

Slovene Lexical Database automatic extraction and crowdsourcing

Simon Krek

„Jožef Stefan“ Institute

Iztok Kosem

Trojina, Institute for Applied Slovene Studies

Polona Gantar

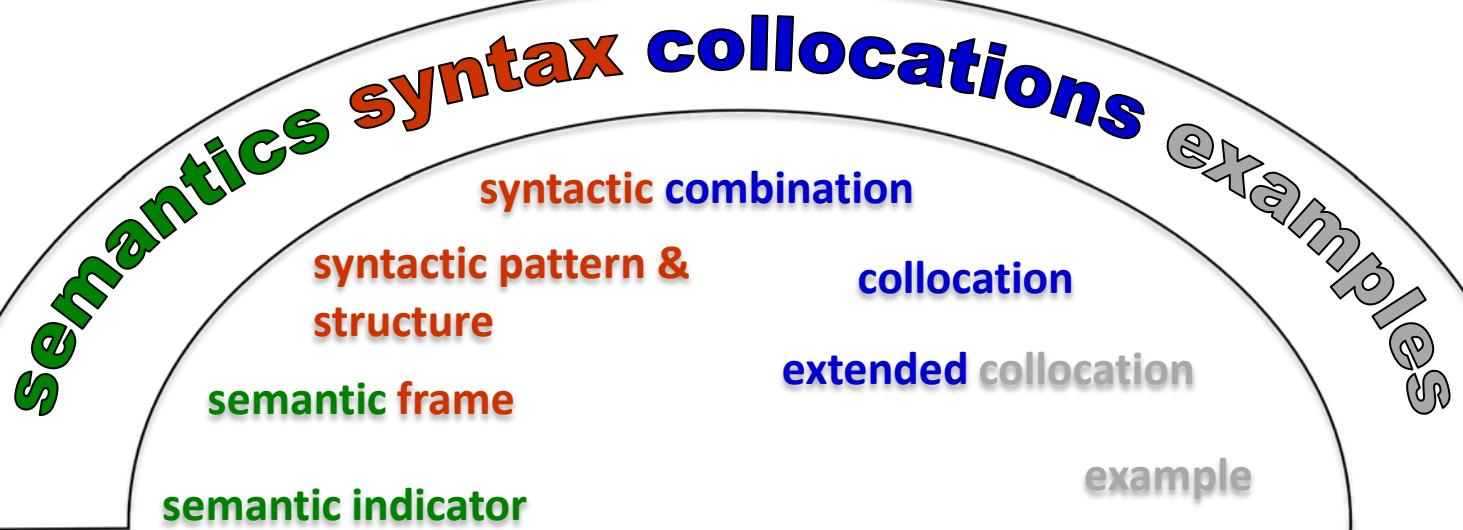
Fran Ramovš Institute of the Slovenian Language

Plan

- Slovene Lexical Database
- Extraction of data (Sketch Engine)
 - Sketch Grammar
 - GDEX (Good Dictionary EXamples)
- Workflow / crowdsourcing
- ACDC (Automatically Constructed Dictionary Content)

SLD Basics

- corpus data analysis
- lexicogrammatical approach
 - semantics and syntax are not separated
- meaning = meaning potential
 - is not stable (norms & exploitations)
- lumpers vs. splitters = splitters
- lexicography first, NLP second



I. LEMMA

- headword
- part-of-speech

svitati se (to dawn)

verb

II. SENSE

- indicator

1. daniti se (day)

2. dojemati (understand)

- semantic frame

ko se svita DAN.

začne vzhajati sonce

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

III. SYNTAX

- lable

only in 3rd pers.

- structure

gbz Inf-GBZ

- pattern

*kaj se svita
(sth is dawning)*

- synt. combin.

rbz GBZ

komu se svita o čem

(sth is dawning to sb about sth)



IV. COLLOC.

- collocation

[začeti. pričeti] se svitati

[počasi. malo. malce] se svita

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 pričenja svitati o nekdanji zvezi ned Chadom in Heather.

- multi-word unit

VI. PHRASEOLOGY

- phraseological units

I. Lexical Unit

- link to the lexicon
 - morphosyntactic information
 - corpus frequency
 - pronunciation etc.
- additional grammatical information
 - un/countability, part-of-speech subtypes etc.

II. Semantic Level

- Semantic Indicators
 - simple EFL-like explanations or synonyms forming a sense menu
 - self-explanatory in relation to each other
- Semantic Frames
 - COBUILD / FrameNet / Corpus Pattern Analysis
 - combination of the systems

Semantic Indicators – koža (skin)

koža *samostalnik*

- 1. vrhnji del telesa**
 - 1.1 pri človeku**
 - 1.2 pri živali**
- 2. odstranjen vrhnji del živalskega telesa**
- 3. ovoj ali lupina**

Semantic Frames

- identification of verb/semantic arguments
 - prototypical pattern – “the norm” (Hanks)
 - the headword in its syntactic environment
- identification of semantic types in particular syntactic positions
- the semantic scenario
 - a full-sentence definition making a link between the arguments and the situation (FN) typical for a particular sense

Semantic Frame

2. dojemati

2.1 nekaj vedeti

če se ČLOVEKU svita o nekem DEJSTVU, potem o tem
nekaj ve ali sluti

- semantic types in capital letters (ID-ed)
- linked with collocates via syntax

III. Syntactic Level

- **syntactic structures** (formal)
 - clause and phrase level (all POS; only for NLP)
 - the number of syntactic structures is finite
 - source: word sketches (Sketch Engine)
- **syntactic patterns**
 - valency (mainly verbs; for lexicography and NLP)
- **syntactic combinations**
 - more than basic patterns: „pasti za X stopinj“

Syntactic Structures – koža

4 vrhnji del telesa

1.1 pri človeku

- pbz0 **SBZ0** [občutljiva, suha, mastna] koža
- **SBZ0** sbz2 koža [obraza, telesa, rok, lasišča]
- **SBZ0** pod sbz6 koža pod [pazduho, očmi]
- gbz **SBZ4** [dražiti, pomirjati, hladiti] kožo

Syntactic Patterns – svitati se

2. dojemati

2.1 nekaj vedeti

če se ČLOVEKU svita o nekem DEJSTVU, potem o tem
nekaj ve ali sluti

- komu se svita se o čem
- komu se svita kaj

IV. Collocation Level

- **SEMANTIC FRAME:**

če se ČLOVEKU svita o nekem DEJSTVU, potem o tem nekaj ve ali slutí

- **SYNTACTIC STRUCTURES AND PATTERNS:**

NOUN – koža

pbz0 SBZ0

SBZ0 sbz2

SBZ0 pod sbz6

gbz SBZ4

VERB – svitati se

komu se svita se o čem

komu se svita kaj

If a part of syntactic patterns are collocational, they are shown on the collocation level.

- **COLLOCATIONS**

- [občutljiva, suha, mastna] koža
- koža [obraza, telesa, rok, lasišča]
- koža pod [pazduho, očmi]
- [dražiti, pomirjati, hladiti] kožo

I. Examples

• COLLOCATIONS

- [občutljiva, suha, mastna] **koža**
- **koža** [obraza, telesa, rok, lasišča]
- **koža** pod [pazduho, očmi]
- [dražiti, pomirjati, hladiti] **kožo**

• EXAMPLES

- *Tonik je namenjen je **občutljivi koži** in ne vsebuje alkohola.*
- *Koža rok postane pozimi občutljivejša.*
- *Opažate na **koži pod očmi** prezgodnja znamenja staranja?*
- *Se vam že kaj **svita**, o čem govorim?*
- *Petru pa se pričenja **svitati** o nekdanji zvezi med Chandlerjem in Heather.*
- *Holly je na svojem stolu v klubu Diva zastokala in se prijela za glavo, ko se ji je začelo **svitati**, kaj se bo zgodilo.*

Sketch Engine (word sketch)

user: Simon Krek corpus: Fida PLUS 620m (SLD sketch grammar)

[Concordance](#)
[Word List](#)
[Word Sketch](#)
[Thesaurus](#)
[Sketch-Diff](#)
[? Help menu](#)

[Save](#)
[Change options](#)
[Turn on clustering](#)
[More data](#)
[Less data](#)
[Switch menu position](#)

ljubiti (gлагол) Fida PLUS 620m (SLD sketch grammar) freq = 39356

kako-kdaj?	7198	5.3
<input checked="" type="checkbox"/> strastno	134	64.24
<input checked="" type="checkbox"/> neizmerno	150	63.12
<input checked="" type="checkbox"/> neskončno	93	51.73
<input type="checkbox"/> nesmrtno	36	51.49
<input type="checkbox"/> znova	290	49.92
<input type="checkbox"/> vedno	562	47.95
<input type="checkbox"/> nadvse	145	47.83
<input type="checkbox"/> brezpogojno	52	47.35
<input type="checkbox"/> iskreno	83	45.13
<input type="checkbox"/> noro	40	42.6
<input type="checkbox"/> ponovno	157	41.24
<input type="checkbox"/> srčno	41	40.47
<input type="checkbox"/> zelo	431	40.35
<input type="checkbox"/> tako	566	38.53
<input type="checkbox"/> resnično	58	37.63
<input type="checkbox"/> preprosto	95	36.85
<input type="checkbox"/> brezmejno	18	36.81
<input type="checkbox"/> bolj	309	36.52

veznik	1903	1.6
<input type="checkbox"/> in	221	43.18
<input type="checkbox"/> a	131	38.35
<input type="checkbox"/> vendar	130	37.95
<input type="checkbox"/> čeprav	70	34.25
<input type="checkbox"/> ker	107	33.38
<input type="checkbox"/> toda	56	32.63
<input type="checkbox"/> če	101	31.75
<input type="checkbox"/> zato	73	30.73
<input type="checkbox"/> kakor	34	30.44
<input type="checkbox"/> da	299	28.95
<input type="checkbox"/> ali	54	28.58
<input type="checkbox"/> ampak	45	27.18
<input type="checkbox"/> ko	81	26.6
<input type="checkbox"/> kar	50	23.55
<input type="checkbox"/> kadar	16	22.89
<input type="checkbox"/> dokler	17	22.12
<input type="checkbox"/> kajti	20	21.81
<input type="checkbox"/> saj	50	20.98

predlog	134	0.0
<input type="checkbox"/> kot	46	26.26
<input type="checkbox"/> od	23	18.13

[>>](#)

z-d	640	1.7
<input type="checkbox"/> srce	87	45.79
<input type="checkbox"/> ljubezen	42	34.13
<input type="checkbox"/> ženska	49	31.67
<input type="checkbox"/> maček	16	31.55
<input type="checkbox"/> žena	19	24.53
<input type="checkbox"/> moški	21	23.41
<input type="checkbox"/> moč	15	20.0

[>>](#)

predl-za	1874	0.5
<input type="checkbox"/> z	756	32.52
<input type="checkbox"/> kot	122	21.28
<input type="checkbox"/> zaradi	71	20.56
<input type="checkbox"/> brez	39	17.48
<input type="checkbox"/> ob	55	15.15
<input type="checkbox"/> do	70	15.08
<input type="checkbox"/> iz	74	14.66
<input type="checkbox"/> med	50	13.82
<input type="checkbox"/> v	268	13.65
<input type="checkbox"/> na	170	13.0
<input type="checkbox"/> od	39	9.82
<input type="checkbox"/> po	35	6.78
<input type="checkbox"/> pri	23	5.95

[>>](#)

po-d	75	0.6
<input type="checkbox"/> smrt	18	32.34

[>>](#)

kot-d	135	3.1
<input type="checkbox"/> Slovenec	17	29.5

[>>](#)

v_dajal	313	3.3
<input type="checkbox"/> oseba	75	42.67

[>>](#)

y_rodil	1624	10.4
<input type="checkbox"/> nada	238	73.21
<input type="checkbox"/> oseba	158	43.35
<input type="checkbox"/> škarje	21	34.79
<input type="checkbox"/> ženska	60	28.96
<input type="checkbox"/> Bog	22	27.09
<input type="checkbox"/> opera	16	24.28
<input type="checkbox"/> nit	29	24.19
<input type="checkbox"/> žena	24	23.12
<input type="checkbox"/> žival	26	22.23
<input type="checkbox"/> težava	27	17.8
<input type="checkbox"/> stvar	18	16.94

[>>](#)

Good dictionary examples (GDEX)

user: Simon Krek corpus: Fida PLUS 620m (SLD sketch grammar)

[Concordance](#)
[Word List](#)
[Word Sketch](#)
[Thesaurus](#)
[Sketch-Diff](#)
[?Help on main menu](#)

[Sketch menu position](#)

Tickbox Lexicography - Select Examples

Lemma: **ljubiti**
Gramrel: kako-kdaj?

Template: fidaplus_slovene Alternative GDEX configuration: None

strastno

Ali pa v hudi moški nuji svojo ljubico strastno **ljubi** kar zadaj, za težko žametno zaveso.
 Skupaj sva kar precej popila, nakar sva se strastno **ljubila** dve uri.
 Suh in dolgolas je bil intelektualec, strastno je **ljubil** športne avtomobile, a ni imel vozniškega izpita.

neizmerno

Bil je pravi Italijan in je poleg dobrih avtov in cigar neizmerno **ljubil** ženske.
 Bil je zares iskren in dober prijatelj, ki je neizmerno **ljubil** in varoval naravo.
 Kako bi mogel drugače, saj je Evo neizmerno **ljubil**.

neskončno

Morje zares neskončno **ljubim**, tam sem » bebavo blažena ».
 Posnemati Boga, ki neskončno **ljubi** končne stvari kot končne.
 Solarna oseba, veselega značaja, ki neskončno **ljubi** življenje in barve!

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

gramrels

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

III. SYNTAX

word
sketches

- label
- structure
- pattern
- synt. combin.

only in 3rd pers.

gbz Inf-GBZ

kaj se svita
(sth is dawning)

rbz GBZ

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

• collocation

[začeti. pričeti] se svitati

[počasi. malo. malce] se svita

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 pričenja 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	(imetи) 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
- new „directives“
 - *SEPARATEPAGE
 - *CONSTRUCTION
 - *COLLOC

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: *SEPARATEPAGE

- 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_e2	206	17.9
iz-d_e2	107	5.0
do-d_e2	22	1.6
brez-d_e2	21	3.7
v-d_e2	21	12.0
poleg-d_e2	21	8.9
zaradi-d_e2	18	2.0
z-d_e2	15	2.1
za-d_e2	14	4.4

Example: *SEPARATEPAGE

dobiti (glagol) FidaPlus (20M) freq = 11162 (735.5 per million)					
displaying only: koga-česa_g2 whole word sketch					
<u>brez-d_g2</u>	<u>21</u>	3.7	<u>iz-d_g2</u>		
laganje	<u>1</u>	9.61	Nigerija	<u>2</u>	8.82
recept	<u>6</u>	8.34	onostranstvo	<u>1</u>	8.21
odredba	<u>1</u>	6.69	arest	<u>1</u>	8.14
natečaj	<u>1</u>	6.06	katran	<u>1</u>	8.13
težava	<u>9</u>	4.87	Ligojna	<u>1</u>	8.13
razpis	<u>1</u>	3.95	Juršinci	<u>1</u>	8.07
problem	<u>1</u>	3.02	sukcesija	<u>1</u>	8.03
razlog	<u>1</u>	2.89	limfa	<u>1</u>	8.01
			vrečica	<u>1</u>	7.92
			ZPIZ	<u>1</u>	7.91
			<u>proračun</u>	<u>20</u>	7.74
			NT	<u>1</u>	7.71
			Montreal	<u>1</u>	7.71
			Kremelj	<u>1</u>	7.63
			mozeg	<u>1</u>	7.62
			Carigrad	<u>1</u>	7.6
			<u>na-d_g2</u>	<u>5</u>	3.1
			priložnost	<u>1</u>	3.18
			stran	<u>3</u>	3.12
			sredstvo	<u>1</u>	2.41
			<u>o-d_g2</u>	<u>1</u>	8.0
			informacija	<u>1</u>	2.91
			<u>zaradi-d_g2</u>	<u>18</u>	2.0
			kolegialnost	<u>1</u>	10.68
			črevesje	<u>1</u>	8.06
			panika	<u>1</u>	7.32
			taktika	<u>1</u>	7.17
			obnašanje	<u>2</u>	7.15
			noša	<u>1</u>	7.08
			prostornina	<u>1</u>	6.63
			pnevmatika	<u>1</u>	6.32
			<u>za-d_g2</u>	<u>14</u>	4.4
			spletka	<u>1</u>	8.64
			pušt	<u>1</u>	8.33
			karta	<u>1</u>	5.09
			mleko	<u>1</u>	4.25
			denar	<u>4</u>	3.8
			naloga	<u>1</u>	3.56
			stanovanje	<u>1</u>	3.29
			vloga	<u>1</u>	2.64
			pravica	<u>1</u>	2.12
			cesta	<u>1</u>	1.92
			<u>od-d_g2</u>	<u>206</u>	17.9
			dalajlama	<u>2</u>	8.05
			gradbenik	<u>2</u>	7.93
			prednik	<u>2</u>	7.43
			Avstrijec	<u>2</u>	7.34
			Pelhan	<u>1</u>	7.3
			Adanič	<u>1</u>	7.29
			Izabela	<u>1</u>	7.21
			deklič	<u>1</u>	7.17
			Lek	<u>2</u>	7.15
			FIBA	<u>1</u>	7.13
			NEK	<u>1</u>	7.12
			Uefa	<u>1</u>	7.11
			bršljan	<u>1</u>	7.09
			Portugalska	<u>1</u>	6.98
			vol	<u>1</u>	6.93
			Nik	<u>1</u>	6.93

*CONSTRUCTION

- 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: *CONSTRUCTION

=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

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 Vrhni. 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). </p><p> Sveti **Martin** je **dal** svoje ime cerkvi s prepoznavnim baročnim privoščiš ſe 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) </p><p> Kongresno **testiranje** varnosti **dalo** porazno sliko </p><p> Kot švicarski sir </p> 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 odnosu do (koga/česa)" (in relation to (sb/sth))
- Mainly nominal headwords
- Under (sub)sense after syntactic structures as a separate category

Example: *COLLOC

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

preposition

noun

preposition

Example: "in relation to"

corpus: FidaPlus (20M)

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

GDEX – Good Dictionary Examples

- system for evaluation (ranking) of sentences with respect to their suitability to serve as dictionary examples
- sorting sentences so that good examples do not have to be searched for in hundreds of unusable sentences
- initially trained on English, but it did not give good results for other languages

GDEX – configuration

- parameters in a GDEX configuration file
- GDEX Tools web-interface to create and use custom GDEX configurations
- the GDEX evaluation process
 - ranking of out-of-corpus sentences
 - evaluation of TBLex logs
 - cooperation with WEKA

GDEX classifiers

- procedures that quantify measurable features of sentences or tokens
- sentence classifiers: sentence length, keyword position, etc.
- token classifiers: token frequencies, matches to RE, etc.

Evaluation of TBLex logs

user: Simon Krek corpus: Fida PLUS 620m (SLD sketch grammar)

[Concordance](#)
[Word List](#)
[Word Sketch](#)
[Thesaurus](#)
[Sketch-Diff](#)
[?Help on menu menu](#)

[Switch menu position](#)

Tickbox Lexicography - Select Examples

Lemma: **ljubiti**
Gramrel: **kako-kdaj?**
Template: **fidasplus_slovene** Alternative GDEX configuration: **default**

GDEX: Slovene3 **GDEX: default configuration**

strastno

Ali pa v hudi moški nuj svojo ljubico strastno **ljubi** kar zadaj, za težko žametno zaveso.

Skupaj sva kar precej popila, nakar sva se strastno **ljubila** dve uri.

Suh in dolgolas je bil intelektualec, strastno je **ljubil** športne avtomobile, a ni imel vozniskega izpita.

strastno

Sophie me zdaj strastno **ljubi** .

Nate, ki te strastno **ljubim** .

Vse dneve je pridno delal in jo vse noči strastno **ljubil** .

neizmerno

Bil je pravi Italijan in je poleg dobrih avtov in cigar neizmerno **ljubil** ženske.

Bil je zares iskren in dober prijatelj, ki je neizmerno **ljubil** in varoval naravo.

Kako bi mogel drugače, saj je Evo neizmerno **ljubi** .

neizmerno

Dekle, ki neizmerno **ljubi** .

Sam je Marijo neizmerno **ljubil** .

Bog nas neizmerno **ljubi** .

neskončno

Morje zares neskončno **ljubim** , tam sem » bebavo blažena ».

Posnemati Boga, ki neskončno **ljubi** končne stvari kot končne.

Solarna oseba, veselega značaja, ki neskončno **ljubi** življenje in barve!

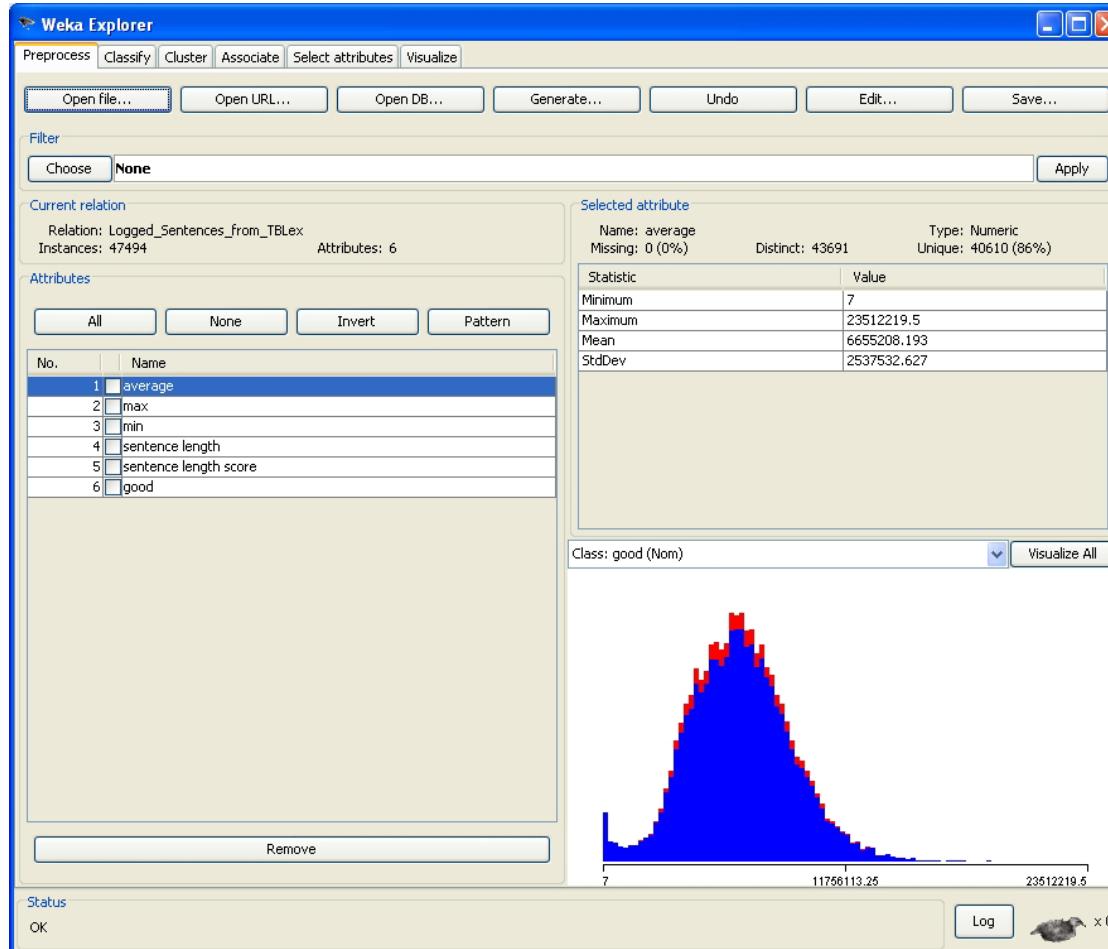
neskončno

Tudi ona ga neskončno **ljubi** .

Neskončno te **ljubim** .

Zavedel se je, kako neskončno jo **ljubi** .

Cooperation with WEKA



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>

```
<skladenjska struktura>
```

```
<struktura>kakšen?</struktura>
```

```
<kolokacije>
```

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

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

```
    <kolokacija id="840017"><k>spletен</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>. </zgled>
```

ADJECTIVE + NOUN

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 neuvršcene države <i>dale</i> vso
prednost jedrski razorožitvi.</zgled>
</zgledi>
```

Gramrel to <oznaka>

<oblika>

<iztocnica>**mesto**</iztocnica>

</oblika>

 unary to label: "with proper names"

<zaglavje>

<besvrs>samostalno</besvrs>

<oznaka>z_lastnim_imenom</oznaka>

</zaglavje>

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 (each PoS)
 - 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)

Lemmalist

- -I LEMMALIST, --lemmalist=LEMMALIST
 - The file containing a list of lemmas 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>

automatic data extraction + visualisation

sense division, definitions, compounds and phraseology

editing

computer

crowd-sourcing

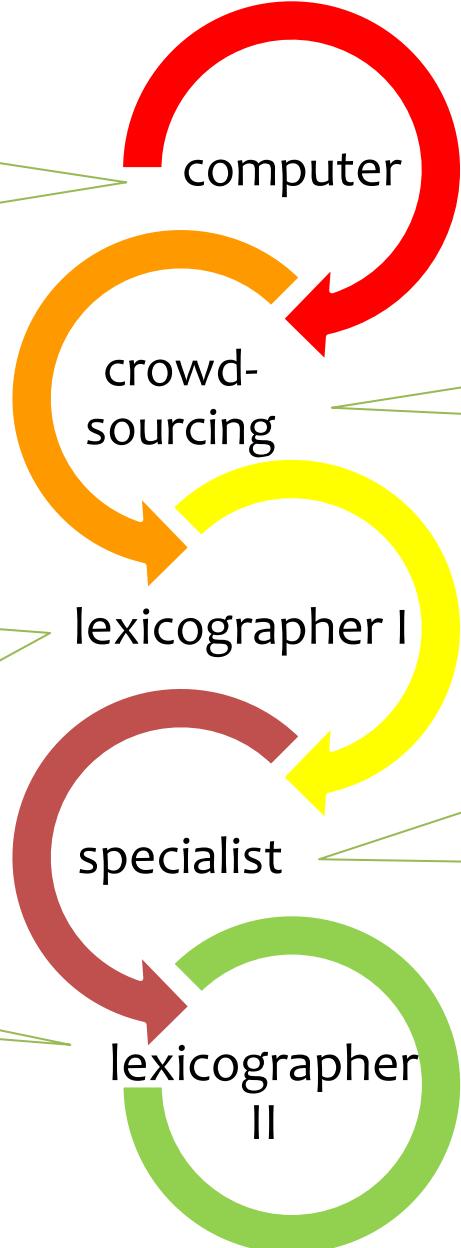
lexicographer I

specialist

lexicographer II

data clean-up and sorting

Terminology, pronunciation, tonality etymology



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 ustrezano 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 ozioroma poslanstvo, ki vam lahko ustreza ali pa ne.

DA

NE

Ne vem

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.)

Lexicographer I.

clanek SSSJ-baza (Entry Document)

glava

oblika

zapis:**bančnik**
iztocnica:**bančnik**

zaglavje

besvrs:**samostalnik**

geslo

1.pomen^I

indikator:

pomenska shema:

skladenjske skupine □

a) struktura: pbz0 SBZ0

kolokacije □

- [sivolas]
- [privaten]
- [centralen]
- [investicijski]
- [upokojen]
- [ugleden]
- [dolgočasen]
- [premožen]
- [vodilen]
- [izkušen]
- [švicarski]
- [vpliven]
- [oseben]
- [dolgoleten]
- [bivši]
- [obsojen]

zgledi □

- Za Irca Patricka Hickeyja, elegantnega svolasega **bančnika** iz Dublina lahko rečemo, da je eden izmed tistih, katerega olimpijska specifična teža v zadnjem obdobju vidno narašča.
- JE ČEZ MNOGO ČASA REKEL SIVOLASI **BANČNIK**
- Dolgoletni svolasi **bančnik** je odklonil ponujeno funkcijo in je svojo odločitev obrazložil z moralno-političnimi argumenti, vlada, ki je takšno vodstvo banke skuhal (in nadzorni svet, v katerem sta dva politična funkcionarja), pa se je znašla v položaju, da bo morala pod drobnogledom javnega smehjanja srebat skuhan.
- Odnos do strank in Kultura poslovanja naslohi sta v privatni banki popolnoma drugačna, poudarjajo tudi drugi privatni **bančniki**, s katerimi smo se pogovarjali.
- Vrednote, ki vodijo privatnega **bančnika**, so strokovna odličnost, diskretnost in zaupnost.
- Ponudbo in svetovanje izvajajo izobraženi finančni svetovalci s posebno licenco, ki jih imenujemo privatni **bančniki** in so strankam na voljo kadar koli in kjer koli.

ACDC

[Explain](#)[Combine](#)[Exemplify](#)[Soundify](#)[Streamify](#)[Visualize](#)[Translate](#)**baboon** noun

plural: baboons

Definitions**Found**

Baboons are terrestrial monkeys found in open or rocky areas, including open woodland, savannah, grassland, and rocky hills in Africa.

International Primate Protection League (IPPL)

Context – synt. structures + ex.

Wikipedia

Baboons - large terrestrial monkeys having doglike muzzles

WordNet

The baboons are some of the largest non-hominid members of the

Context – collocations + ex.**Generated**

Baboons are animals.

Baboons live in Africa.

Baboons are boisterous and cunning.

Baboons live up to 45 years.

Multi-word expressions (Parseme?)**Definitions found – def extraction****Collocations**

- ▶ as subject
baboons [bark]
If baboons are barking and impala are snorting, is it lions or a leopard on the prowl?
- ▶ as object
baboons are [chased, spotted]

One lioness started chasing a pregnant female baboon amongst almost 100 baboons troop.

A lone baboon is spotted sitting on snow-covered rocks looking about.

noun+of+baboon

[troop, tribe, colony] of baboons
A troop of baboons tried to drag both mother and cub from their den.
One tribe of baboons ate tainted meat, and it ended up killing most of those at the higher levels.
These San Antonio based colonies are populated with several unique animal models of human disease as well as a large, unparalleled, colony of baboons

Generated definitions – NL generation**• chacma baboon**

Swamp Troop is a tale of tragedy and triumph, friendship and fatherhood in a troop of chacma baboons that inhabits Botswana's legendary Okavango Delta.

• gelada baboon

some female primates, such as the gelada baboon , advertise their sexual status via sexual swellings of the chest nodules which flush red

• hamadryas baboon

Captive hamadryas baboon can similarly learn to use tools through trial and error.

• dog-faced baboon

Beside us, on the surrounding rocks and forming a complete semi-circle, sat great dog-faced baboons with their wives and children

• mandrill baboon

His drawing of a mandrill baboon - wrongly identified as a 'hyena' in the original - appears in successive books including those of William Topsell (1607) and Johann Johnston (1653).

Explain

Combine

Exemplify

Soundify

Streamify

Visualize

Translate

baboon

noun

plural: baboons

Sounds, Graphics and Visuals

Sounds

Recorded
baboon ►

Speech Synthesis
baboon ►

Graphics



Visuals

Images



<http://www.image-net.org/> - Synset: baboon



<http://www.image-net.org/> - Synset: chacma, chacma baboon

Videos



www.youtube.com

[Explain](#)[Combine](#)[Exemplify](#)[Soundify](#)[Streamify](#)[Visualize](#)[Translate](#)

baboon

noun

plural: baboons

Streaming

Twitter

Sam Kanyi @snkanyi Its a shame how my standards have dipped. I have to settle for those whose spelling is as bad as IR baboon's! Smh

37m Lesley Marshall @lesleyrocksface Socotra Island
Blue Baboon... @devonmassyn is the best boyfriend ever! I can't believe I got her for Christm

1h ____^ XinYee ♥ @dream180793 thanks for talking to me hahaha thanks baboon ♥

1h kyle† @KylieJarvis people are arguing on here on about who got the best presents.. Be grateful you baboon!

2h Yaya Darwish ∞ hahahah you know me too well. I literally can't wait, counting down the seconds & milliseconds till you're here baboon! ♥

2h Maximus aurelius @don_jide How can this baboon next to me say the queen would get the D i am officially dead

News Feeds

BABOONS CAN RECOGNIZE WORDS The findings might weigh in on debates about how best to teach children to read...
<http://news.discovery.com/animals/baboons-words-120412.html>

Penn Researchers Connect Baboon Personalities to Social Success and Health Benefits...
<http://www.upenn.edu/pennnews/news/penn-researchers-connect-baboon-personalities-social-success-and-health-benefits>

IT PAYS TO BE A NICE BABOON Nice baboons enjoy stronger social circles and...
<http://news.discovery.com/animals/baboons-nice-121002.html>

Foraging Baboons Are Picky Punters: Baboon Foraging Choices Depend On Their Habitat and Social Status
<http://www.sciencedaily.com/releases/2012/09/120913203917.htm>

Showdown over baboons
Cape Town - Baboon conservationists, scientists, animal welfare activists and wildlife managers hold a no-holds-barred meeting on Monday...
<http://www.iol.co.za/news/south-africa/western-cape/showdown-over-baboons-1.1421348#.UNnwFOQmbh4>

Word-detecting baboons are a tough read
Studies differ on whether monkeys learn to recognize letter combinations
http://www.sciencenews.org/view/generic/id/346591/description/Word-detecting_baboons_are_a_tough_read

[Explain](#)[Combine](#)[Exemplify](#)[Soundify](#)[Streamify](#)[Visualize](#)[Translate](#)

baboon

noun

plural: baboons

Translate

Language

Basque
Czech
Estonian
Slovenian

Translation(s)

pavián
paavian
pavijan

Czech

Tak já ti něco řeknu , ty mechanická opice .

Představa lady Nicholsonové jako jeho choti , skutečnost , že její jemná krása byla někdy vydána zvrhlostenem tohoto degenerovaného paviána , naplnila Doylea niternou zuřivostí .

Dále Slater cituje příklady , kdy se vampýři dělí o krev , kterou obstaral jen jeden z nich , a později si zas vyměňují úlohy , nebo kdy subdominantní paviáni samec vyvádí takové kousky , že odláká pozornost hlavního samce v tlapě a další samec se mezitím může pářit z jednou ze samic . Poté dodává :

Medvědí 'oplátkový systém ' neprobíhá v podmínkách , kdy jeden oplácí druhému , jako tomu bylo u Slaterových příkladů s vampýry a paviány ;

Jejich zjevná radost nebyla tak úplně nepodobná známému štěstí návštěvníků zoologické zahrady , kterým se podaří zastihnout ochotnatého právě v době krmení .

English

Just let me tell you something , my brass **baboon** . At very high doses, a poor gastric tolerability (gastritis, gastric erosions and/or vomiting) of clopidogrel was also reported in rat and **baboon**.

The idea of Lady Nicholson as his spouse , that her handsome refinement had ever been subject to the vicissitudes of this **baboon** 's degeneracy , filled Doyle with moral outrage .

During non clinical studies in rat and **baboon**, the most frequently observed effects were liver changes.

Citing examples of vampire bats sharing blood that only one has obtained , then reversing roles later , and of sub-dominant male **baboons** taking turns distracting the alpha male so the other can mate with one of its females , Slater goes on to say :

At very high doses, a poor gastric tolerability (gastritis, gastric erosions and/or vomiting) of clopidogrel was also reported in rat and **baboon**.

A bear 's reciprocity is n't one-on-one like Slater 's examples of bats and **baboons** ;

During non clinical studies in rat and **baboon**, the most frequently observed effects were liver changes.

Their obvious delight was not entirely unlike the often seen elation of visitors to a zoo on coming upon the crested **baboon** at feeding time .

Slovenian

Pri zelo visokih odmerkih so pri podganah in **pavijanih** poročali tudi o slabem želodčnem prenašanju (gastritis, želodčne erozije in/ ali bruhanje).

V predkliničnih študijah na podganah in **pavijanih**, so bile jetrne spremembe najpogosteje opaženi učinek.

Pri zelo visokih odmerkih so pri podganah in **pavijanih** poročali tudi o slabem želodčnem prenašanju (gastritis, želodčne erozije in/ ali bruhanje).

V predkliničnih študijah na podganah in **pavijanih**, so bile jetrne spremembe najpogosteje opaženi učinek.

Plan

- Slovene Lexical Database
- Extraction of data (Sketch Engine)
 - Sketch Grammar
 - GDEX (Good Dictionary EXamples)
- Workflow / crowdsourcing
- ACDC (Automatically Constructed Dictionary Content)