><k>prost</k></kolokacija>

<kolokacija id="839830"><k>parkiren</k></kolokacija>

<kolokacija id="839601"><k>velik</k></kolokacija>

<kolokacija id="839952"><k>vodilen</k></kolokacija>

<kolokacija id="839625"><k>pravi</k></kolokacija>

<kolokacija id="839814"><k>prodajen</k></kolokacija>

</kolokacije>

<zgledi>

<zgled seek="839596" position="1">Zavod za zdravstveno varstvo Novo
<i>mesto</i></zgled>

<zgled seek="839601" position="1">" V glavnem v vseh večjih <i>mestih</i>
</i>.</zgled>

collocations and corresponding examples

Gramrel to <vzorec>

Construction to <vzorec>

```
|<skladenjska_struktura>  
|<vzorec>S_V_03_04</vzorec>  
<zgledi>
```

```
<zgled seek="16213" position="1">Tako združene moci  
so tovarni <i>dale</i> nov zagon in postala je  
najboljša tovarna klobukov.</zgled>
```

```
<zgled seek="16215" position="1">Njen predsednik Borut  
Pahor je Drnovšku ponovno <i>dal</i> košarico ne sicer  
za vecno, ampak vsaj do prvih naslednjih  
volitev.</zgled>
```

```
<zgled seek="16215" position="2">Južnoafriški zunanji  
minister Alfred Nzo je bil po konferenci še posebej  
zadovoljen, ker so neuvrščene države <i>dale</i> vso  
prednost jedrski razorožitvi.</zgled>
```

```
</zgledi>
```

Gramrel to label

<oblika>

<iztocnica>mesto</iztocnica>

</oblika>

unary to label: "with proper names"

<zaglavje>

<besvrs>samostalni</besvrs>

<oznaka>z_lastnim_imenom</oznaka>

</zaglavje>

Sketch Grammar for Slovene (v.16)

- 105 gramrels
- 50 macro definitions
- 25 DUAL
- 36 TRINARY
- 36 SEPARATEPAGE
- 8 UNARY
- 2 SYMMETRIC
- 19 CONSTRUCTION
- 3 COLLOC
- 18 no directive
- SEPARATEPAGE=TRINARY
- 6 CONSTRUCTION-UNARY

API and settings

- API script to extract data from word sketch information in the Sketch Engine
- a list of lemmas for extraction: lemmas with frequency between 1000 (0.85 per million words) and 10,000 (8.5 per million words)
- settings for extraction
 - lemmas divided into five frequency groups
 - different setting for each group

Selection of lemmas

- Frequent enough to offer a good-sized word sketch
 - less than 600 hits in Gigafida did not provide enough relevant data
 - we divided lemmas of each word class into five different frequency groups
- Monosemous lemmas or having up to
 - two synsets/senses in sloWNet, a Slovene version of Wordnet
 - exceptionally, in the Dictionary of Standard Slovenian (SSKJ)
- Found in sloWnet, preferably, but not in SSKJ, as we wanted to focus on new words and/or senses

Distribution of lemmas

- The final selection included
 - 515 nouns
 - 260 verbs
 - 275 adjectives
 - 117 adverbs

 - lemmas with frequency between 1000 (0.85 per million words) and 10,000 (8.5 per million words)

API script (Python)

- Python for Windows (v. 2.7.3)
- httplib2-0.7.4 library
- simplejson-2.5.2 library
- **python tblscript.py fidaplus-gf -U
http://ske.slovenscina.eu/ -f 8 -l
spisek_lempos.txt**

Lemmalist

- -l LEMMALIST, --lemmalist=LEMMALIST
 - The file containing a list of lemposes for which the examples are to be extracted (stdin by default).

General (Gramrellist)

- **-f MINFREQ, --frequency=MINFREQ**
 - Default minimum frequency of a collocate(default=0.0).
- **-s MINSAL, --salience=MINSAL**
 - Default minimum salience of a collocate(default=0.0).
- **-F MINFREQREL, --Freqrel=MINFREQREL**
 - Minimum frequency of a relation (default=25).
- **-S MINSALREL, --Salrel=MINSALREL**
 - Minimum salience of a relation (default=0.0).

Gramrellist

- -r GRAMRELLIST, --relations=GRAMRELLIST
 - The file containing a set of grammatical relations from a given sketch grammar for inclusion (all by default).
 - One record consists of:
 - gramrel regular expression
 - min. collocation frequency
 - min. col. salience
 - min. gramrel frequency
 - min. g. salience
 - gramrel type
 - The gramrel type should be one of: 'SVOZ' in order: 'struktura', 'vzorec', 'oznaka' and 'zveza'. If no type is provided than the first letter of gramrel name decides. For example:
 - (sub|ob)ject 3 2.5 30 20 S

Maximums & GDEX

- **-n NUMBER, --number=NUMBER**
 - Maximum number of sentences per collocation (default=6).
- **-m MAXITEMS, --maxCollocs=MAXITEMS**
 - Maximum number of collocations per grammatical relation (default 10).
- **-g GDEXCONF, --gdexconf=GDEXCONF**
 - Name of the gdex configuration to use.

Gramrellist example

gramrel regular expression	min. coll. freq	min. coll. salience	min. gramrel freq	min. gramrel salience	gramrel type
...					
O_tretja_oseba	8	0.5	60	0.5	O
O_z_lastnim_imenom	8	0.5	8	2.5	O
O_zanikanje	8	0.5	8	20.0	O
S_.*_p2	4	0.5	8	25.0	S
S_.*_p3	4	0.5	8	100.0	S
S_.*_p4	4	0.5	8	20.0	S
...					

We started with...

- 10 collocates per relation
- 6 examples per collocate
- Minimum salience of a relation/collocate = 0
- Minimum frequency of a collocate = 0
- Minimum frequency of a relation = 25

- Statistical & manual analysis
- identifying the lowest values where the collocation still yielded relevant results

And ended with...

- Minimum number of collocates per relation was increased to 25
- Selection of relevant collocates was 'left' to minimum frequency and salience settings
- Number of examples per collocate was reduced to three
- We divided lemmas into frequency groups, and prepared separate settings for each group

XML template

- DOC_TEMPLATE = (""""<?xml version="1.0" encoding="UTF-8"?>
 - <clanek>
 - <glava>
 - <oblika><zapis>%(headword)s</zapis>
 - <iztocnica>%(headword)s</iztocnica></oblika>
 - <zaglavje>
 - <besvrs>%(pos)s</besvrs>
 - """" ,# here come all O_ """"
 - </zaglavje>
 - </glava>

Output

- `?xml version="1.0" encoding="UTF-8"?>`
- `<clanek>`
- `<glava>`
- `<oblika><zapis>anoreksija</zapis><iztocnica>anoreksija</iztocnica></oblika>`
- `<zaglavje><besvrs>samostalnik</besvrs></zaglavje>`
- `</glava>`
- `<geslo>`
- `<pomen>`
- `<indikator></indikator><pomenska_shema></pomenska_shema>`
- `<skladenjske_skupine><skladenjska_struktura>`
- `<struktura>S_predl-pred</struktura>`
- `<kolokacije><kolokacija kid="100344429"><k>proti</k></kolokacija></kolokacije>`
- `<zgledi><zgled kid="100344429" pozicija="1">Francoska manekenka, ki je leta 2007 s fotografijo v okviru kampanje boja proti <i id="1338652551">anoreksiji</i> dvignila veliko prahu, je umrla.</zgled></zgledi>`

Gramrel conversion

- `?xml version="1.0" encoding="UTF-8"?>`
- `<clanek>`
- `<glava>`
- `<oblika><zapis>anoreksija</zapis><iztocnica>anoreksija</iztocnica></oblika>`
- `<zaglavje><besvrs>samostalnik</besvrs></zaglavje>`
- `</glava>`
- `<geslo>`
- `<pomen>`
- `<indikator></indikator><pomenska_shema></pomenska_shema>`
- `<skladenjske_skupine><skladenjska_struktura>`
- `<struktura>zveze s predlogi</struktura>`
- `<kolokacije><kolokacija kid="100344429"><k>proti</k></kolokacija></kolokacije>`
- `<zgledi><zgled kid="100344429" pozicija="1">Francoska manekenka, ki je leta 2007 s fotografijo v okviru kampanje boja proti <i id="1338652551">anoreksiji</i> dvignila veliko prahu, je umrla.</zgled></zgledi>`

“it is more efficient to edit out the computer’s errors than to go through the whole data-selection process from the beginning”

(Rundell & Kilgarriff, 2011)

“too many choices early in the data-selection process leave more room for error”

(Kosem, Gantar & Krek)

Work left for lexicographers

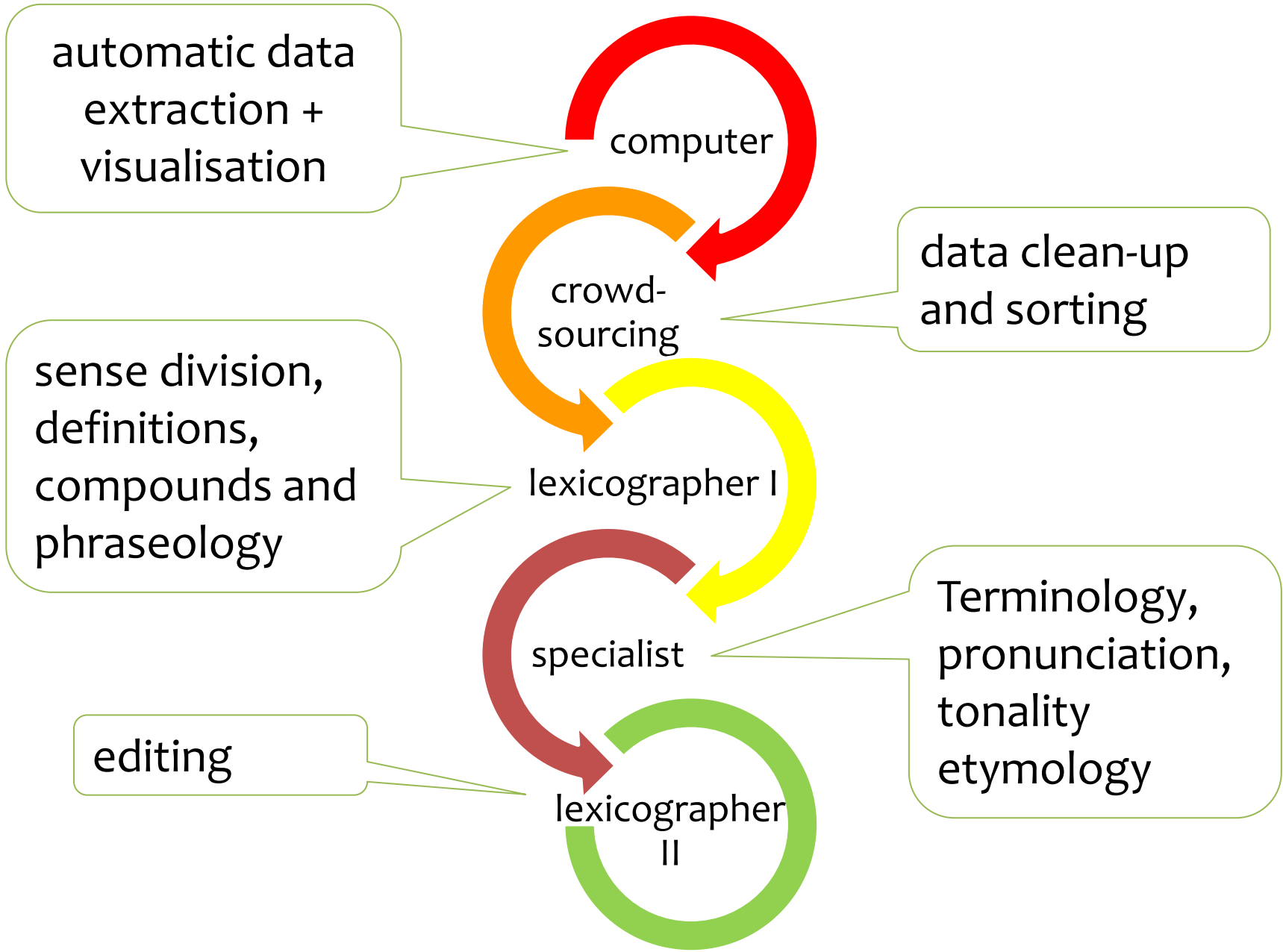
- Analytical
 - sense division
 - writing definitions, sense indicators
 - identification of multi-word units, phrases, pragmatics
 - adding certain labels
- Editorial
 - distributing information according to sense division
 - copying grammatical relations and collocates typical for more than one sense
 - deleting irrelevant info (collocates, examples etc.)



Proposal for a dictionary of contemporary Slovene

- starting from scratch
- <http://www.sssj.si/>)
- multi-phase dictionary compilation (and presentation)
- data extraction used for the first phase





Ocenjevanje slovnične ustreznosti besednih kombinacij

V tej nalogi vas prosimo, da ocenite, ali kombinacija besed v zgledu ustreza navedeni slovnični strukturi. S pravilnimi odgovori boste iz spletnega slovarja odstranili zglede, v katerih besedne kombinacije ne ustrezajo slovničnim strukturam, pod katere so bile uvrščene na podlagi avtomatskega postopka. Pozorni morate biti predvsem na pripis besedne vrste, sklona in stavčne vloge pri kateri od obarvanih besed v zgledu.

Ali kombinacija besed v zgledu ustreza navedeni slovnični strukturi?

Beseda

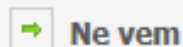
franšiza - *samostalnik*

Slovnična struktura

glagol + **za** + **samostalnik v tožilniku**

Zgled

Vsak poslovni sistem - ne glede na to, ali gre za franšizo ali ne - ima svoj cilj oziroma poslanstvo, ki vam lahko ustreza ali pa ne.



SLOVAR SODOBNEGA SLOVENSKEGA JEZIKA

1. 4. 2013

Pomen

Oblike

Sinonimi

Govor

Vizualizacija

Statistika

globalen pridevnik

PRIDEVNIK + SAMOSTALNIK

1. 4. 2013

Pomen

Oblike

Sinonimi

Govor

Vizualizacija

Statistika

Print | Save | Email | Cite

Text size: A A

Options: [Show all](#) | [Hide all](#)

This entry has been updated (OED Third Edition, September 2013).

[Publication history](#)

[Entry profile](#)

[Previous version](#)

z novimi trendi zahodne oblačilne mode, je oblačenje vedno vpeto v ni, transnacionalnimi in **globalnimi** normami lepote, spola in

družitev lokalnega in **globalnega**, avstralsko nacionalnega in

- 🔊 Njih se tičejo "nacionalna, regionalna in **globalna** varnost", ki ju zagotavlja vojaški stroj.
- 🔊 Proces izobraževanja in usposabljanja zdravnikov postaja vedno bolj celovit, univerzalen in **globalen**.
- 🔊 Poleg te » **globalne** in pavšalne odškodnine« je bila jugoslovanska stran dolžna Italiji s »seznama za prosto razpolaganje« vrniti še 179 nepremičnin, nekoč odvzetih italijanskim upravičencem.

1. 4. 2013

Pomen

Oblike

Sinonimi

Govor

Vizualizacija

Statistika

globalen pridevnik



P 3000

PRIDEVNIK + SAMOSTALNIK

globalno [segrevanje, zasedanje, ogrevanje]

globalna [ekonomija, kriza, delniska, konkurenčnost, recesija, raven]

globalna [razsežnost, korporacija, otoplitev, spremenljivka, vas]

globalna [pravičnost]

globalni [kapitalizem, vodič, nomad, trend, klepet]

globalni [trg, fenomen, izziv, terorizem, motorist]

- 🔊 Danes ni več politika ali politične stranke, ki v svojih nastopih in programih ne bi omenjala tudi klimatskih sprememb in zaskrbljenosti zaradi **globalnega** segrevanja ozračja, skrbi za pitno vodo in zdravo okolje.
- 🔊 Za podjetja je postala zaradi soodvisnosti in povezanosti **globalne** ekonomije nujna za preživetje.
- 🔊 Spremljamo dogajanje na ljubljanski borzi, ki je zaradi **globalne** finančne krize teden začela v rdečih številkah.
- 🔊 Brez podkupnine, ne mešajte s provizijo, v **globalnem** kapitalizmu ni več resnega posla.
- 🔊 Zaradi pospeševanja **globalne** konkurenčnosti slovenskega gospodarstva pa bi bilo treba to strateško usmeritev uresničiti brez dodatnih fiskalnih bremen.

PRIDEVNIK + in + PRIDEVNIK

globalen in [transnacionalen, lokalni, regionalni, univerzalni, pavšalni]

globalen in [meddržaven, nacionalni, celovit, radikalen, operativni]

globalen in [regijski, dolgoročni, evropski, konkurenčni, trajni]

globalen in [krajveni, domači, mednarodni, splošni]

- 🔊 Medtem ko lepota stalno tekmuje z novimi trendi zahodne oblačilne mode, je oblačenje vedno vpeto v zapletena pogajanja med lokalnimi, transnacionalnimi in **globalnimi** normami lepote, spola in

1. 4. 2013

Pomen

Oblike

Sinonimi

Govor

Vizualizacija

Statistika

globalen pridevnik



P 3000

PRIDEVNIK + SAMOSTALNIK

globalno [[segrevanje](#), [ogrevanje](#)]

globalna [[ekonomija](#), [kriza](#), [konkurenčnost](#), [recesija](#), [raven](#)]

globalna [[razsežnost](#), [korporacija](#), [otoplitev](#), [spremenljivka](#), [vas](#)]

globalna [[pravičnost](#)]

globalni [[kapitalizem](#), [nomad](#), [trend](#)]

globalni [[trg](#), [fenomen](#), [izziv](#), [terorizem](#)]

- 🔊 *Danes ni več politika ali politične stranke, ki v svojih nastopih in programih ne bi omenjala tudi klimatskih sprememb in zaskrbljenosti zaradi **globalnega** segrevanja ozračja, skrbi za pitno vodo in zdravo okolje.*
- 🔊 *Svet in z njim **globalna** ekonomija dejansko vodita planet Zemljo v fazo prekomernega izkoriščanja naravnih virov.*
- 🔊 *Spremljamo dogajanje na ljubljanski borzi, ki je zaradi **globalne** finančne krize teden začela v rdečih številkah.*
- 🔊 *Brez podkupnine, ne mešajte s provizijo, v **globalnem** kapitalizmu ni več resnega posla.*
- 🔊 *Zaradi pospeševanja **globalne** konkurenčnosti slovenskega gospodarstva pa bi bilo treba to stratešk usmeritev uresničiti brez dodatnih fiskalnih bremen.*

PRIDEVNIK + in + PRIDEVNIK

globalen in [[transnacionalen](#), [lokalen](#), [regionalen](#), [univerzalen](#), [pavšalen](#)]

globalen in [[meddržaven](#), [nacionalen](#), [celovit](#), [radikalen](#), [operativen](#)]

globalen in [[regijski](#), [dolgoročen](#), [evropski](#), [konkurenčen](#), [trajen](#)]

globalen in [[krajeven](#), [domač](#), [mednaroden](#), [splošen](#)]

- 🔊 *Medtem ko lepota stalno tekmuje z novimi trendi zahodne oblačilne mode, je oblačenje vedno vpeto*

1. 4. 2013

Pomen

Oblike

Sinonimi

Govor

Vizualizacija

Multimedija

Statistika

globalen pridevnik



P 3000

1. svetovni; mednarodni

- 1.1 splošno veljaven; razširjen
- 2. zemeljski; planetarni
- 3. ki zadeva celoto; celostni

1. svetovni; mednarodni

globalni procesi, zlasti gospodarski in politični, zajemajo ves svet

- 🔊 V New Yorku naj bi državniki in podjetniki razpravljali o **globalni** varnosti.
- 🔊 Tudi največje svetovne firme, ki danes obvladujejo **globalni** trg, so se razvile iz malih podjetij in obratov.
- 🔊 Motorola je zaradi **globalne** recesije v visokotehnoloških gospodarskih panogah lani odpustila 48.400 zaposlenih.

1.1 splošno veljaven; razširjen

če postanejo neke dejavnosti ali lastnosti globalne, jih upošteva vedno več ljudi ali držav po svetu

- 🔊 Merila, kakšna ženska je lepa, postajajo vse bolj **globalna**.

2. zemeljski; planetarni

ekologija

globalne spremembe v okolju vplivajo na celoten zemeljski planet

- 🔊 Eden najbolj preprostih in praktično izvedljivih načinov za zmanjšanje **globalnega** segrevanja je sajenje novih dreves.
- 🔊 Krčenje ledenikov in spremembe v pokrajini sta zelo očitni in lahko razumljivi posledici **globalnega** ogrevanja.
- 🔊 Vodik bo tudi občutno skrčil emisije ogljikovega dioksida in učinke **globalne** otoplitve.

1. 4. 2013

Pomen

Oblike

Sinonimi

Izvor

Govor

Vizualizacija

Multimedija

Statistika

globalen pridevnik



/globálen/

P 3000

1. svetovni; mednarodni

1.1 splošno veljaven; razširjen

2. zemeljski; planetarni

3. ki zadeva celoto; celostni

1. svetovni; mednarodni

globalni procesi, zlasti gospodarski in politični, zajemajo ves svet

- 🔊 V New Yorku naj bi državniki in podjetniki razpravljali o **globalni** varnosti.
- 🔊 Tudi največje svetovne firme, ki danes obvladujejo **globalni** trg, so se razvile iz malih podjetij in obratov.
- 🔊 Motorola je zaradi **globalne** recesije v visokotehnoloških gospodarskih panogah lani odpustila 48.400 zaposlenih.

1.1 splošno veljaven; razširjen

če postanejo neke dejavnosti ali lastnosti globalne, jih upošteva vedno več ljudi ali držav po svetu

- 🔊 Merila, kakšna ženska je lepa, postajajo vse bolj **globalna**.

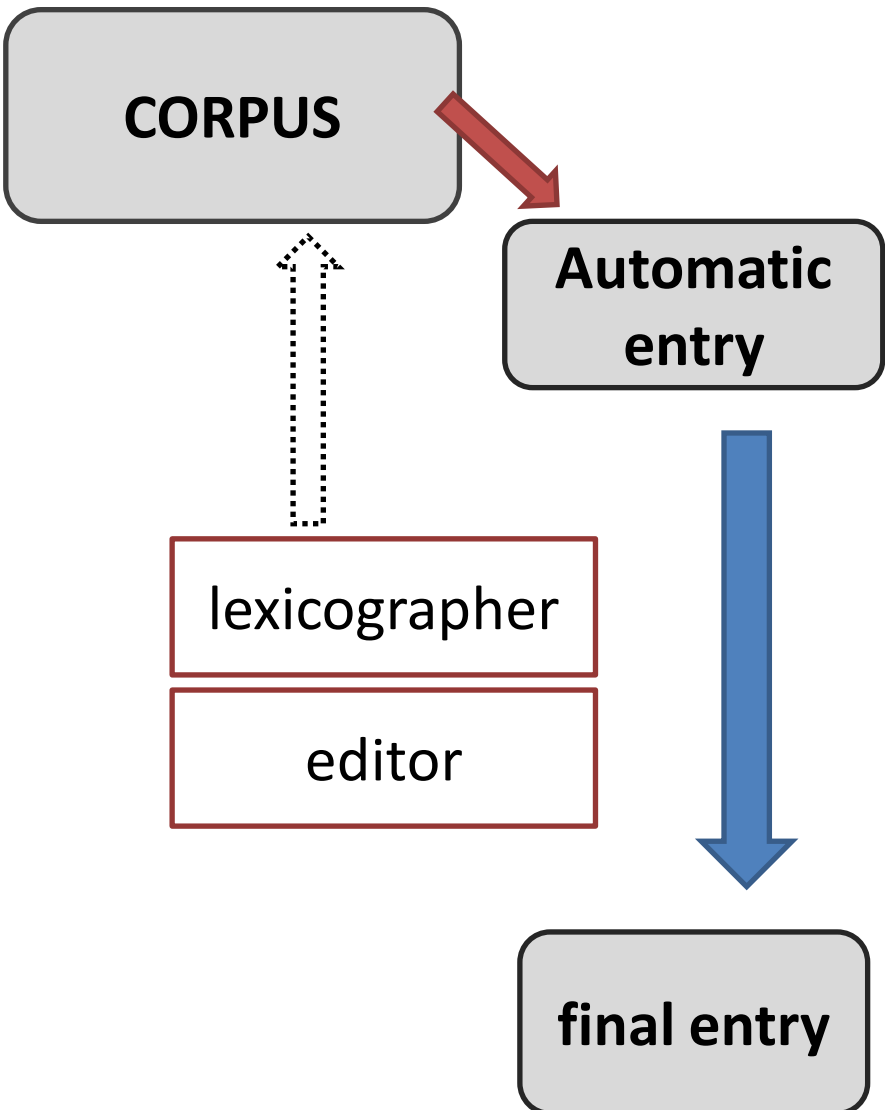
2. zemeljski; planetarni

ekologija

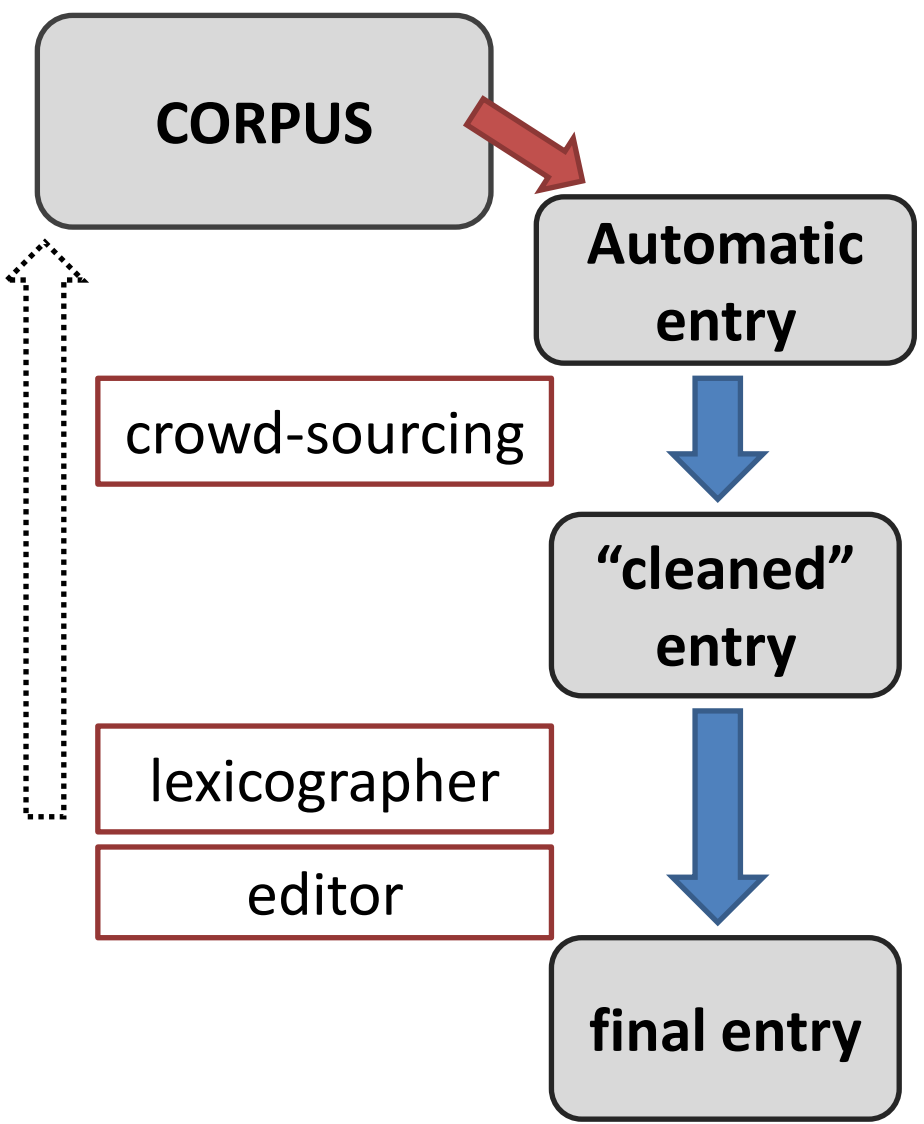
globalne spremembe v okolju vplivajo na celoten zemeljski planet

- 🔊 Eden najbolj preprostih in praktično izvedljivih načinov za zmanjšanje **globalnega** segrevanja je sajenje novih dreves.
- 🔊 Krčenje ledenikov in spremembe v pokrajini sta zelo očitni in lahko razumljivi posledici **globalnega** ogrevanja.
- 🔊 Vodik bo tudi občutno skrčil emisije ogljikovega dioksida in učinke **globalne** otoplitve.

Rundell & Kilgarriff (2011)



our proposal



Consortium

- University of Ljubljana
 - Faculty of Arts
 - Faculty of Social Sciences
 - Faculty of Pedagogy
 - Faculty of Computer Sciences and Informatics
 - Faculty of Electrotechnics
- University of Maribor
- University of Primorska
- "Jožef Stefan" Institute
- Amebis, Alpineon, Trojina
- <http://www.cjvt.si/projekti>

Plan

- 50 researchers
- Upgraded concept & proposal: spring 2015
- Already confirmed
 - crowdsourcing
 - automatic data extraction
 - multi-phase publication
- Work on
 - technical aspects
 - lexicographical aspects
 - design