

DEB II – Lexikografická platforma pro vývoj slovníkových aplikací

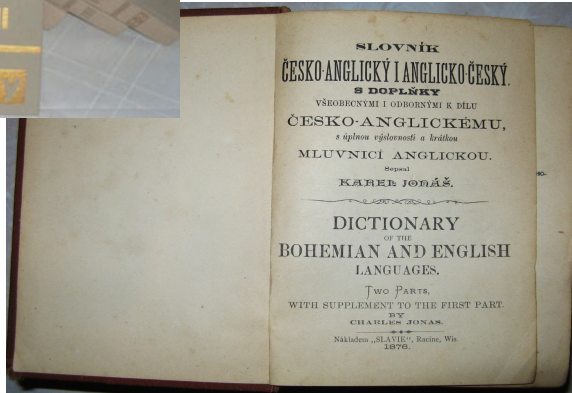
Adam Rambousek

Centrum zpracování přirozeného jazyka
Fakulta informatiky, Masarykova univerzita

`xrambous@fi.muni.cz`

`deb@aurora.fi.muni.cz`

`http://deb.fi.muni.cz`







writer.odt - OpenOffice.org Writer

File Edit View Insert Format Table Tools Window Help

Heading 3 Arial 14 B / U

17:40 18:30	Panel Discussion: Open Office Community	Looking into a career - can the community get the green light?	The Freelance Council	Live in Paris? OpenOffice.org
18:45 19:30		Community Roundtable on Marketing: Marketing 2.0: What's of the Future of Jobs	Visualize the User Experience: Developer's Environment	Handling records: a new in our field: a new approach to data strategies & tools
19:30-	Welcome Dinner			

myr Community ODF XML #OpenOffice.org #OpenOffice.org

Community

What's New in OpenOffice.org 3.0?

Track: Community Type: Presentation Audience: Community, all

Abstract: From Mac OS X Support to tables in Impress, OpenOffice.org 3.0 offers a wide range of improvements. This demo will showcase the highlights of the new and improved features of OpenOffice.org 3.0.

Presenter: Christian Jansen, **Bettina Haberler** (User Experience Engineer, Sun Microsystems, Inc.) Biography: Christian Jansen is a User Experience Engineer at Sun Microsystems. During the last 3 years, he played a major role in designing the user interfaces of

Page 4 6 / 4 Default English (UK) INSERT STD Level 3

The screenshot shows the 'Entry Editor' interface for a dictionary entry. The main window displays the entry for 'Wednesday' with its definition and etymology. The left sidebar contains navigation and search options. The right pane shows the underlying TEI XML structure for the entry.

Project Information:
 Project: Wednesday
 Headerword: 1
 Senses: 1
 Words: 74 Characters: 443

Entry Content:
Wednesday ■ **noun** the day of the week before Thursday and following Tuesday. ■ **adverb** chiefly N. Amer. on Wednesday. ► (**Wednesdays**) on Wednesdays, each Wednesday.
 –ORIGIN OE *W[odnes]dæga*, named after the Germanic god Odin; translation of late L. *Mercurii dies*.

TEI XML Structure:

```

    <entry>
      <word>Wednesday</word>
      <partofspeech>noun</partofspeech>
      <definition>the day of the week before Thursday and following Tuesday.</definition>
      <partofspeech>adverb</partofspeech>
      <usage>chiefly N. Amer. on Wednesday.</usage>
      <usage>► (Wednesdays) on Wednesdays, each Wednesday.</usage>
      <etymology>–ORIGIN OE W[odnes]dæga, named after the Germanic god Odin; translation of late L. Mercurii dies.</etymology>
    </entry>
  
```

Slovníky a počítače

- 60. léta - používají se počítače, lexikografové píší na papír, specialisté přepisují do databáze, Brown Corpus
- 1978, *Longman Dictionary of Contemporary English*
 - první s omezeným slovníkem definicí, kontrolováno stroje
 - kódování pro NLP výzkum
- 1980, *COBUILD*, University of Birmingham + Collins
 - korpus současných textů (Bank of English)
 - 1987, *Collins COBUILD English Language Dictionary*
 - první slovník založený na korpusových datech
 - nový styl definice - celé věty
 - *If a person, animal, or other living thing is killed, something or someone causes them to die.*
- 90. léta - vývoj specializovaných systémů pro tvorbu slovníků
- 1987, Text Encoding Initiative

XML

- PB138 Moderní značkovací jazyky
- eXtensible Markup Language - značkovací (meta)jazyk
- pravidla, jak má vypadat správně vytvořený dokument - snadné strojové zpracování a výměna informací
- konkrétní názvy značek určuje uživatel (standards, vlastní)

Zobrazení

- **XSLT** – eXtensible Stylesheet Language (Transformations)
- převod XML na jiné formáty
 - jiné XML značkování, text, HTML, LaTeX, PDF
- šablony pro části XML dokumentu, postupné procházení dokumentu
- funkcionální programovací jazyk

ssjc Slovník spisovného jazyka českého

lov

-u m. (6 j. -u)

1. *stíhání a zmocňování se zvířete (nejč. odstřelem); chytání ryb*: l. jelenů, divokých kachen, velryb; l. lososů; l. perel; doba lovu; uspořádat l. na medvědy; vyjet na l.; právo lovu; l. odstřelem, chytáním, lapáním; l. lesní, poliň, vodní; hromadný l. hon; liška vyšla na l.; lovu zdar! (*lovecký pozdrav*)
2. *expr. chytání, shánění čehokoliv, vůbec získávání, při kterém se uplatní obratnost a náhoda*: l. vzácného lumyzu; sběratelé se vydali na l. lidových písní; policie podnikla l. na zloděje; *expr.* to je l! *řízavý nálezk. výhodná koupě ap.*
3. *výsledek lovu; úlovek, kořist*: vrátit se s bohatým lovem z ulovenou zvířítí ap., *přen. expr.* z věcmi získanými obratností n. šťastnou náhodou

SSC Slovník spisovné češtiny

lov

-u m

1. *lovení zvířete a ryb* koroptví, lov na zajíce, liška vyšla na lov,
2. *úlovek (syno) kořist (syno)* má bohatý lov,

Ukládání

- XML databáze
- ukládají se přímo XML dokumenty
- vyhledávání - XPath, XQuery
- např. eXist, BaseX, Sedna

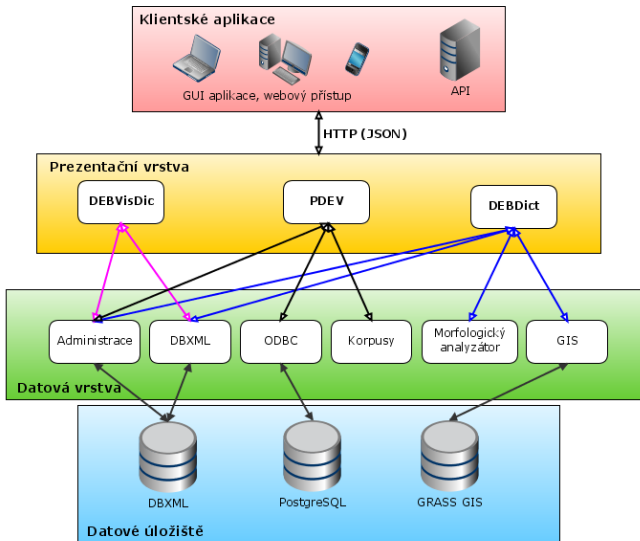
TEI

- *Text Encoding Initiative*, <http://www.tei-c.org/>
- *TEI Guidelines* (aktuálně verze 5 z roku 2007)
- XML formát pro sémantický popis textových dokumentů
- velký rozsah značek
- *TEI Lite* – osekaná verze, "90 % potřeb 90 % uživatelů"
- romány, poezie, divadelní hry, dokumentace, slovníky, korpusy, grafy, rukopisy, zarovnání, odkazy, změny textu, notové zápisy...
- nástroje - sada XSLT pro převod na LaTeX, docx, EPUB, HTML

Úvod

DEB – Dictionary Editor and Browser

- platforma pro vývoj slovníkových aplikací
 - všechna data ve formátu *XML*, *Unicode*
 - libovolná struktura, jakýkoliv jazyk
- architektura klient-server
- *server*
 - práce s daty, velká část funkcí
 - databázové úložiště
 - správa uživatelů, slovníků, spolupráce
 - rozděleno do modulů, spojování podle potřeb
- *klient*
 - omezená funkcionality
 - grafické nebo webové rozhraní



DEBDict

- prohlížeč slovníků
- 7 slovníků češtiny pro veřejnost, další přístupné jen pro část uživatelů
- napojení na morfologický analyzátor
- data z externích zdrojů
- přes 700 uživatelů (podepsané prohlášení)

SSK slovník spis. jaz. českého

SSS slovník cizích slov

SSČ slovník spis. češtiny

SSS slovník českých synonym

SFS slov. fráze a idiomy

SFI neslov. fráze a idiomy

PSK příruční sl. jaz. českého

všechny slovníky

Český WordNet

Concise Oxford English D.

Oxford Dictionary of English

Oxford Thesaurus of English

Slovinský slovník

Komplexní slovník ruštiny

Diderot

Google

Answers.com

Wikipedia

Seznam Encyklopedie

CIA World Factbook

mapa ČR

morf. analyzátor ajka

ssjc Slovník spisovného jazyka českého

kód

-u m. (6. j. -u) (z fr.) odb. předpis pro převod jedné soustavy znaků pro sdělování zpráv do jiné takové soustavy; používat kód; tajný k. známý jen určitému okruhu osob telegrafní k. předpis pro převod elektrických, optických n. jiných signálů do soustavy jazykových znaků; **kódový** příd.: k. název; k-á zkratka

kódovací

příd. odb. týkající se kódování, sloužící ke kódování; k. soustava

kódovati

ned. i dok. odb. (co) sestavovat, sestavit podle kódu; k. telegram

kódový

v. kód

SSC Slovník spisovné češtiny

kód

-u m <f>

1. systém znaků pro přenos informace telegrafní, dálnopisný kód, jazykový kód
2. výp. tech. pravidla pro jednoznačné přiřazení významu k znakům n. signálům

kódovat

ned. i dok. sestavovat, sestavit podle kódu kódovat zprávu,

kódovací příd. sloužící ke kódování, kódovací systém,

psjc Příruční slovník jazyka českého

Praled

- příprava lexikální databáze češtiny v Ústavu pro jazyk český

PDEV (CPA)

- editace slovníku vzorů anglických sloves
- varianty také pro češtinu, italštinu a španělštinu

TeDi – Terminologický slovník

- společný projekt s Fakultou výtvarných umění VUT
- glosář výtvarných pojmů
- multimediální prvky
- nově také Divadelní fakulta JAMU, Agronomická fakulta MU

jazyk
 Zdroj pro heslář fsc+ssjc
 Typ hesla jednosl. podstatné jméno
 STAT. Sg/Pl 12792/3985 FRQ 16777 ARF 6122 Zdroj

SEZNAMY (rozbalit/sbalit)
SLOVOTVORBA
 Derivovaná slova
 Fundace/motivace
 Defraz. lexémy

Zpracovatel VVeseley |
 Vytvořeno 2008-12-04 10:44 změněno 2009-09-14 12:26

Jazyk: ZÁHLAVÍ

Heslo **jazyk** Homonymie

Zdroj pro heslář fsc+ssjc

Typ hesla jednosl. Slovní druh/typ podstatné jméno

U zkratk: Plné znění zkratky

Pozn. k celému heslu

V1

Původ definice: SSJC **svalnatý, velmi pohyblivý orgán v dutině ústní (u zvířat v tlamě, zobáku atd.); orgán chuti, mluvy**

E1 k V1 Adj+SUBST (rozbalit/sbalit)
E2 k V1 SUBST+Adj. (rozbalit/sbalit)
E3 k V1 SUBST+Subst-gen (rozbalit/sbalit)
E4 k V1 Subst+SUBST-gen (rozbalit/sbalit)
 vyplazování jazyka;
 vyfíznutí jazyka;
 špička jazyka, kořen jazyka
 Pozn. k E4 zpracováván substantivum je samo genitivním
E5 k V1 SUBST+Prep+Subst/SUBST+Subst-ji

Pozn. k E5 následkovámi přídavnými ke zpracovávanému sub

Jazyk: FRAZÉMY

mit jazyk (ostrý) jako meč

Poznámka

 mit jazyk jako na obrtlíku

Poznámka

 mit jazyk jako poleno

Poznámka

 mlčí jako by /muj přimrzl jazyk

Poznámka

Entry	Filter	OEC	E
amuse	2	17444	4
anaesthetize	2	60	3
analyse	1	12828	4
anger	2	7467	8
angle	4	3189	2
anglicize	2	122	1
anchor	5	4693	4
animate	3	4929	1
anneal	1	221	5
annex	3	2277	2
annihilate	1	1831	1
annotate	1	1015	9
announce	4	92547	1
annoy	2	28130	5
annul	1	852	1
anoint	1	1285	1
answer	13	129214	9
antagonize	1	782	8
antedate	1	152	2
anthologize	1	93	2
anthropomorph...	1	80	3
anticipate	2	20741	2
ape	1	892	7
apologize	1	17983	4
apostrophize	1	16	7
appal	1	106	1
appeal	4	40303	4

answer: CPA Patterns

Patterns for: answer Add Copy Corpora Preview Renumber Delete Close

Save Sample size Semantic class Aspectual class

- 1 **[[Human]] answer [NO OBJ] (that [CLAUSE] | [QUOTE])**
[[Human]] says {that [CLAUSE] | [QUOTE]} in response to a question or statement by someone else
- 2 **[[Human]] answer [[Ask Activity]]**
[SUBJ][Human]] says or writes something intended to provide relevant information in response to someone else's [OBJ][Ask Activity]]
- 3 **[[Human]] answer [[Telephone]]**
[SUBJ][Human]] speaks into [OBJ][Telephone]] after it rings
- 4 **[[Human 1]] answer [NO OBJ] {to [[Human 2]] | to [[God]]}**
[SUBJ][Human]] has an obligation to account for his/her actions to [[Human 2 | God]]
- 5 **[[Human]] answer [[Mail]]**
[SUBJ][Human]] writes a letter in response to [OBJ][Mail]] from someone else
- 6 **[[Human]] answer [[Speech Act = Accusation]]**
[[Human]] says or writes something intended to refute [[Speech Act = Accusation]]

answer Pattern 1 show: Save Save & close Close Test

subject Human Role Lexset Attr.

verb form answer

object
 no object

adverbial add
 no adverbial

clausal optional to/INF [V] -ING that [CLAUSE] WH-[CLAUSE] [QUOTE]

primary implicature semantics
[[Human]] says {that [CLAUSE] | [QUOTE]} in response to a question or state: idiom

Count: 13

Patterns: 1383 Verbs: 356

TeDi, editace hesla adresa (grafická, dedikační) - SeaMonkey

File Edit View Go Bookmarks Tools Window Help

https://apollo.fi.muni.cz:8010/teDi?action=edit&id=adresa_g Search

Home Bookmarks

Autor: ID: **adresa (grafická, dedikační)-11804343501**

obor

heslo česky anglicky

německy francouzsky

varianty +

styl, příznak

definice

příklady +

nadpojem +

podpojem +

morfologie

Done

DEBVisDic

- editor slovníků typu WordNet
- samostatný modul pro každý jazyk (modifikace)
- použit pro tvorbu několika wordnetů
- poskytuje API – napojení externích aplikací
- možnost rozšíření a modifikací

