

Dictionary Writing Systems and Corpus Query Systems

Session 1, part 1

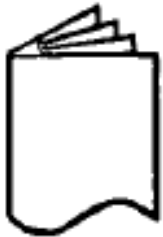
Laurence Urdang

- Dictionaries No. 6 (1984)
- A Lexicographer's Adventures in Computing
 - „I clearly recall my first brush with computers: in 1959, I was working on a new dictionary for Random House, and it occurred to me that a computer would be ideal for the sorting and manipulation of the various kinds of data I was responsible for generating.“

1966

- "As far as I know, the book we produced, *The Random House Dictionary of the English Language – The Unabridged Edition*, was the first to use computers extensively in its editorial preparation."
- "With all the data appropriately coded, programs enabled the computer to sort all of the bits and pieces into dictionary order."
- "In short, we created a powerful database, consisting of a dictionary with about 260,000 entries which we could access at will."

Euralex bulletin 1990



COMPULEXIS DICTIONARY SYSTEM

The only data-base system specially designed for lexicographers and dictionary publishers.

- Combines data-capture with editorial tools and typesetting.
- Developed by a team of lexicographers and computer specialists, and *individually configured to suit the requirements of each dictionary project.*
- Quick to learn and easy to use.
- Already successfully applied by Oxford University Press and other major European dictionary publishers.
- Runs on IBM AT or PS/2 or compatibles.

Specialized dictionary software

With the COMPULEXIS dictionary system ...

- The editor is free to concentrate on the "real" work rather than checking punctuation, typography, alphabetical order, etc.
- Proofs which look like the final product can be produced at any time, without external cost.
- Selective check-lists with cross-references and other problem areas can be printed.
- Structural punctuation, frequently used text elements, and numbers and letters used for subdivision can be generated automatically.
- Consistency is dramatically improved.
- Data can be re-used in new editions and spin-off products.

Oxford University Press

What the users think ...

"... We soon realized what fanatastic help we got from the COMPULEXIS system.

... For the editors the COMPULEXIS system has opened the door to all sorts of possibilities.

... The most conservative department in the house has in less than a year completely changed its working methods.

... The more we have worked with the system, the more our enthusiasm has grown."

*(Quotes from an unpublished report by
the Dictionary Department of a major European publisher)*

Gestorlex 1991

- Nov/Dec 1991 issue of Language Industry Monitor
 - official launching of GestorLEX at the October Buchmesse in Frankfurt
 - developed by Danish software house TEXTware
 - hybrid product, a strong database engine providing the platform for a structured editing environment based on SGML
 - it runs under OS/2

SGML & floppy

- facilities to allow more than one person to work on a dictionary at a time
- can select and export a selection of entries on to a floppy disk for further editing on another machine
- selection is “locked” out in the main database although it can be viewed
- GESTORLEX costs DKR 15,000 (c.\$2000)

Oxford-DZS 1994

- Oxford-Hachette French Dictionary
 - Oxford University Press
 - 1994, First Edition

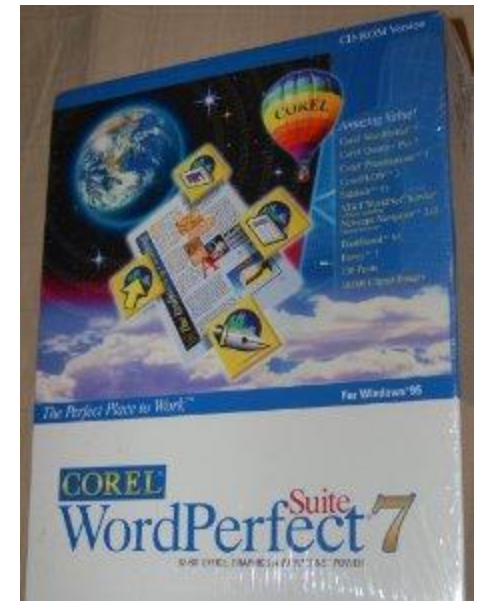
```
<se>
  <hw>abacus</hw>
  <pr><ph>"&b@k@s</ph></pr>
  <hg><ps>n</ps> (<gr>pl</gr>
  <f>-cuses</f>) </hg>
  <s2 num=1>
    <ann><la>Math</la></ann> (<ic>counting
  frame</ic>)
  </s2>;
  <s2 num=2><la>Archit</la> </s2>.
</se>
```

Gestorlex / WP7-8 (1995-)

- Problems – data:
 - database not in SGML
- Problems – Gestorlex:
 - expensive
 - lexicographers without OS/2
 - OS/2 demanding on hardware
- Solution – WordPerfect 7 (later 8, 11, X1)
 - worked on Windows (95)
 - affordable (≈60 people working at home – ~~LAN~~)
 - handled SGML format (!) with different layouts

WordPerfect

- general-purpose word processor
 - no management (!) – workflow, distribution etc.
- at the end of the project (custom tools)
 - consistency check (closed sets etc.)
 - proof-reading
 - punctuation
- during project
 - valid SGML (semi-restrictive DTD)
 - macros



Better solutions!

- Publishers home-grown systems: no
- Dedicated DWS candidates:
 - Digital Publishing System (DPS)
 - IDM (Ingénierie Diffusion Multimedia), Paris
 - iLex
 - Erlandsen Media Publishing AB, Malmö
 - TLex Suite (aka TshwaneLex)
 - TshwaneDJe Human Language Technology, Stellenbosch

Summary – DPS

- „perfect system“
- relatively complex to manage
- expensive
- for big projects (many collaborators)
- for big producers (many projects)
- XML less prominent (for lexicographers)
- very powerful management system

Summary – iLex

- „handy system“
- advanced computer literate manager
- not cheap but managable (start at ≈ 400 EUR)
- for collaborative projects
- also for solo projects
- „cafeteria style“
- explicitly XML-based and proud of it

Summary – TLex

- „one-man band system“
- designed for computer literate lexicographers
- affordable (start at 150 EUR)
- mainly for solo projects
- less for collaboration (?)
- „suite style“, "all-in one"
 - dictionary + corpora + terminology etc.
- database with XML output

Other

- Andrea Abel in the coming: *Electronic Lexicography*. Oxford University Press
 - off-the-shelf
 - ABBYY Lingvo Content
 - home-grown
 - DicSy, Norstedts
 - Wissensnetz deutsche Sprache, Duden
 - EELex, Institute for Estonian Language
 - ANW, Institute for Dutch Lexicology
 - Onoma, Woordeboek van die Afrikaanse Taal
 - Dictionary Editor and Browser (DEB), Masaryk University
 - Termania, Amebis (Slovenia)



User : polonca.kocjancic

Project SLOLEKS_A
Headword pasti
Senses 5
Words 0 Characters 0

FILE EDIT CONFIG TOOLS HELP



XREF CHECKING

NAVIGATION PANEL

Not in batch

Select all Lose edits

pasti nov polona.gantar
pasti
spomin polona.gantar

SEARCH

Download Select all

boleti 2.1 outstore ES
boleti 2.2 outstore ES
boleti 2.3 outstore ES
boleti 2.4 outstore ES
čas mojca.sorli ES
čiščenje katja.grabnar ES
deček outstore ES
deklica outstore ES
domá katja.grabnar ES
glasben olga.pobirk ES
igrati polonca.kocjancic ES
iskanje outstore ES
javnost outstore ES
kisel mojca.sorli ES
kònj outstore ES
kós outstore ES
krompir outstore ES
leteti mojca.sorli ES
mešati_Olga_konk outstore ES
mešati_Olga_konk_NEW outstore ES
mešati_Olga_SkE_NEW outstore ES
mésto katja.grabnar ES
mir outstore ES
nepriлагоjèn outstore ES
otròk katja.grabnar ES
ozek olga.pobirk ES

pasti *glagol*

1 premikati se od zgoraj navzdol navadno v 3. osebio predmetu

PREDMET zaradi sile težnosti pade KAM

a) Struktura:

- ▶ kaj pade kam
- ▶ kaj pade na kaj
- ▶ kaj pade v kaj

- [bomba, granata, drevo] pade
- pasti na [dno, Zemljo, zemljo, tla]
- pasti v [prepad, globino, morje]

- Se še spominjate dne, ko je **padla** bomba čisto blizu vaše gostilne.
- Svetleča granata je **padla** kakih petdeset metrov stran.
- Pri odcepu za Kostanjevico je **padlo** na dolenjko drevo.
- Mnoge živali in rastline poginejo in padejo na vodno dno.
- Vsako leto pade na Zemljo iz vesolja obilica kamenin.
- Debele kostanjeve ježice so padle na zemljo.
- V stanju breztežnosti nič ne pade na tla.
- Dve gondoli sta padli v globok prepad.
- Propeler je razsekalo in helikopter je kot kamen padel v globino.
- Po izjavah posadke naj bi bomba padla v morje daleč proč od letalonosilke
- Kaj pa, če helikopter pade dol?

skladenjske zveze:

pasti pod kotom

pasti pod kotom [x] stopinj

- Kovinski meteorit s premerom 40 metrov, ki pade pod kotom več kot 30 stopinj, bi tla še dosegel, preden bi se povsem stalil oziroma izhlapel.
- Pred 50 milijoni let je na Luno padel pod blagim kotom težak asteroid.

1.1 ne biti več kje nameščen, pritrjen; izpasti navadno v 3. osebi

PREDMET pade od KOD KAM

CLANEK

- GLAVA
 - OBLIKA
 - ZAPIS pasti
 - IZTOCNICA pasti
 - ZAGLAVJE
 - BESVRS glagol
- GESLO
 - POMEN
 - INDIKATOR premikati se od zgoraj navzdol
 - OZNAKA navadno v 3. osebi
 - OZNAKA o predmetu
 - POMENSKA_SHEMA PREDMET zaradi sile težnosti pade KAM
 - SKLADENJSKE_SKUPINE
 - SKLADENJSKA_STRUKTURA
 - STRUKTURA
 - VZOREC kaj pade kam
 - VZOREC kaj pade na kaj
 - VZOREC kaj pade v kaj
 - KOLOKACIJA K bomba, granata, drevo /K pade
 - KOLOKACIJA pasti na K dno, Zemljo, zemljo, tla /K
 - KOLOKACIJA pasti v K prepad, globino, morje /K
 - ZGLEDI
 - ZGLED Se še spominjate dne, ko je I padla /i bomba čisto blizu vaše gostilne.
 - ZGLED Svetleča granata je I padla /i kakih petdeset metrov stran.
 - ZGLED Pri odcepu za Kostanjevico je I padlo /i na dolenjko drevo.
 - ZGLED Mnoge živali in rastline poginejo in padejo na vodno dno.
 - ZGLED Vsako leto pade na Zemljo iz vesolja obilica kamenin.

ATTRIBUTES

Element: CLANEK

Name Value

ANNOTATIONS

EDIT

No comments for this node



Attributes (F1) Attributes (F2) Search (F3) Format (F4) Filter (F5) Corpus (F6)

clanek: <input type="checkbox"/> Incomplete	
LemmaSign	squeeze
Pronunciation	
Deriv	
Etymology	
Notes	
Frequency	0
glava:	
obluka:	
zapis:	
pc [PCDATA]	squeeze
iztocnica:	
pc [PCDATA]	squeeze
zaglavje:	
besvrs:	
pc [PCDATA]	verb

[A][B][C][D][E][F][G][H][I][J][K][L][M][N][O][P][Q][R][S][T][U][V][W][X][Y][Z] (Default styles)

squeeze verb

1 hold firmly

Structure: with adverb

■ squeeze [gently, tightly, firmly]

- I walk round behind his chair, place my hands on his shoulders, and squeeze them gently.
- I took a firm grip and squeezed tightly, but nothing happened

A with your hand

a PERSON squeezes an OBJECT with his/her hand

Structure: intransitive

► sb squeezes

- She was walking past with her granny friend, and as she was passing, a hand cupped my left butt cheek and squeezed.
- I took a firm grip and squeezed tightly, but nothing happened.

Structure: transitive

► sb squeezes sth

- She smiled as he squeezed her hand.
- Brent reaches out to squeeze Justice 's thigh with his left hand, but as the the camera picks up the actual squeeze, its Brent 's right hand squeezing her thigh.

B with your arms LABEL: AFFECTION?

a PERSON squeezes a PERSON with his/her arms

Structure: transitive

► sb squeezes sth/sb

- 'I'll come with you,' he said. He put his arms round her thin little body and squeezed it.
- I just put my arms around her and squeezed her so tight I thought I 'd break her ribs, ' said Rosemary Lodge.

Structure: intransitive

► sb squeezes

- Michael Collins wrapped his arm around her shoulder and squeezed affectionately. "You 've worked wonders, as usual."
- She threw her arms around him and squeezed tightly.

C with your fingers

a PERSON squeezes an OBJECT with his/her fingers

Structure: transitive

► sb squeezes sth

- Hold the rolled up condom at the tip of his erect penis squeezing the teat as before.

D activate a device

a PERSON squeezes an activating PART OF DEVICE using his/her fingers

Structure: transitive

► sb squeezes sth

- to squeeze [a trigger]

SSSJ-baza

Documents

- Look up Design
- globalen
- gleženj
- globalen**
- glodalec-red-k
- gmota
- gmoten
- gnan
- gnati
- gneča
- gnezditi
- gniti-red-k
- gojen
- gojenje
- gojišče
- gojitelj
- gojiti
- golf
- gora
- gorat
- goronzola
- gorski
- gosenica
- gospodinja
- gospodinjiti
- gospodinski
- gospodinstvo
- gostja
- gostoljubnost
- gotovina
- govedo
- govor
- govorjen
- gozdarstvo
- grabiti
- grablje
- grad
- gradnja
- grafika
- grafik-red-k
- grajski
- gramofon
- gramofonski
- grba
- grebsti
- gred
- greda
- grenkoba

globalen SSSJ-baza (Entry Document)

clanek

glava

oblika

zapis:globalen

iztocnica:globalen

zaglavje

besvrs:pridevnik

Pomenski meni

- 1. svetovni; mednarodni
- 1.1. splošno veljaven; razširjen
- 1.2. o spremembah v okolju
- 2. ki zadeva celoto; celostni

geslo

1.pomen

indikator:**svetovni; mednarodni**

pomenska shema:globalni PROCESI, zlasti EKONOMSKI in KOMUNIKACIJSKI, potekajo po vsem svetu

skladske skupine

a) struktura: PBZ0 sbz0

kolokacije

- globalna [ekonomija, konkurenčnost, recesija, korporacija]
- [razpravljati] o globalni ekonomiji
- globalna [komunikacija]
- globalni [trg]
- [obvladovati] globalni trg
- [obvladovati] globalno tržišče
- [konkurirati] na globalnem trgu
- globalno [pozicioniranje]
- globalni [problem, proces]
- globalna [razsežnost]

zglede

- Končal je s pričakovano izjavo, da **globalna** ekonomija ne more izničiti nacionalnih interesov.
- Kadar pravim, da moramo razpravljati o bazičnih stvareh v svetu **globalne** ekonomije, so ljudje prestrašeni.
- V New Yorku naj bi državniki in podjetniki razpravljali o **globalni** vamosti.
- Madžarska kot nekoč socialistična država je vstopila v **globalni** kapitalizem z velikimi koraki.
- Razvoj informacijske družbe ima precejšen pomen tako za **globalno** konkurenčnost kot za gospodarsko rast.
- Trgi v svetu pa se združujejo, tako da je edina perspektiva lahko le **globalni** trg.
- Tudi največje svetovne firme, ki danes obvladujejo **globalni** trg, so se razvile iz malih podjetij in obratov.
- Letos je bilo prodanih za 182 odstotkov več tovrstnih izdelkov kot pa lani, tako da Dell sedaj obvladuje 16 odstotkov **globalnega** tržišča.
- Ali pa so slovenska podjetja morda le dovolj usposobljena, da bodo učinkovito konkurirala na **globalnih** trgih?
- Seveda se v Trimu, ki konkurira na **globalnem** trgu v pogojih popolne konkurence, še kako zavedajo stalne dobre pripravljenosti na vse izzive konkurence.
- Teroristi so očitno imeli sisteme za **globalno** pozicioniranje in vnaprej pripravljene podatke o stavbah.
- GPRS bo mobilne telefone spremenil v kreditne kartice, z njimi bo možno tudi **globalno** pozicioniranje.
- Seveda pa take lokalne krize ne vodijo vedno v krizo **globalnih** razsežnosti.
- Ljudje kupujejo manj, kot so pred krizo, ki je poleti dosegla vrhunec in se skorajda razbohotila v krizo **globalnih**

Element panel

deepLink

Attributes

?	docref	globalen
?	senseref	svetovni;&mednarodni

Operations

Forward Validation

deepLink

Templates

Issues

- management tool
 - control workflow, monitor progress, distribute batches, backups etc.
- publishing tool
 - reusability, direct transfer (online, DTP)
- data management tool
 - reversing (bilinguals), data extraction etc.
- future
 - integration with browsers
 - online editing (mobile, cloud computing)

Termania – quick demo

- www.termania.net

The screenshot shows the Termania website interface. At the top, there is a navigation bar with the Termania logo, a search bar, and links for 'Search', 'Dictionaries', 'Welcome simplysimy', 'My account', 'Logout', and 'English'. Below the navigation bar, a search input field contains the word 'bulgaria' and a 'Find' button. A checkbox for 'Use the virtual keyboard to enter special characters' is checked, and a link for 'Advanced Search' is visible. The search results section shows 'bulgaria - 17 results found (0.28 seconds)'. On the left, there are filters for 'Source language' (English (12), Latin (2), Portuguese (2), Spanish (1), More...) and 'Target language' (Slovenian (10), German (8), Croatian (6)). The main results list includes: 'es Bulgaria' with a link to 'de Bulgarien ... More...' and source 'Spanish-German dictionary - Zeno Gantner, Matthias Buchmeier and others'; 'en Bulgaria (samostalnik)' with links to 'sl Bolgarija', 'de Bulgarien', 'sq Bullgari', and 'fr Bulgarie ... More...' and source 'Presisov večjezični slovar - Amebis, d. o. o., Kamnik'. On the right, there are two promotional boxes: 'Termania za Android' with a link to 'Naloži si ...' and 'Amebis Besana' with a link to 'http://besana_amebis.si/'.

termania | Search | Dictionaries | Welcome simplysimy | My account | Logout | English ▾

bulgaria

Use the virtual keyboard to enter special characters [Advanced Search](#)

bulgaria - 17 results found (0.28 seconds).

Source language:

- ▶ [English \(12\)](#)
- ▶ [Latin \(2\)](#)
- ▶ [Portuguese \(2\)](#)
- ▶ [Spanish \(1\)](#)
- ▶ [More...](#)

Target language:

- ▶ [Slovenian \(10\)](#)
- ▶ [German \(8\)](#)
- ▶ [Croatian \(6\)](#)

es Bulgaria
de Bulgarien ... [More...](#)
source: Spanish-German dictionary - Zeno Gantner, Matthias Buchmeier and others

en Bulgaria (samostalnik)
sl Bolgarija
de Bulgarien
sq Bullgari
fr Bulgarie ... [More...](#)
source: Presisov večjezični slovar - Amebis, d. o. o., Kamnik

Termania za Android
Slovarji na vašem telefonu
[Naloži si ...](#)

Amebis Besana
Slovenski slovnici pregledovalnik
za pomoč pri odpravljanju napak
http://besana_amebis.si/

Corpus Query Systems

Session 1, part 2

Topics

- We will talk about:
 - CQS & lexicographers
 - personal experience from Sara to the Sketch Engine
 - CQS & language learning/teaching
 - Gigafida interface
- We won't talk about
 - corpus design & annotation
 - corpora & NLP

Lexicographers

- Data collection:
 - corpora and other forms of evidence
- Data analysis:
 - sifting the evidence to discover relevant facts
- Synthesis:
 - creating dictionary entries

Data analysis in 2010s

- Lexicographer's changing role
 - from scanning data, to identify lexicographically-relevant facts
 - to validating (or rejecting) decisions made by computer
- New role
 - Identify/describe what can be automated & expand set of automatable processes
 - Identify weaknesses in support software

Data analysis (OED 1910s)

1896 Cosmopolitan XX. 356/2
Near Herbert Island I secured a
goodly number of walruses — cows, calves,
yearlings and two-year-olds

(See entry walrus - calf)

Data analysis (COBUILD 1980s)

m br 132
 m br 25 look and cranny in the vessel where even a stray seagull could be hiding and you can take my word
 er. "every wave on the atlantic was like a dead seagull drag- ing its driftwood artillery from h
 8 seal

r br 127
 r br 129 e king's taster?' i looked at the unbroken lead seal. 'not uness you think some- one has brought
 r br 140 church until 1835. years later, galileo put his seal on copernicus's discovery, wvas hauled up b
 m br 172 her lover to assuage her inner doubts, set the seal on her femininity, provide her with psychic
 m br 14 brooding darkncss is lifted? could the seventh seal or winter light have been conceived in anot
 a br 34 that never cleaned anything away, heavy thermal seal over diesel fuel, mildew, garbage, excremen
 a br 127 s foot in it. lynn tried to be gracious but the seal was set on her dislike of him. and somethin
 a br 82 plant aboard.' 'i've checked.' smithy broke the seal. "we talked last night. at least, i did. yo
 ce she discovered that, lynn thought, the final seal would be set on jane's hatred and rebellion
 10 sealed

r br 133
 r br 129 ingenu. both their fates were, to some extent, sealed. after "bunty" closed he went sady back t
 a br 151 place strips of the paper in a thin rubber tube sealed at one end and connected to suction at th
 a br 138 d as superior and knowledgeable. a partnership. sealed by why? so many exquisite little symmetri
 a br 135 g was led.' the europea party swept to dover, in sealed cars through back streets. 'you were a lo
 a br 86 ss asked "thank you. i am not fond of salad.' a sealed envelope passed to the prime minister wit
 a br 132 c forms and filled them out. i put those in one sealed envelope, the signed affidavit - i just
 a br 80 ottago he would flee to when all was signed and sealed. he hadn't had a proper night's sleep for
 r br 84 ed. "on a night like this? no fear. the gash is sealed in polythene bags, then they're punctured
 r br 199 lions of years but his doom, paradoxically, was sealed in the very fact that he became too perfe
 ote out the telegram, put it into its envelope, sealed it and handed it over to dolly. the four
 2 sealing

m br 48
 r br 21 gon stream thinning and trickling out: frontier sealing, cencus grievance, black operations (pre
 m each other. our once one-flesh divided again, sealing me into me, him into him. he is now a te
 3 seals

m br 3
 m br 35 saw a row of old houses, huddled together like seals on a rock. then there was a long field tha
 r br 161 ang we'd get stone together and keep the lurps, seals, recondos, green-beret bushmasters redunda
 omen serves only their own artificial needs and seals them off in their folie a deux from the re

Hanks on COBUILD

- "For the first time ever, lexicographers had evidence that enabled them to begin to see how words actually fitted together."
- "Opening a fresh concordance for all uses of a given word in a 7.3 million word corpus was like opening a window on a landscape of fresh snow on a sunny winter's morning."
- "Very often, a simple right-sort of concordances revealed patterns of usage that were at the same time unexpected and yet obvious (once seen)."

BNC – BoE (1995-)

- SARA (SGML-Aware Retrieval Application)
- Bank of English (HarperCollins)
 - 50 million (<http://titania.cobuild.collins.co.uk/form.html>)

Word Query

sing Pattern

Word	TAG	Frequency
sing	SUBST	28
sing	UNC	1
sing	VERB	6005
sing-	ADJ	1
sing-a-long	ADJ	13
sing-a-long	VERB	1

Form	POS	Frequency
sang	VVD	1119
sing	VVB	678
sing	VVI	1302
sing	VVB-NN1	47

555 words Controls Forms

Download Lemmata

A lemmatisation scheme determines how individual words are grouped under headwords. Select from the options below.

lancaster

The Lancaster scheme was developed at the University of Lancaster.

Example

doing

Headword do=VERB

<w TO0>to <w VBI>be <w VVN>sung<c PUN>, <w AV0>thereby <w AV0>then <w AV0>half <w VVN>sung <w PNP>it<c PUN>:<poem> TO>the <w NN2>ladies <w VVN>sung <w NN1-NP0>Erse <w NN2>w AV0>only <w VBI>be <w VVN>sung <w CJC>but <w AV0>also <w X0>not <w AV0>often <w VVN>sung<c PUN>, <c PUQ>'<w AJ0>novel <w CJC>and <w VVN>sung <w PRP>by <w AJ0>Junior <w NN1>solo<c PUN>, <w VVN>sung <w PRP>to <w DPS>his <w I>)>vocal <w NN1>music <w VVN>sung <w PRP>by <w AT0>the <w girlie <w CJT-DT0>that <w VVN>sung <w AT0>the <w NN1>song <w PNP>them <w CIS>as <w VVN>suna <w PNP>them<c PUN>.</u>

FIDA corpus (1997-2000)

- online (dial-up, 56 kbit/s modems)
- 100 million words
- POS-tagged (rule-based tagger)
- (complex) metadata
- consortium: **publisher**, LT company, university, technical institute
 - free for partners, restricted access
- purpose: lexicography, education, LT



KORPUS SLOVENSKEGA JEZIKA

[Enovrstično iskanje](#)

[Razširjeno iskanje](#)

[Rezultati iskanj](#)

[Spiski \(datoteke\)](#)

[Pomoč](#)

[Nastavitve](#)

[Novice](#)

[Servisna stran](#)

ODJAVA

Kontaktni naslov:

Korpus FIDA

DZS, Založništvo literature

Mestni trg 26

1538 Ljubljana

fida@dzo.si

Kratek opis točk

Enovrstično iskanje: Enostaven ali zahteven pogoj se napiše v eni vrstici. Uporabljajo se lahko logični in besedni operatorji. Postavijo se lahko tudi najbolj zapleteni pogoji, seveda pa morate biti vešči uporabe raznovrstnih operatorjev.

Razširjeno iskanje: Dodatni iskalni pogoji.

Rezultati iskanj: Shranjeno je do 10 prejšnjih iskanj.

Spiski (datoteke): Delo s spiski konkordanc in statistike (pregled, brisanje), ki smo jih shranili v besedilne datoteke.

Nastavitve: Uporabniške nastavitve.

Novice: Novice o uporabniškem vmesniku.



Enovrstično iskanje:

#1 zlikovec

išči briši

- kode MSD
- zapis posebnih znakov pri iskanju

Vrsta iskanja	Primer	Opis primera
enostavno iskanje besede	mize	poišče vse pojavitve besede "mize"
iskanje z nadomestnimi znaki (? nadomesti eno črko, * nadomesti poljubno zaporedje črk)	miz* miz?	poišče vse pojavitve besed, ki se začnejo z nizom "miz" poišče vse pojavitve štiričrkovnih besed, ki se začnejo z nizom "miz"
iskanje po kanalih (#1 lemma, #2 msd, #3 lemmas, #4 msds)	#1miza #2pkomein	poišče vse pojavitve besed z osnovno obliko "miza" poišče vse pojavitve pridevnikov (kakovostnih, nestopnjevanih...)
iskanje po frazah	#1okrogel_#1miza ki_*_je	poišče vse pojavitve, kjer sta zapored besedi z lemama "okrogel" in "miza" poišče vse pojavitve, kjer je med besedama "ki" in "je" še natanko ena beseda
iskanje po bližini (// za privzeto bližino ali /0 do /9 za število vmesnih besed)	#3stol//#3miza se/0je	poišče vse pojavitve, kjer sta lemi "stol" in "miza" blizu skupaj poišče vse pojavitve, kjer med "se" in "je" ni vmesnih besed (se_je ali je_se)
notranji in	#3vodavod	poišče vse pojavitve besed, pri katerih je možna tako lema "voda" kot "vod"
notranji ne	#3vod&~#1vod	poišče vse pojavitve besed z možno lemo "vod", kjer je ta možnost pri analizi odpadla
in	se je	poišče vse pojavitve besed "se" in "je" v odstavkih, v katerih nastopata obe besedi
ali	se je	poišče vse pojavitve besed "se" in "je" v odstavkih, kjer



Razširjeno iskanje:

Pogoj:

- kode MSD
- zapis posebnih znakov pri iskanju

Od leta: Do leta:

Prenosnik:

vsa besedila

Ft.P prenosnik
 Ft.P.G govorni
 Ft.P.E elektronski
 Ft.P.P pisni
 Ft.P.P.O objavljeno
 Ft.P.P.O.K knjižno
 Ft.P.P.O.P periodično
 Ft.P.P.O.P.C časopisno
 Ft.P.P.O.P.C.D dnevno
 Ft.P.P.O.P.C.V večkrat tedensko
 Ft.P.P.O.P.C.T tedensko
 Ft.P.P.O.P.R revialno
 Ft.P.P.O.P.R.T tedensko
 Ft.P.P.O.P.R.S štirinajstdnevno
 Ft.P.P.O.P.R.M mesečno
 Ft.P.P.O.P.R.D redkeje kot na mesec
 Ft.P.P.O.P.R.O občasno
 Ft.P.P.N neobjavljeno

Zvrst:

vsa besedila

Ft.Z zvrst
 Ft.Z.U umetnostna
 Ft.Z.U.P pesniška
 Ft.Z.U.R prozna
 Ft.Z.U.D dramska
 Ft.Z.N neumetnostna
 Ft.Z.N.S strokovna
 Ft.Z.N.S.H humanistična in družboslovna
 Ft.Z.N.S.N naravoslovna in tehnična
 Ft.Z.N.N nestrokovna

Lektorirano:

vsa besedila

Ft.L lektorirano
 Ft.L.D da
 Ft.L.N ne



IZVOR	ODSTAVEK	KONKORDANCA
0021271.0000004	toda če kje obstaja kdo, ki je zares pokvarjen	zlikovec , nevaren človek, izdajalec d
0021404.0000053		Postojna - Neznani zlikovec je v Postojni vlomil v osebni
0021404.0000055	oškodoval za približno 8000 tolarjev. Nedaleč stran pa je	zlikovec pred neko stanovanjsko hišo
0021450.0000893	da si ga bodo za vedno zapomnili tako kolegi kot	zlikovci . Ima dar, da se med akcijam
0021450.0000895		Nekega dne dva zlikovca oropata neko bergamsko ba
0021486.0000037		Ilirskobistriški policisti paiščejo zlikovca , ki je med zadnjim vikendor
0021608.0000237	smrti ne gre nihče v pekel. V Duhu so	zlikovce kakopak čakale hude, peko
0021766.0000003	lastnino, odkrivajo mamila, vlečejo nam sani, preganjajo	zlikovce , lovcem prinašajo divjad, v
0021830.0000034	o dogodku, ko je z železniške postaje Štore neznani	zlikovec v lokomotivo potniškega vli
0021865.0000577	Meek je še en mojster borilnih veščin, ki mu	zlikovci ubijejo starše. Oh, ko bi prej
0013788.0000008	na prvo hišo? Naj si mislijo, da je	zlikovec pač nekje na tem območju, :
0013788.0000042	-ja in ga bomo v zgodbi o policijskem iskanju	zlikovca
0014192.0000084		Lopati, last A. L. iz Kočevja. Zlikovec je tudi odlomil pipo na večji
0014352.0000287		Nekakšen zlikovec z melono in psom mi je izma
0011028.0000002		Sezona zlikovcev je tu
0011028.0000004	polnem razmahu, temu primerno pa so se organizirali tudi	zlikovci . 'Ordinirajo' po plažah, vlam
0011028.0000005	bilo škode za 30 tisočakov. V začetku tedna je	zlikovec iz kioska na plaži v Portorož
0011319.0000013		Po dveh letih in pol je neznanega zlikovca spet zaplela vest: šel je pon
0011319.0000016		Nekulturno dejanje neznanega zlikovca je prijavljeno policiji v Šiški,
0011324.0001691	opisuje zaradi humane razneženosti zaslepljeno tretmansko oko. Gre za	zlikovce , hudodelce, žeparje, tatove
0011410.0000072		ima F. K. iz Novega mesta barako. Zlikovec je razbil vse mize in klopi pr
0011662.0000262	ekscsesni negativci, na smrt obsojene zverine, morilski krvoločni	zlikovci najhujše sorte, padejo v nem
0011662.0000262	negativcev, še hujših zverin, ja, morilskih krvoločnih	zlikovcev še hujše sorte, okej, neke v



Statistika

Iskanje po:

Besed pred zadetkom:

Besed za zadetkom: [✓ izvedi](#)

Urejanje

Opozorilo: urejanje daljših seznamov lahko traja precej dolgo. Ne priporočano urejanje seznamov daljših od 5000 besed.

Urejanje po:

Upoštevaj označeno besedo

[✓ izvedi](#)

Sito



KONKORDANCA

pa vem le to, da nič ne vem. **Bojim se Boga kot hudiča** . Ali v vas lahko zaupam?

, je po papeževi želji leto Očeta. Pa tudi **hudiča se ni treba bati** , ker ga ima Bog v oblasti in nam ne od antidepresivov začneš rediti. Vse od takrat se jih **bojim kot hudič križa** .)

Tako slovenstvo nekoč in danes kot vsak človek sta stalen **boj med Bogom in hudičem** .

prastari so tudi ljudje, ki se trave še vedno **bojijo kot hudič križa**.

izstopiti iz koalicije s SDS (tega se Alojz Peterle **boji kot hudič križa**, zato je že napovedal možnost, da se

. Razlog? Izjemno visoke kazni, katerih se vsi **bojijo kot hudič križa**.

emnoge pustolovščine: junaška družčina se na svojih vandranjih ni **bala ne boga ne hudiča** , možje so trepetali samo pred meglo, iz roke za sedeče kinoušesno-drugogodbno občinstvo, ki se **boji dratarjev in tehnologov kot hudič križa?**

noben ni upal iti čez, ker so se vsi **bali hudiča** .

prodajalci, na primer tile iz ALG, se javnosti **bojijo kot hudič križa**. Najprej, globoko v osemdesetih, je bila

! Le kako sem zašel v ta pekel " **Hudič mu je začel v predsmrtni boj** prišepetavati besede, ki jih ji

bi zame poiskali knjige o vedeževanju. Tarota sem se **bala kot hudiča** . Na začetku sem vedeževala iz običajnih kart.

reprezentance BiH, vele mojster, ki se za šahovnico ne **boji niti hudiča** , igralec, ki že od mladostrnih nog ne spoštuje vrsto industrijske politike, ki se je sicer ameriška ideologija **boji kot hudič križa**. Ne glede na motiv: če je država

, ki se, če smo malce ironični, kot **hudič križa boji** javne diskusije o človekovi prvinskosti. Strip je .

in je z njo naravnost obseden. Jaše ga " **hudič, nečista žival**", **boji** se glasu krvi. Čim bolj ta glas duši,

Rim - Italijani se kot **hudič križa bojijo** množičnega prihoda beguncev s Kosova. Za v

tekstom, a tudi v takih primerih se avtorji kakor **hudič križa bojijo** pridodati denimo glagole, s katerimi je često v -reprezentanco ali mlade, predvsem Primoža pa so se **bali "kot hudič križa"**. Njegova velika kristalna globusa sta bila n

To je krik sprevržene vere, vere **hudičev in brezbožnežev: Boga se bojijo** .

v dveh tekmah. Morebitnega tretjega srečanja v Ljubljani se **boji kot hudič križa**, tudi zato, ker bo sodil Litvanec Brazauskas

Bolnišnice in dohtarjev se **bojim kot hudič križa**, takrat sem se jih veselil v čudni vznesenos

FidaPLUS (2003-2006)

- online (slowly to broadband)
- 620 million words
- POS-tagged (rule-based tagger, upgraded)
- metadata still complex (same taxonomy)
- consortium: publisher, LT company, **university**, technical institute
 - freely available online, authentication
- purpose: lexicography, education, LT

Aims

- From lexicographers to students
 - analysis, survey
 - booklet
 - simple query window
 - avoid flood of intimidating information
- BUT!
 - limited by the existing technology
 - DEMO

First glimpse of WS – 2002

- Adam Kilgarriff and Michael Rundell
- *Lexical Profiling Software and its Lexicographic Applications – a Case Study*
- Euralex 2002, Copenhagen

- Just for English
- Just a combination of collocation patterns

PP_about: ratio = 12.07 : 1, counts = 86

PP_with 665 18.65 :1	PP_between 135 13.2 :1	PP_at 61 1.74 :1	object_of 1637 1.69 :1	before_prepn 2651 1.33 :1
friend 19	princess 8	table 4	overhear 33	topic 93
stranger 5	Charles 4	dinner 3	steer 25	snatch 16
passenger 5	woman 6	time 7	record 46	lull 14
people 22	people 8	party 4	tape 11	deep 34
man 17	officer 3	moment 3	tap 16	listen 57
artist 5	child 4	school 3	resume 14	eavesdrop 9
pupil 5			hold 63	engage 30
girl 7			interrupt 14	hum 11
woman 11			continue 42	buzz 10
mother 7			prolong 8	babble 7
visitor 4			make 120	transcript 11
colleague 4			hear 38	deep 26
speaker 4			conduct 13	gist 6
lady 4			finish 19	recording 17
wife 5			end 20	murmur 7
minister 6			start 34	participant 13
husband 4			remember 27	fragment 11
queen 3			recall 12	indulge 7
adult 3			dominate 10	extract 8
doctor 4			keep 34	course 31

The real deal – 2004

- Adam Kilgarriff, Pavel Rychly, Pavel Smrz, David Tugwell
- *The Sketch Engine*
- Euralex 2004, Lorient

- Sketch Grammar (cross-lingual)
- Sketch Difference, Thesaurus

Slovene Word Sketches (2006)

- Simon Krek and Kilgarriff
- *Slovene Word Sketches*
- Proceedings of the 5th Slovenian/First International Language Technology Conference, Ljubljana, Slovenia
- First version in 2005

čas

<u>coord</u>	<u>7334</u>	<u>0.5</u>
prostor	<u>837</u>	44.63
datum	<u>144</u>	37.89
trud	<u>75</u>	35.92
kraj	<u>247</u>	34.76
Vera	<u>33</u>	31.35
energija	<u>155</u>	30.85
denar	<u>206</u>	27.91
se	<u>192</u>	26.5
bit	<u>35</u>	25.96
zaslonka	<u>13</u>	25.71
potrpljenje	<u>21</u>	24.85
trimesečen	<u>17</u>	23.36
da	<u>97</u>	23.18
on	<u>131</u>	21.66
svoj	<u>26</u>	20.95
ne	<u>63</u>	19.83
napor	<u>25</u>	19.0
tudi	<u>43</u>	18.33
kraja	<u>24</u>	18.02
njegov	<u>35</u>	17.97

glava

<u>is obj4 of</u>	<u>2506</u>	<u>3.9</u>
skloniti	<u>98</u>	54.21
beliti	<u>83</u>	49.81
dvigniti	<u>218</u>	49.15
razbijati	<u>71</u>	44.74
odsekati	<u>60</u>	43.2
pomoliti	<u>45</u>	42.62
nagniti	<u>59</u>	40.97
tiščati	<u>46</u>	40.92
stakniti	<u>39</u>	39.75
obrniti	<u>90</u>	34.92
odrezati	<u>44</u>	33.54
nasloniti	<u>29</u>	32.21
sploščiti	<u>20</u>	32.18

<u>a modifier</u>	<u>19068</u>	<u>1.4</u>
sklonjen	<u>157</u>	59.54
obrit	<u>66</u>	47.71
kronan	<u>35</u>	39.24
dvignjen	<u>62</u>	36.12
zeljnat	<u>23</u>	35.75
odsekan	<u>33</u>	35.5
Hermanov	<u>34</u>	35.11
razgret	<u>40</u>	34.63
mrtvaški	<u>34</u>	31.12
trezen	<u>41</u>	30.65
bikov	<u>22</u>	29.33
ovnov	<u>16</u>	29.12
video	<u>106</u>	29.03
koničast	<u>34</u>	28.96
pobrit	<u>13</u>	27.94
zeljen	<u>18</u>	27.92
bister	<u>35</u>	27.26

<u>prec po</u>	<u>1090</u>	<u>9.0</u>
rojiti	<u>91</u>	65.67
udariti	<u>142</u>	55.22
motati	<u>49</u>	51.25
popraskati	<u>24</u>	42.95
treščiti	<u>38</u>	39.83
poditi	<u>32</u>	38.99
tolči	<u>34</u>	37.54
lopniti	<u>12</u>	32.97
tepsti	<u>20</u>	32.27
praskati	<u>13</u>	28.88
blođiti	<u>12</u>	28.06
plesti	<u>14</u>	27.2
trepljati	<u>9</u>	26.25
poškodovati	<u>23</u>	25.79
čohati	<u>6</u>	25.61
pobriti	<u>5</u>	24.37
srati	<u>8</u>	23.96

Sketch

Difference

cona = zone

območje = area

green = cona

red = območje

a_modifier	4941	36173	2.6	1.9
erogen	<u>75</u>	<u>13</u>	58.6	21.4
obrten	<u>288</u>	<u>9</u>	58.3	5.6
ekonomski	<u>411</u>	<u>19</u>	54.5	5.2
carinski	<u>271</u>	<u>214</u>	53.4	33.9
industrijski	<u>284</u>	<u>76</u>	52.2	19.4
prost	<u>317</u>	<u>79</u>	50.1	16.7
prostocarinski	<u>61</u>	<u>21</u>	47.3	21.8
obmejen	<u>7</u>	<u>313</u>	11.5	46.7
moder	<u>139</u>	<u>20</u>	41.4	8.2
okupacijski	<u>35</u>	<u>10</u>	38.6	14.6
nekdanji	<u>41</u>	<u>710</u>	16.6	37.3
tamponski	<u>17</u>	<u>12</u>	35.8	23.1
demilitariziran	<u>9</u>	<u>30</u>	25.3	33.8
zaprt	<u>5</u>	<u>162</u>	7.6	32.4
brezcarinski	<u>26</u>	<u>17</u>	30.3	16.5
konvergenčen	<u>12</u>	<u>6</u>	29.7	15.0
operativen	<u>34</u>	<u>47</u>	28.5	21.8
bivši	<u>15</u>	<u>199</u>	11.9	27.4
posamezen	<u>9</u>	<u>338</u>	5.1	26.8
svoboden	<u>42</u>	<u>43</u>	25.5	14.9
turističen	<u>16</u>	<u>183</u>	11.0	23.6
visok	<u>8</u>	<u>353</u>	3.1	23.5
koprski	<u>12</u>	<u>110</u>	11.8	23.2
mesten	<u>11</u>	<u>190</u>	7.2	21.7
siv	<u>23</u>	<u>14</u>	20.9	8.1

19 [spor](#) 0.273 [spopad](#) 0.267 [vojna](#) 0.266 [nasilje](#) 0.202 [boj](#) 0.17

03 [težava](#) 0.288 [razmera](#) 0.279 [situacija](#) 0.264 [dogajanje](#) 0.228 [dogode](#)
[problematika](#) 0.18 [potreba](#) 0.172 [stvar](#) 0.172

25 [katastrofa](#) 0.243 [padec](#) 0.198 [zlom](#) 0.185 [razpad](#) 0.179 [izbruh](#) 0.17

24 [posledica](#) 0.21 [pomanjkanje](#) 0.209 [nevarnost](#) 0.188 [pritisk](#) 0.182 [vp](#)

29 [politika](#) 0.196 [proces](#) 0.189 [razvoj](#) 0.184 [sila](#) 0.182 [gibanje](#) 0.182 [o](#)

23 [recesija](#) 0.211 [revščina](#) 0.178

21

07 [nesreča](#) 0.196 [bolečina](#) 0.169

05 [reforma](#) 0.203 [napad](#) 0.196 [volitve](#) 0.178 [poseg](#) 0.176 [akcija](#) 0.171

93 [napetost](#) 0.19

87

82 [poraz](#) 0.175

Thesaurus

kriza = crisis

Gigafida (2012)

- online (broadband)
- 1,2 billion words
- POS-tagged (statistical tagger, better!)
- metadata simplified
- consortium: ~~publisher~~, LT company, university, technical institute, language institute
 - freely available online, ~~authentication~~
- purpose: education, lexicography / LT

The split

- Communication in Slovene project
 - 2008-2013
- Gigafida corpus
 - Slovene Lexical Database (dictionary)
 - Sketch Engine
 - for lexicographers
 - online concordancer(s)
 - for education and general public

Web concordancers

- Log analysis of FidaPLUS concordancer
- FidaPLUS web survey
- Analysis of existing corpus tools
- Analysis of popular web tools (Google etc.)

Survey – findings

- Simple search – regularly used by 72% users
- Advanced search – rarely used (only 8% use it regularly)
- Lack of intuitiveness
- the manual is almost key to learning how to use a corpus tool
 - “...if you are not using the interface for a while, you forget what the search commands are, and you don’t (want to) bother with looking into the manual”
 - “...the interface should have a **modern design**, it should be more **user-friendly**, and its use should be clear and transparent”

Main design principles

- similarity to the well-known non-linguistic tools (e.g. Google)
- No registration
- Minimum navigation
- No redundant functions (less is more)
- Simplicity of searches
- Help and tips in pop-up windows
- Simple descriptions of functionality (no terminology)

Concordancers

- For lexicographers
 - NoSketchEngine
 - <http://nl.ijs.si/noske/index-en.html>
 - Sketch Engine
 - <http://ske.amebis.si/>
- For education / general use
 - Gigafida
 - <http://www.gigafida.net/>
 - Kres
 - <http://www.korpus-kres.net/>
- DEMO