# Semantic Web Mining and the Representation, Analysis, and Evolution of Web Space
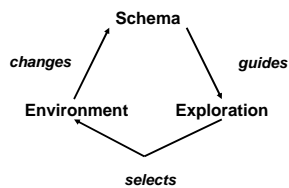
**Bettina Berendt,**

**Andreas Hotho, & Gerd Stumme**

Humboldt University Berlin / University of Kassel, Germany

More info: www.berendt.de

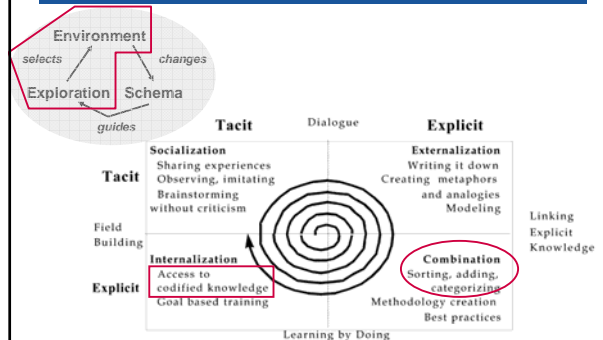**The Web is mankind´s largest repository of knowledge ...**

**... but knowledge isn´t something that can be "put in a container and then used as the need arises".**

**Knowledge is constructed in learning activities.**

**Knowledge$_1$ ("in people's minds") is created by interaction with the environment (e.g., Neisser, 1967)**

Schema

*changes*        *guides*

Environment        Exploration

*selects*

**Knowledge$_2$ (codified) is part of the environment; Learning accesses this knowledge (e.g., Nonaka, 1991)**

Environment

*selects*        *changes*

Exploration    Schema

*guides*

Tacit        Dialogue        **Explicit**

**Tacit**

**Socialization**
Sharing experiences
Observing, imitating
Brainstorming
without criticism

**Externalization**
Writing it down
Creating metaphors
and analogies
Modeling

Field
Building

Linking
Explicit
Knowledge

**Internalization**
Access to
codified knowledge
Goal based training

**Explicit**

**Combination**
Sorting, adding,
categorizing
Methodology creation
Best practices

Learning by Doing

1

## Slide 7

**Semantic Web Mining**



Tacit | Dialogue | Explicit

Tacit

**Socialization**
Sharing experiences
Observing, imitating
Brainstorming
without criticism

**Externalization**
Writing it down
Creating metaphors
and analogies
Modeling

Field Building

Linking Explicit Knowledge

Explicit

**Internalization**
Access to
codified knowledge
Goal based training

Learning by Doing

---

## Slide 8

**Approaches to the current Web's biggest challenges:**
**lots of data, human-understandable**

**Web Mining**
**extracts implicit knowledge**

**Semantic Web Mining**
• **use semantics to improve mining**
• **use mining results to generate semantics**

**The Semantic Web makes knowledge machine-understandable**

[Berendt, Hotho, & Stumme, *Proc. ISWC* 2002]
[-"- (Eds.), *Proc. WS Semantic Web Mining at ECML/PKDD* 2001 and 2002]
[Berendt, Hotho, Mladenic, van Someren, Spiliopoulou, Stumme (Eds.),
*Web Mining: From Web to Semantic Web*, 2004]

---

## Slide 9

**Agenda**

**Web Mining**

**(Semantic) Web**

---

## Slide 10

**Agenda**

**Web Mining**

**Semantic Web**



---

## Slide 11

**Extracting semantics from Web content & structure – ideas and examples**

**Using syntactic structure, semi-automatically learn**

- **Ontologies (build or extend Yahoo-like taxonomies; Web-scale example *KnowItAll*: „... such as ...", see Etzioni et al. 2004)**
- **Instances of concepts and relations in a given ontology (ontology population)**
  - **Technique: Information extraction**
  - **From textual information including tables**
    - **Krátky, Andrt, & Svátek**
  - **From visual information including text layout**
    - **Burget; Gatterbauer, Krüpl, Holzinger, & Herzog; Hassan & Baumgartner; Labský, Vacura, & Praks**
  - **From structure (hyperlinks)**
    - **Frivolt & Bieliková**
- **Interactive learning**
  - **Ceresna; Schindler, Arya, Rath, & Slany**
- **Re-using existing conceptualizations**
  - **Švihla & Jelínek**

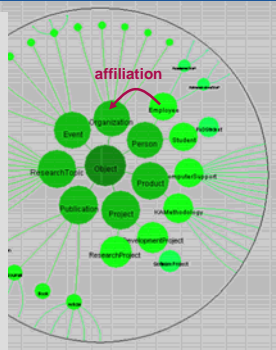*PS: This is just my understanding of your papers – please send me email if you find I've missed something!*

---

## Slide 12

**Agenda**

**Web Mining**

**(Semantic) Web**

## Agenda

**Web Mining**

**Semantic Web**

---

## Exploiting semantics for textual Web resources – ideas and examples

- **Allowing expert users to contribute ontologies for semantics-enhanced IE**
  - Schindler, Arya, Rath, & Slany
- **Using ontologies to build templates for the composition of Web services**
  - Svátek & Vacura
- **Use ontologies as additional structure on the tokens in a text to**
  - disambiguate meaning (e.g., word sense disambiguation: Navigli & Velardi, 2005)
  - reveal additional structure (e.g., clustering: Hotho, 2004)
  - help in the discovery of new ontological structures, or in instance learning (e.g., Navigli & Velardi, 2005; Hotho, 2004; KDD Cup 2005)

R. Navigli & P. Velardi. Structural Semantic Interconnections: a knowledge-based approach to word sense disambiguation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (27-7), 2005.

---

**WordNet 2.0 Search**

Search word: [          ]  [Find senses]

## Overview for "bus"

The noun "bus" has 4 senses in WordNet.

1. **bus**, autobus, coach, charabanc, double-decker, jitney, motorbus, motorcoach, omnibus -- (a vehicle carrying many passengers; used for public transport; "he always rode the bus to work")
2. bus topology, **bus** -- (the topology of a network whose components are connected by a busbar)
3. busbar, **bus** -- (an electrical conductor that makes a common connection between several circuits; "the busbar in this computer can transmit data either way between any two components of the system")
4. **bus**, jalopy, heap -- (a car that is old and unreliable; "the fenders had fallen off that old bus")

Search for [Synonyms, ordered by estimated frequency ▼] of senses [    ]
☑ Show glosses
☐ Show contextual help
[Search]

The verb "bus" has 3 senses in WordNet.

---

**WordNet 2.0 Search**

Search word: [          ]  [Find senses]

## Overview for "bus"

The noun [Synonyms, ordered by estimated frequency]
Coordinate Terms
Hypernyms (bus is a kind of...)
1. **bus**, aut [Hyponyms (...is a kind of bus), brief]  jitney, motorbus, motorcoach, omnibus -- (a vehicle carrying many passengers; used for public t [Hyponyms (...is a kind of bus), full]
2. bus topo [Holonyms (bus is a part of...), regular]  whose components are connected by a busbar)
3. busbar, [Meronyms (parts of bus), regular]  s a common connection between several circuits; "the busbar in this computer can transmit da [Meronyms (parts of bus), inherited]  nts of the system")
4. **bus**, jalo [Domain Terms]  le; "the fenders had fallen off that old bus")
Search for [Synonyms, ordered by estimated frequency ▼] of senses [    ]
☑ Show glosses
☐ Show contextual help
[Search]

---

Basic idea: Graphs induced by WordNet + domain labels for synsets + cooccurrence information from annotated corpora + collocations

$S_G \rightarrow S_s \mid S_g$
$S_s \rightarrow S_1 \mid S_2 \mid S_3$  (all the heuristics)
$S_1 \rightarrow E_1 S_1 \mid E_1$  (simple heuristics)
$E_1 \rightarrow e_{kind-of} \mid e_{part-of}$  (hyperonymy/meronymy)
$S_2 \rightarrow E_2 S_2 \mid E_2$  (hyponymy/holonymy)
$E_2 \rightarrow e_{has-kind} \mid e_{has-part}$  (parallelism)
$S_3 \rightarrow e_{kind-of} S_3 e_{has-kind} \mid e_{kind-of} e_{has-kind}$  (gloss)
$S_4 \rightarrow e_{gloss} S_4 \mid S_4$  (gloss, context)
$S_5 \rightarrow e_{gloss} e_{in-in-gloss}$  (gloss+gloss')

FigFig. 2. An excerpt of the context-free grammar for the recognition of semantic interconnections.

**(+ weights for each pattern and thus interconnection)**

Using SSI for word sense disambiguation ("The driver turned on his heel and went back to the truck.")

---

=> program, programme, computer program, computer programme -- ((computer science) a sequence of instructions that a compute
    HAS PART: routine, subroutine, subprogram, procedure, function -- ((a set sequence of steps, part of larger computer program)
    HAS PART: instruction, command, statement, program line -- ((computer science) a line of code written as part of a computer p
=> software, software system, software package, package -- ((computer science) written programs or procedures or rules and ass
    => code, computer code -- ((computer science) the symbolic arrangement of data or instructions in a computer program or the
        => coding system -- (a system of signals used to represent letters or numbers in transmitting messages)
            => writing -- (letters or symbols written or imprinted on a surface to represent the sounds or words of a language; "he tur
                => written communication, written language -- (communication by means of written symbols)
                    HAS PART: leaf, folio -- (a sheet of any written or printed material (especially in a manuscript or book))

Sense 5
driver, number one wood -- (a golf club (a wood) with a near vertical face that is used for hitting long shots from the tee)
    => wood -- (a golf club with a long shaft used to hit long shots; originally made with a wooden head; metal woods are now available)
        => golf club, golf-club, club -- (golf equipment used by a golfer to hit a golf ball)
            HAS PART: golf-club head, club head, club-head, clubhead -- ((golf) the head of the club which strikes the ball)
                HAS PART: face -- (the striking or working surface of an implement)
                HAS PART: **heel** -- ((golf) the part of the clubhead where it joins the shaft)
                HAS PART: sole -- (the underside of footwear or a golfclub)
                    HAS PART: shank, waist -- (the narrow part of the shoe connecting the heel and the wide part of the sole)
                HAS PART: toe -- ((golf) the part of a clubhead farthest from the shaft)
                => whole, whole thing, unit -- (an assemblage of parts that is regarded as a single entity; "how big is that part compi
                    HAS PART: part, portion -- (something less than the whole of a human artifact; "the rear part of the house"; "g
                    HAS PART: section, segment -- (one of several parts or pieces that fit with others to constitute a whole object;

Return to overview for driver
Return to WordNet home

---

3

## SSI for ontology learning

1. **Extract pertinent domain terminology:**
   - **Simple and multiword expressions that consistently occur in domain-related corpora and are not found in other domains (e.g., packet switching network)**
2. **Web search of available NL definitions from glossaries or documents**
3. **Use context-free grammar to**
   1. **filter out non-relevant definitions, based on statistical domain model**
   2. **parse definitions to extract kind-of information**
4. **Arrange terms in hierarchical trees**
5. **Link sub-hierarchies to the concepts of a core ontology (general-purpose: WordNet)**
6. **Provide the output to domain specialists for evaluation and refinement**

```
S → PP ',' NP SEP
NP → N1 KIND1
KIND1 → MOD1 NOUN1
MOD1 → Verb | Adj | Verb ',' MOD1 | Adj ',' MOD1
NOUN1 → Noun
N1 → Art | Adj
SEP → ',' | ';' | Prep | Verb | Wh
PP → Prep NP
```

## SSI: Word sense disambiguation and ontology learning

- *artifical language*: **monosemous in WordNet**
- *temporary or permanent termination*: **does not exist in WordNet**
- *termination*: **need to apply WSD (#1: end of a time span, #2: expiration of a contract)**
  → **use terms that occur in the subtree and have a lexical correspondent in WordNet**

$$P = \{disengagement, failure, block, transfer, dropout, temporary, permanent\}.$$

## Agenda

**Web Mining**

**(Semantic) Web**

## Agenda

**Web Mining**

**Semantic Web**

**Under-stand**

...
p3ee24304.dip.t-dialin.net
[19/Mar/2002:12:03:51 +0100
/search.html?t=jane%20austen
3785&ord=asc HTTP/1.0" 200
p3ee24304.dip.t-dialin.net
[19/Mar/2002:12:05:06 +0100
/search.html?t=jane%20austen&m=vide
o&SID=023785&ord=desc HTTP/1.0" 200
8450
p3ee24304.dip.t-dialin.net - -
[19/Mar/2002:12:06:41 +0100] "GET
/view.asp?id=3456&SID=023785
HTTP/1.0" 200 3478
...

## Application: Search in knowledge portals

## Ontology-based modelling of behaviour: URLs and application events

URL ⟶ Web page with content

Desired service      Obtained content

**Definition 24.3.1** *A site model* $M := (\mathcal{S}, \mathcal{C}[, \mathcal{D}_S][, \mathcal{D}_C])$ *is a tuple consisting of a service ontology* $\mathcal{S}$*, a content ontology* $\mathcal{C}$*, and optionally tuples* $\mathcal{D}_S$ *and/or* $\mathcal{D}_C$ *consisting of the service/content information dimensions.*

**Definition 24.3.2** *An atomic application event for a site model* $M = (\mathcal{S}, \mathcal{C}[, \mathcal{D}_S][, \mathcal{D}_C])$ *is a tuple* $AAE := (S, C)$ *consisting of (1) S, the service: a concept or relation from the service ontology* $\mathcal{S}$*, or (if and only if* $\mathcal{D}_S$ *is present in M) a tuple whose components are concepts/relations from S describing the service along the dimensions of* $\mathcal{D}_S$*, and (2) C, the content: a concept or relation from the content ontology* $\mathcal{C}$*, or (if and only if* $\mathcal{D}_C$ *is present in M) a tuple whose components are concepts/relations from C describing the service along the dimensions of* $\mathcal{D}_C$*.*

**Definition 24.3.3** *A complex application event CAE is a nonempty (1) set, (2) sequence (i.e., order relation), or (3) another sequential structure (a regular expression, a context-free grammar, ...), whose elements are atomic application events.*

Berendt, B., Stumme, G., & Hotho, A. (2004). Usage mining for and on the Semantic Web. In H. Kargupta, A. Joshi, K. Sivakumar, & Y. Yesha (Eds.), *Data Mining: Next Generation Challenges and Future Directions*. Menlo Park, CA: AAAI/MIT Press.

## Slide 25

### Semantics of requests
#### Step 1: Domain ontology

- **community portal**
  ka2portal.aifb.uni-karlsruhe.de

- **ontology-based:**
  - **Knowledge base** in F-Logic
  - **Static pages:** annotations
  - **Dynamic pages:** generated from queries & KB
  - **Queries** also in F-Logic
  - **Logs** contain these queries

affiliation

[Oberle, Berendt, Hotho, & Gonzalez, *Proc. AWIC* 2003]

---

## Slide 26

### Semantics of requests
#### Step 2: Modelling requests as atomic application events

RESEARCHER
PERSON
PROJECT
PUBLICATION
RESEARCHTOPIC
EVENT
ORGANIZATION
RESEARCHINTEREST
LASTNAME
TITLE
ISABOUT
EVENTS
EVENTTITLE
WORKSATPROJECT
AUTHOR
AFFILIATION
ISWORKEDONBY
PROGRAMCOMMITTEE
EMPLOYS
NAME
RESEARCHGROUPS
EMAIL

*An example query with concepts and relations:*

FORALL N,PEOPLE <-PEOPLE:
Employee[affiliation->> "http://www.anInstitute.org"]
and PEOPLE:Person[lastName->>N].

*Query =*
*feature vector of concepts + relations*

➔

**Clustering,**
**Association rules,**
**Classification, ...**

**Session =**
**feature vector of concepts + relations,**
**summed over all queries in the session**

---

## Slide 27

### Semantics of sequences
#### Step 3: Using ontologies of behaviour for info. Extraction – Modelling sequences as composite application events

**Composite application events - Example *customer typology***

- **Based on background theory from marketing: the customer buying cycle**
- **Modelled in terms of regular expressions and employed in Web usage mining**
- **Example:**
  *knowledge builders*

  (as opposed to, e.g.,
  *direct buyers*)

Home
Background information → Detail information
Contact

[Moe, *Journal of Consumer Psychology,* 2002]
[Spiliopoulou, Pohle, and Teltzrow, *Proc. Wirtschaftsinformatik 2002*]

---

## Slide 28

### 2. Semantics of sequences
#### for Step 3: an interactive tool with a query language

```
select t
from node a b, template a * b as t
where a.url startswith "SEITE1-"
and a.occurrence = 1
and b.url contains "1SCHULE"
and b.occurrence = 1
and (b.support / a.support) >= 0.2
```

Tool: www.hypknowsys.de; Data: [Berendt & Spiliopoulou, *VLDB Journal,* 2000]

---

## Slide 29

### Semantics of sequences
#### Step 4: Pattern discovery / instance learning

**An ontology of composite application events (CAEs)**

- **Define templates as regular expressions**
  - of atomic application events
  - of transitions (between atomic application events)

  *Ex.* [.search .* individual]

- **Discover instances by learning a CAE trie**

affiliationSearch, 629
topicSearch, 312
repetition, 402
refinement, 113
individual, 112
repetition, 295

[Berendt & Spiliopoulou, *VLDB Journal,* 2000]
[Berendt, *Data Mining and Knowledge Discovery,* 2002]

---

## Slide 30

### Semantics of sequences
#### Step 5: Pattern evaluation

**Use pattern statistics to**

- **derive descriptive measures of CAEs**
  - support, confidence
  - popularity, effectiveness, efficiency
- **apply inferential statistics to compare CAEs**

[Berendt, *Data Mining and Knowledge Discovery,* 2002]

Slide 31 — Communication – Visual data mining
Step 6: Mapping an ontological relation over concepts to a linear order; mapping to visual variables

Concreteness
Goal: Individual page — Reach goal
Search with more constraints — Refine search
First search page
Abandon search — Remain unspecific
Time



Slide 32 — Communication – Visual data mining
Step 6 – Example

Search criterion *location*
concreteness
Individual page 4
Search with $x$ parameters 3 2 1
Entry page
step1 step2 step3 step4 step5

Search criterion *textual property*
concreteness
4 3 2 1
step1 step2 step3 step4 step5

[Berendt, *Data Mining and Knowledge Discovery*, 2002], [Berendt, *Postproc. WebKDD 2001*]



Slide 33 — Communication – Visual data mining
Step 7: Visual abstraction ➔ new semantic patterns

Binary transitions
Cyclic behaviour
Patterns of leaving

Data:see [Berendt, Günther, & Spiekermann, *Communications of the ACM*, 2005]



Slide 34 — Step 8: Semantic Abstraction ➔ Detail & context

[Berendt, *Proc. WebKDD 2005*]



Slide 35 — Case study: Information search in a medical portal*

alphabetical search: hub-and-spoke → linguistic relations only (6.4%)

diagnoses serve as "hubs" for navigation (5.3%, 4%)

localisation search: linear / depth-first → search refinement & medical knowledge (5%)

* (20333 requests / 1397 sessions from Web log collected in 2001/2002; preprocessing: mining in concept space, see paper in proceedings)



Slide 36 — Agenda

Web Mining

(Semantic) Web

## Agenda

Web Mining

Semantic Web

...
<BIBLIOGRAPHY><FLOAT><PAGENUMBER>136</PAGENUMBER></FLOAT>
<HEAD>Literaturverzeichnis</HEAD>
<CITATION WORKTYPE="journal" PUBLISHED="PUBLISHED">
ID="bib-15-">[1]
</CUT><WORKAUTHOR>Aga... Krueger, B. P.; Scholes, G. D... M.; Yom, J.; Mets, L.; Fle... R.</WORKAUTHOR>U... ARTICLETITLE >ltrafast energy transfer in LHC-II revealed by three-pulse photon echo peak shift measurements</ARTICLETITLE>, <WORKTITLE>J. Phys. Chem. B</WORKTITLE>, <PUBDATE>2000</PUBDATE>, <NUMBER>104</NUMBER>, <PAGES>2908</PAGES>, </CITATION>
...

contribute

---

## Application: Knowledge construction for educational portals / Digital Libraries



---

## Knowledge contributions: Data and metadata

Aguirre-Arteta, Ana Maria: REGULATION OF DNA METHYLA... DEVELOPMENT: ALTERNATIVE ISOFORMS OF DNA METHY...

[Titelseite] [Widmung] [1] [2] [3] [4] [5] [Danksagung...
[Selbständigkeitserklärung] [Abkürzungsverz...

Aus dem Institut für Biologie der Humboldt-Universit...

DISSERTATION

zur Erlangung des akademischen Grades
doctor rerum naturalium
(Dr. rer. nat.)
im Fach BIOLOGIE

REGULATION OF DNA METHYLATION DURIN...
ALTERNATIVE ISOFORMS OF DNA METHY...

eingereicht an der
Mathematisch-Naturwissenschaftlichen Fakul...
der Humboldt-Universität zu Berlin

von Ana Maria Aguirre-Arteta,
geb. am 02.01.1968 in Bilbao (Spani...

<BIBLIOGRAPHY><FLOAT><PAGENUMBER>136</PAGENUMBER></FLOAT>
<HEAD>Literaturverzeichnis</HEAD>
...
<CITATION WORKTYPE="journal" PUBLISHED="PUBLISHED">
<CUT ID="bib-45-">[2]
</CUT><WORKAUTHOR>Albrecht, T. F.; Bott, K.; Meier, T.; Schulze, A.; Koch, M.; Cundiff, S. T.; Feldmann, J.; Stolz, W.; Thomas, P.; Koch, S. W.; Göbel; E. O.</WORKAUTHOR>
<ARTICLETITLE>Disorder mediated biexcitonic beats in semiconductor quantum wells</ARTICLETITLE>,
<WORKTITLE>Phys. Rev. B</WORKTITLE>,
<PUBDATE>1996</PUBDATE>,
<NUMBER>54</NUMBER>,
<PAGES>4436</PAGES>,
</CITATION> ...

---

## Dissertation Markup Language DiML
http://edoc.hu-berlin.de/diml/dtd/xdiml.dtd

...
<!ELEMENT citation (#PCDATA | email | url | note | workauthor | worktitle | articletitle | serialtitle | address | editor | publisher | edition | volume | number | version | pages | pubdate | bible | court | law | cut | pagenumber)*>
<!ATTLIST citation
 id ID #IMPLIED
 label CDATA #IMPLIED
 workType (Book | Journal | Misc) #IMPLIED
 published (yes|no) 'yes'>
<!ELEMENT note (#PCDATA | em | u | strong | br | sup | tt | sub | link | name | email | organization | term | foreign | url | footnote | endnote | glossref | indexref | pagenumber | q | citation | imath | im)*>
<!ATTLIST note
 id ID #IMPLIED>
<!ELEMENT workauthor (#PCDATA | given | surname | suffix | organization)*>
<!ATTLIST workauthor
 role CDATA #IMPLIED
 ref IDREF #IMPLIED
 id ID #IMPLIED>
...

---

## Authoring support for document servers

- **Surveys (ca. 2500 persons; 12-14% response rate) & Web usage mining (ca. 11000 sessions) showed:**
  - **Metadata creation is one of the main barriers for contribution.**
- **Reasons include deficiencies in**
  - information flow
  - understanding and use of structured search
  - education in structured writing
  - HCI aspects

➔ Marketing

) ➔ Education
)
)

**Intelligent Authoring Tools**

[Berendt, Brenstein, Li, & Wendland, *Proc. ETD* 2003]
[Berendt, *Proc. AAAI Spring Symposium KCVC*, 2005]

---

## Consequences of metadata neglect

<BIBLIOGRAPHY><FLOAT><PAGENUMBER>136</PAGENUMBER></FLOAT>
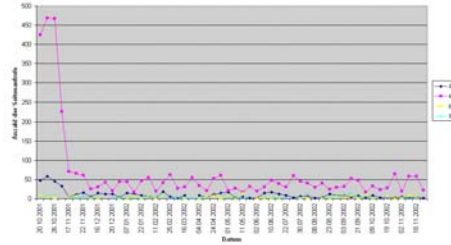
<HEAD>Literaturverzeichnis</HEAD>

<CITATION WORKTYPE="journal" PUBLISHED="PUBLISHED">

<CUT ID="bib-15-">[1] </CUT><WORKAUTHOR>Agarwal, R.; Krueger, B. P.; Scholes, G. D.; Yang, M.; Yom, J.; Mets, L.; Fleming, G. R.</WORKAUTHOR><ARTICLETITLE>ltrafast energy transfer in LHC-II revealed by three-pulse photon echo peak shift measurements</ARTICLETITLE>, <WORKTITLE>J. Phys. Chem. B</WORKTITLE>, <PUBDATE>2000</PUBDATE>, <NUMBER>104</NUMBER>, <PAGES>2908</PAGES>,

</CITATION>

...

## Why is this a problem?

[Cardona & Marx, *Physik Journal* 2004]



[Berendt, in *Neues Handbuch Hochschullehre,* 2003]

## System architecture



citeseer

paratools

TTT

Web service

VBA macro

other WS and info. sources

44

## Usage interface



corrected, XML annotated, and formatted

45

## Information extraction: Reference parsing with 3 tools

## Paratools-Zitations-Parsing
### http://paracite.eprints.org

**A database of templates of the form**

> '_AUTHORS_ (_YEAR_). _TITLE_.
> _PUBLICATION_, _VOLUME_(_ISSUE_):_PAGES_'

**each _XXX_ is associated with a regular expression**

- Ex.: _YEAR_ ➜ ([[:digit:]]{4})

**2 weighting factors**

- reliability: how "syntactically fixed" is a regular expression?
  - Ex.: _URL_ > _TITLE_
- concreteness = number of fixed symbols
  - Ex.: '_AUTHORS_, _PUBLICATION_, in press' > '_AUTHORS_, _PUBLICATION_'

**Templates are matched against the reference.**

**Choose the template with the highest *reliability*, or (if these are equal) with the highest *concreteness*.**

## Outlook 1: Diversity
### (or: Web space and real-life spaces)

## Slide 49 — Which diagnosis is that?

**Which diagnosis is that?**

Request frequency for a specific diagnosis in the investigated eHealth portal, depending on time and request language



[Yihune, 2003]

## Slide 50 — Hypotheses: search preferences

**Hypotheses: search preferences**

| Search option | Characteristics | Presumably preferred by |
|---|---|---|
| Search engine | ■ little context | ■ Low context |
| | | ■ Low Uncertainty Avoidance |
| Search | ■ fast information access | ■ Short-Term oriented |
| | ■ no hierarchies | ■ Low Power Distance |
| Alphabetically organized links | ■ large hierarchies | ■ High Power Distance |
| A   B  E   G | | |
| Content-organized links | ■ highest amount of (context) information | ■ High context |
| | | ■ High Uncertainty Avoidance |
| | ■ more time-consuming information access | ■ Long-Term oriented |
| | ■ large hierarchies | ■ High Power Distance |

(Kralisch & Berendt, *Proc. IWIPS 2004*)

## Slide 51 — Search behaviour: sample results

**Search behaviour: sample results**

1. **Which search options were used?**

UA – Uncertainty Avoidance
Cont – Context Specifity
LTO – Long-Term Orientation
PD – Power Distance

- Expected results
- Unexpected results
- all results significant (p<0.001)

**search engine:**



**content-organized links:**



## Slide 52 — Interactions between language and domain knowledge

**Interactions between language and domain knowledge**



Kralisch & Berendt, *New Review of Hypermedia and Multimedia,* in press)

## Slide 53 — Outlook 2: Community

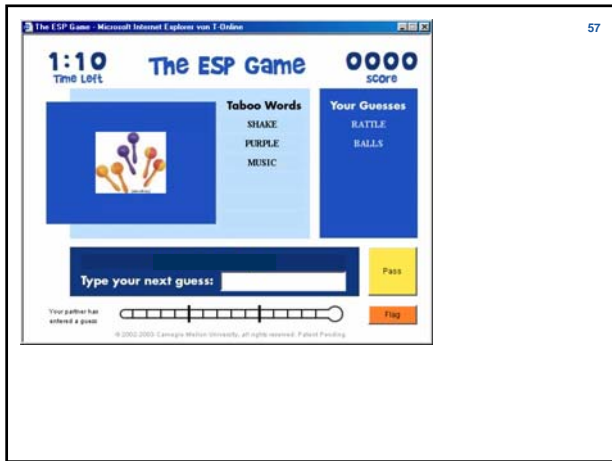**Outlook 2: Community**

## Slide 54 — bibster.semanticweb.org

**bibster.semanticweb.org**



**Definition 1** A user profile is a structure $PR := (E, Q, R, W, t)$ consisting of

- the expertise description $E$,
- a set of recent queries $Q$,
- a set of recent relevant instances $R$,
- a structure $W$ which defines the weights for the similarity function,

**Recommendations based on items' semantics and their**
**... similarity to the user's expertise ➜ measured by previous externalisations (content of personal database)**
**... similarity to relevant items ➜ measured by previous internalisations (answers to a query) and combinations (addition to the personal database)**

Haase, Ehrig, Hotho, & Schnizler, 2004

**www.bibserv.org**

**Outlook 3: Fun!**

The ESP Game - Microsoft Internet Explorer von T-Online

**1:10** Time Left   **The ESP Game**   **0000** score

**Taboo Words**   **Your Guesses**

SHAKE   RATTLE

PURPLE   BALLS

MUSIC

**Type your next guess:**   Pass

Your partner has entered a guess   Flag

© 2002-2003 Carnegie Mellon University, all rights reserved. Patent Pending.

The ESP Game - Microsoft Internet Explorer von T-Online

**1:10** Time Left   **The ESP Game**   **0000** score

**Taboo Words**   **Your Guesses**

The ESP Game - Microsoft Internet Explorer von T-Online

**The ESP Game**

TABOOS

MEN

AMERICA

http://www.perudo.com/

‹ Previous   Play Again   Next ›

© 2000-2003 Carnegie Mellon University, all rights reserved. Patent Pending.

The ESP Game - Microsoft Internet Explorer von T-Online

**The ESP Game**

TABOOS   AGREED ON

MEN   BLACK

AMERICA

http://www.perudo.com/

‹ Previous   Play Again   Next ›

© 2000-2003 Carnegie Mellon University, all rights reserved. Patent Pending.

**?
Outlook 4: Share the initiative –
automated Web service search
and composition for SWM
?**

**Thank you for your attention!**