



# Kolokacijski slovar sodobne slovenščine

Iztok Kosem,<sup>1,2</sup> Simon Krek,<sup>2</sup> Polona Gantar,<sup>1</sup> Špela Arhar Holdt,<sup>1,3</sup>  
Jaka Čibej,<sup>1,2,3</sup> Cyprian Laskowski,<sup>1</sup>

<sup>1</sup>Filozofska fakulteta, Univerza v Ljubljani

<sup>2</sup>Institut "Jožef Stefan"

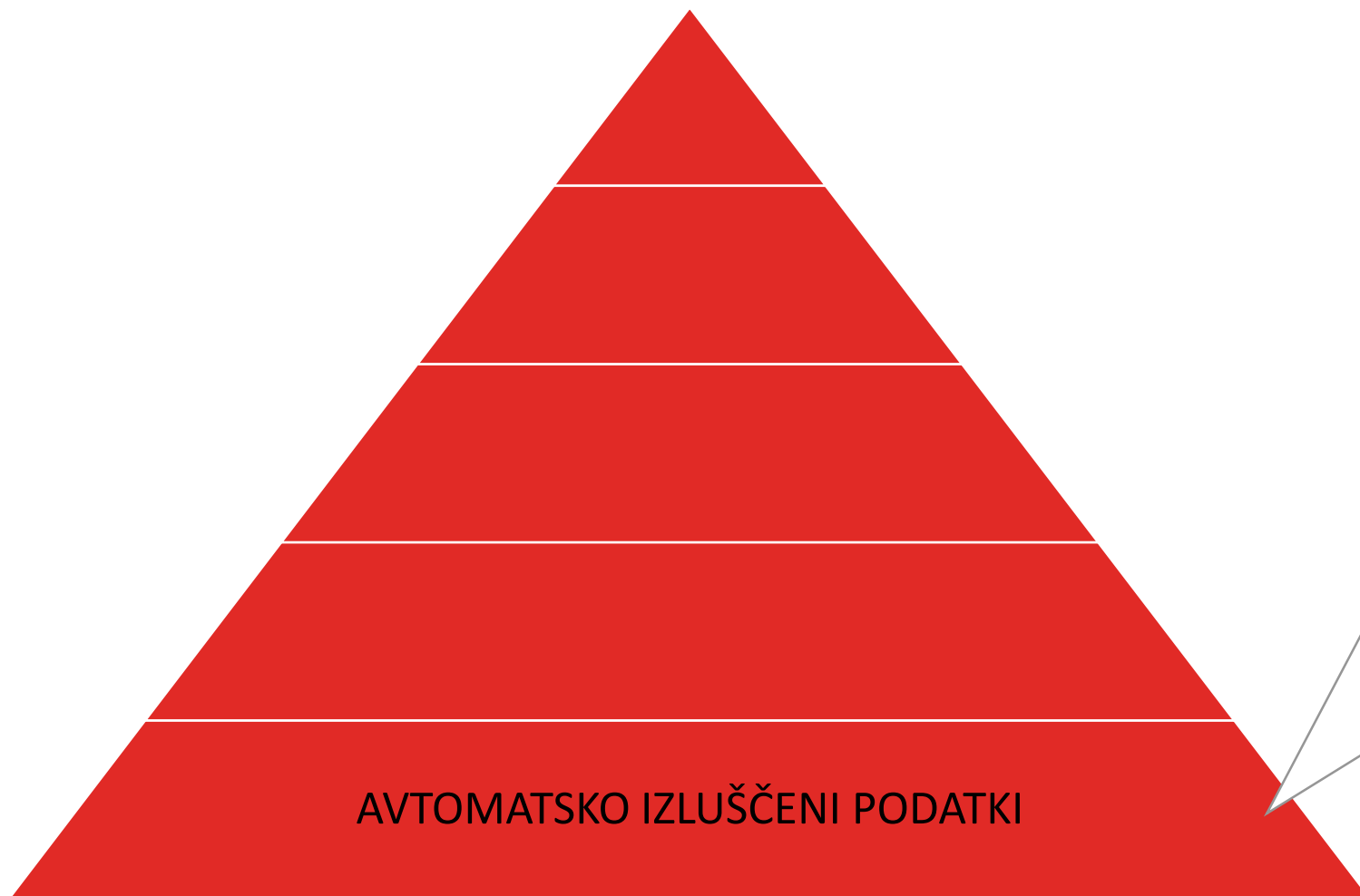
<sup>3</sup>Fakulteta za računalništvo in informatiko, Univerza v Ljubljani

# Modeli objave slovarjev

- Tradicionalni model: dokončaj vse, nato objavi
  - PROBLEM: čas do objave, kdaj je nekaj dokončano?
- Rastoči model: objava dokončanih gesel sproti (npr. Veliki poljski slovar) ali v rednih intervalih
  - PROBLEM: traja več let (ali desetletij?), preden količina gesel doseže uporabno vrednost za uporabnike
- Odzivni stopenjski model: objava vseh gesel v različnih stopnjah dokončanosti (Krek et al. 2013)
  - PROBLEM: zanesljivost?

# Kolokacijski slovar sodobne slovenščine

- Verzija 0.9:
  - 35.989 iztočnic
  - 7.717.561 kolokacij
  - 36.736.168 zgledov
- Podatki izluščeni iz korpusa Gigafida



Sketch Engine API:

- Slovnica besednih skic (Krek 2011, Krek 2015)
- GDEX (Kosem et al. 2011, Kosem et al. 2013, Kosem 2015)

Leksikografski postopek:

Krek et al. 2013

Gantar et al. 2016

Gorjanc et al. 2015, 2016



Leksikalnogramatični in statistični filtri, npr. glagol „biti“, problematične strukture

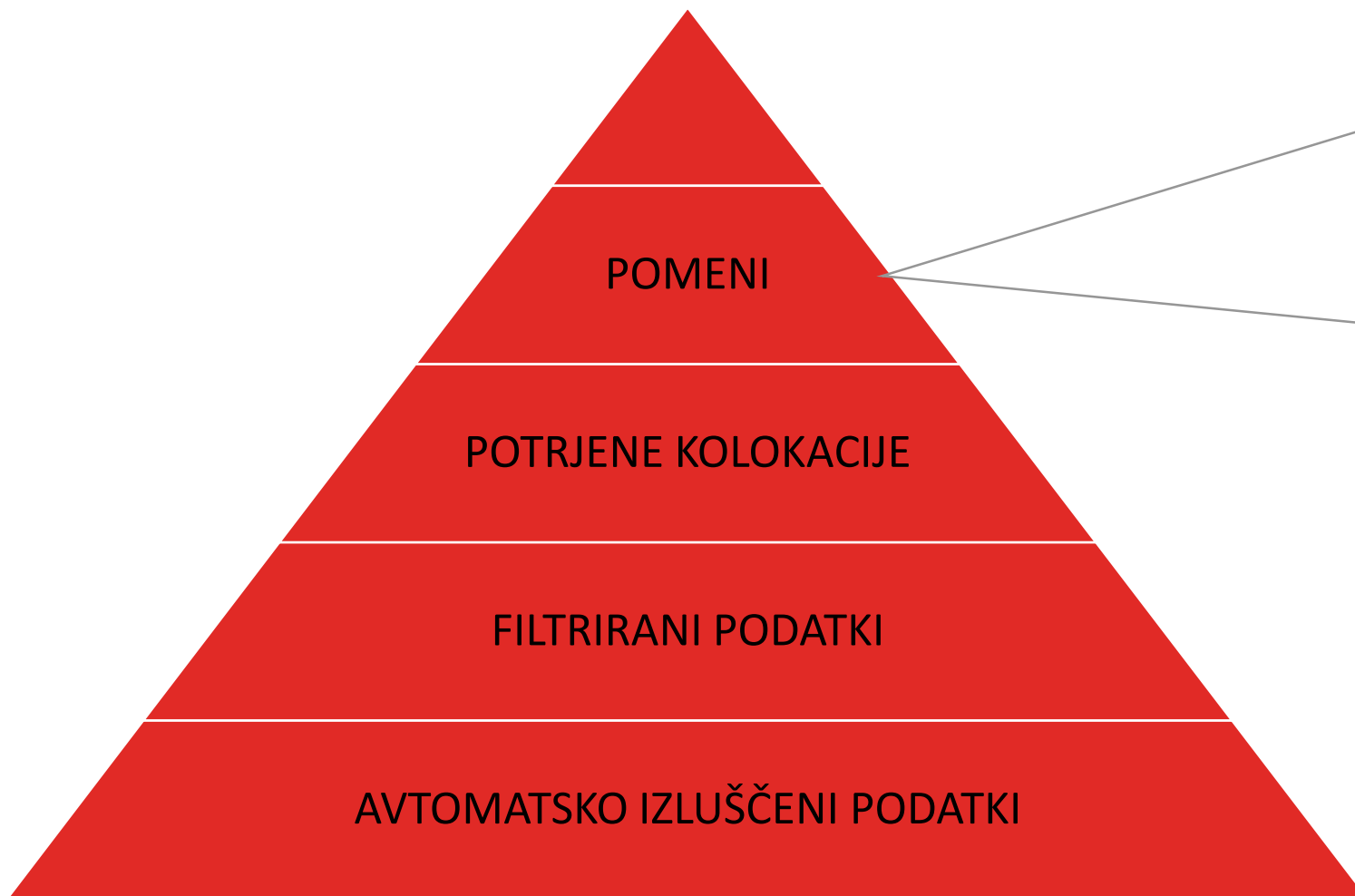


IZLOČAMO V CELOTI:

- nekolokacije
- napake strukture

IZLOČAMO,  
A OBDRŽIMO V BAZI:

- semantično manj relevantne kolokacije



Kolokacije in pripadajoče zglede razvrstimo pod pomene

Uporaba množičenja (deluje!)

**obiskati bazar**

Ogledali si bomo mesto in *obiskali bazar*.

orientalska tržnica

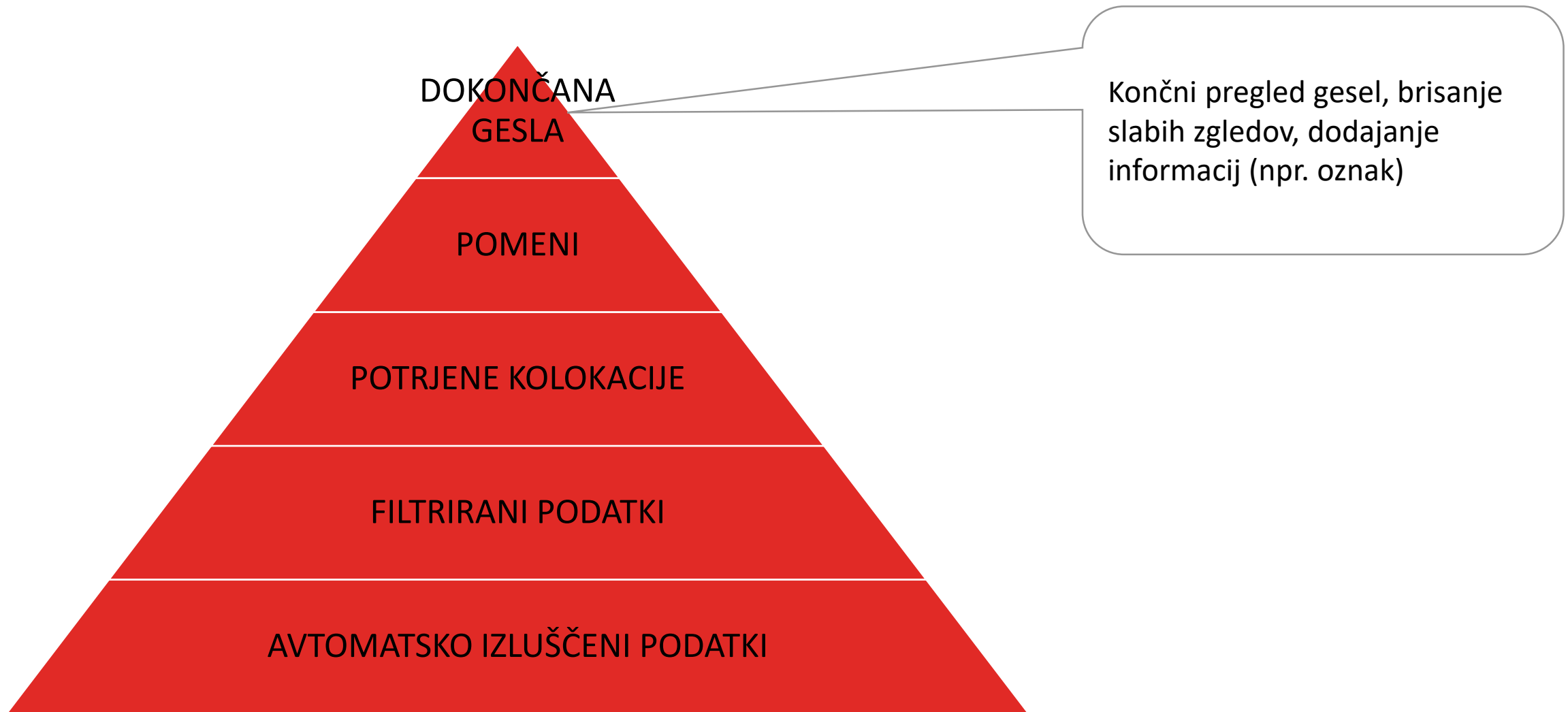
priredivtev

Nič od naštetega

Ne vem

Trenutno rešujete nalogo **12**.

Rešili ste naslednje število nalog: **0** od skupno **1**





# Vmesnik

- Glavne značilnosti:
  - Izhodišče so kolokacije in ne pomeni
  - Jasna informacija o stanju gesla in datumu zadnje posodobitve



Relevantnost | Gruče | A-Ž

Pomeni

**Strukture**

s samostalniki

z glagoli

s pridevniki

**Predlogi**

**cjvt** kolokacije 0.9

uporabnost

stranska



uporabnost (samostalnik) 2018-07-10



aljšati



Relevantnost | Gruče | A-Ž

Pomeni

**Strukture**

nja



praktična uporabnost

kter



uporabnost za prakso

za pripravljenost

za uporabnika

za vožnjo



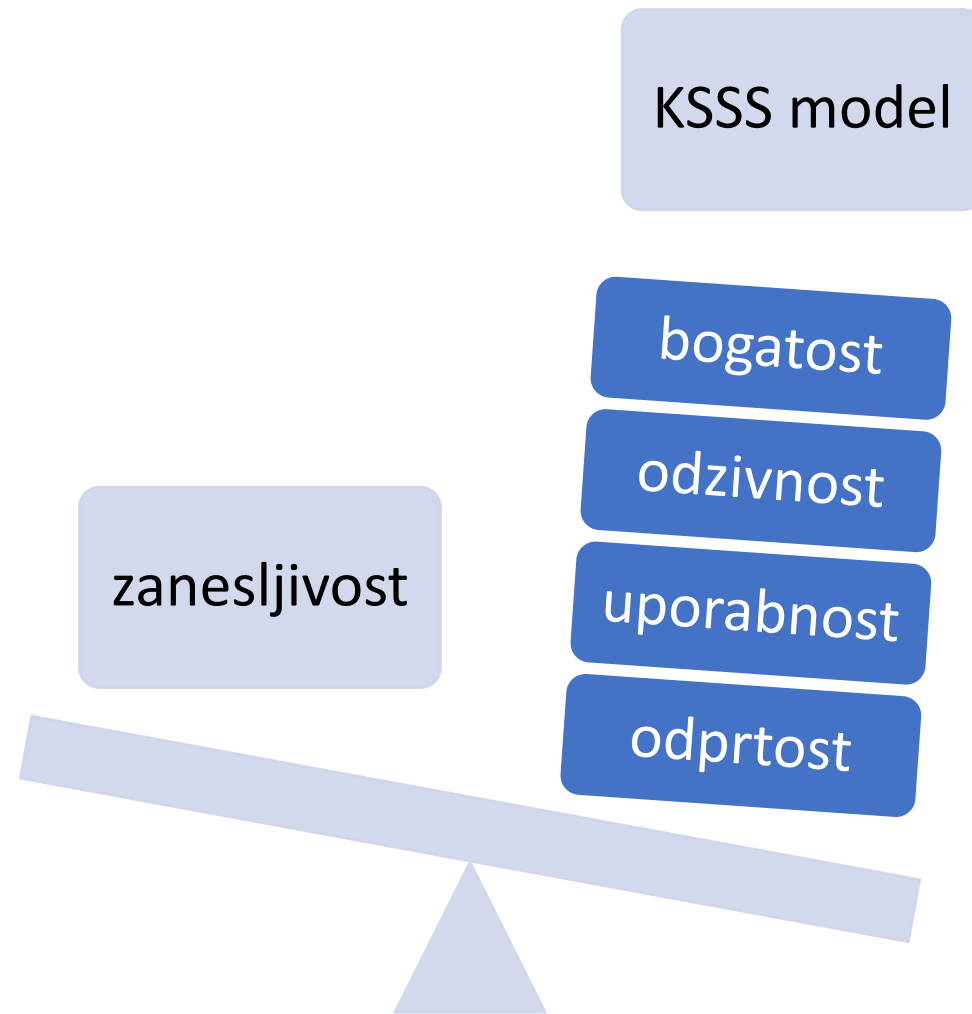
ocena za uporabnost

strokovnjak za

pogoj za

# Vmesnik

- Glavne značilnosti:
  - Izhodišče so kolokacije in ne pomeni
  - Jasna informacija o stanju gesla in datumu zadnje posodobitve
  - Veliko možnosti razvrščanja in filtriranja
  - Prilagojenost za različne digitalne medije



# Projekt Kolokacije v slovenščini (KOLOS)



# Kolokativnost

struktura	% of good	% of noisy	Agreement	same answers	% nestrinjanja ali dva Ne vem
PRID + SAM	88.9	3.1	0.42	0.78	11,8
SAM + SAM <sub>2</sub>	84.9	7.8	0.45	0.63	13,3
GLAG + SAM <sub>4</sub>	87.0	10.0	0.46	0.73	16,1
PRISL + GLAG	87.7	1.3	0.37	0.63	19,5
PRISL + PRID	63.6	20.7	0.46	0.64	15,9
SAM + v + SAM <sub>5</sub>	64.2	0.5	0.35	0.50	39,0
GLAG + PRISL	59.6	7.3	0.47	0.61	19,4
GLAG + v + SAM <sub>5</sub>	86.0	2.1	0.33	0.46	26,4
SAM + s/z + SAM <sub>6</sub>	74.7	0.2	0.42	0.56	23,3
GLAG + s/z + SAM <sub>6</sub>	92.7	1.3	0.46	0.61	10,2
GLAG + SAM <sub>2</sub>	8.6	90.6	0.40	0.79	20,6

# satelit (samostalnik)

## salience

## frequency

# KOLIKO KOLOKACIJ IZLUŠČITI?

Kriterij:

- jakost = 6,00 jakost?
- Prvih 20 kolokacij?

<u>S kakšen?</u>		43.95
vohunski +	<u>357</u>	9.81
Zemljini +	<u>217</u>	9.72
geostacionaren +	<u>179</u>	9.72
Jupitrov +	<u>130</u>	9.09
meteorološki +	<u>111</u>	8.33
umeten +	<u>510</u>	8.13
komunikacijski +	<u>246</u>	7.96
Saturnov	<u>53</u>	7.84
telekomunikacijski +	<u>160</u>	7.78
Nasin	<u>50</u>	7.61
opazovalen	<u>40</u>	7.14
Marsov	<u>33</u>	7.09
izstreljen	<u>28</u>	6.90
navigacijski	<u>47</u>	6.57
Galilejev	<u>20</u>	6.54
sovjetski +	<u>123</u>	6.46
astronomski	<u>34</u>	6.37
Uranov	<u>17</u>	6.33
odslužen	<u>26</u>	6.19
infrardeč	<u>24</u>	6.10
istrski	<u>34</u>	5.88
nedelujoč	<u>13</u>	5.88
utirjen	<u>11</u>	5.74
vremenski	<u>83</u>	5.67
izvidniški	<u>13</u>	5.67

<u>S kakšen?</u>		43.95
umeten +	<u>510</u>	8.13
vohunski +	<u>357</u>	9.81
ameriški +	<u>308</u>	4.28
naraven +	<u>249</u>	5.52
komunikacijski +	<u>246</u>	7.96
Zemljini +	<u>217</u>	9.72
geostacionaren +	<u>179</u>	9.72
vojaški +	<u>171</u>	5.01
telekomunikacijski +	<u>160</u>	7.78
velik +	<u>150</u>	1.07
Jupitrov +	<u>130</u>	9.09
sovjetski +	<u>123</u>	6.46
nov +	<u>116</u>	0.60
meteorološki +	<u>111</u>	8.33
ruski	<u>97</u>	4.38
majhen	<u>88</u>	2.55
evropski	<u>84</u>	1.77
vremenski	<u>83</u>	5.67
komercialen	<u>65</u>	5.67
nekdanji	<u>61</u>	2.39
slovenski	<u>61</u>	0.09
teniški	<u>57</u>	5.07
zemeljski	<u>55</u>	5.45
raziskovalen	<u>54</u>	4.60
edin	<u>54</u>	2.97

# Distribucijska semantika

struktura	baseline	SkE SVM	vložitve SVM	SkE + vložitev SVM
PRID + SAM	0,5	0,563	0,715	<b>0,721</b>
SAM + SAM <sub>2</sub>	0,5	0,615	0,841	<b>0,843</b>
GLAG + SAM <sub>4</sub>	0,5	0,644	0,816	<b>0,822</b>
PRISL + GLAG	0,5	0,716	<b>0,879</b>	<b>0,879</b>
PRISL + PRID	0,5	0,794	0,852	<b>0,856</b>
SAM + v + SAM <sub>5</sub>	0,5	0,651	0,740	<b>0,746</b>
GLAG + PRISL	0,5	0,696	0,801	<b>0,806</b>
GLAG + v + SAM <sub>5</sub>	0,5	0,510	0,666	<b>0,668</b>
SAM + s/z + SAM <sub>6</sub>	0,5	0,533	<b>0,786</b>	0,785
GLAG + s/z + SAM <sub>6</sub>	0,5	0,638	0,660	<b>0,683</b>
GLAG + SAM <sub>2</sub>	0,5	0,862	0,795	<b>0,810</b>



# Zaključki

- Uporabnost
  - KSSS je odzivni stopenjski slovar z bogatim naborom podatkov o kolokacijah
- Odprtost
  - Kolokacijski podatki iz KSSS bodo na voljo kot baza podatkov v repozitoriju CLARIN.SI pod licenco Creative Commons 4.0 CC-BY
- Podprtost
  - Povezani projekti nudijo možnost za izboljšanje metodološkega, vsebinskega in predstavitvenega dela

# Zahvala



- Predstavitev je podprl projekt Kolokacije kot temelj jezikovnega opisa: semantični in časovni vidiki, ki ga financira Javna agencija za raziskovalno dejavnost Republike Slovenije (KOLOS, J6-8255).
- Projekt sta podprla infrastrukturna programa:
  - Center za jezikovne vire in tehnologije, Univerza v Ljubljani
  - Centre za uporabno jezikoslovje pri zavodu Trojina (I0-0051)
- Projekt je podprl raziskovalni program “Slovenski jezik – bazične in kontrastivne raziskave” (P6-0215)