*WebExplorer*
*A TOOL FOR ONTOLOGY-BASED INFORMATION EXPLORATION*

**Feza Baskaya, Anne Keskimaa,**
**Jaana Kekäläinen, Kalervo Järvelin**
***University of Tampere, Finland***

# *Background*

- Vast online information environments
  - billions of digital **documents**
  - many different natural **languages**
  - **distributed** document production and publication: no generally agreed rules
  - general **lack of control** in the process
  - much **spam** and other unwanted material

# *Background (2)*

- Old problem - vocabulary mismatch
  - hard to guess the best search keys; leads to loss of search effectiveness
  - especially in foreign languages
  - hard to know word forms (tokens / lemmas)
- Other problems – depending on one's search environment
  - collection dependency, metadata dependency
  - engine and query language dependency

TRIM
TAMPERE RESEARCH CENTER
FOR INFORMATION AND MEDIA

# *Background (3)*

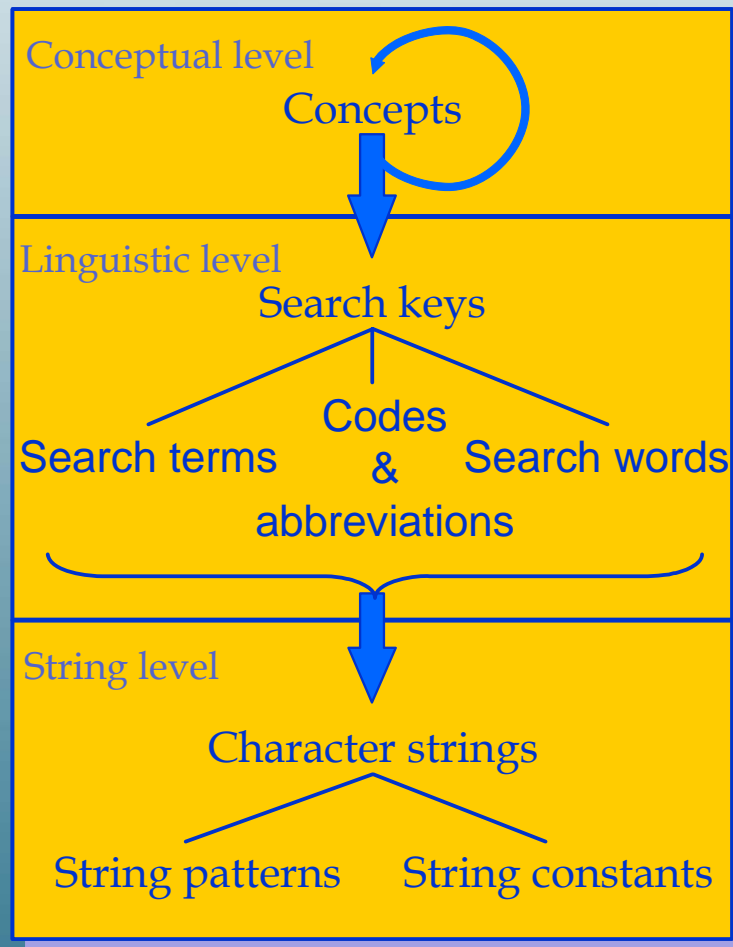to annotate

or

not to annotate

# *Design of WebExplorer*

Requirements on the Ontology

- –**Personal**
- –**Small-scale**
- –**Mapping**
- –**Multilingual**
- –**Editable**

## *Ontologies of WebExplorer based on*
# *Three-level architecture*

Conceptual level

Concepts

Linguistic level

Search keys

Codes
&
abbreviations

Search terms          Search words

String level

Character strings

String patterns     String constants

- Forest industry

  – forest industry
  – paper industry
  – saw mill
  – ...

  – pl(saw, mill)
  – al(industry)
  – pl(paperi, tehdas)

# *Design of WebExplorer*

Additional functionalities are required for information exploration:

- **Keyword search**
- **Query-biased summaries**
- **Classification**
- **Clustering**
- **Bookmarking**
- **Smooth integration of functions**

# WebExplorer Interface



The WebExplorer user interface

# WebExplorer Interface *(2)*

Search | Summary | Cluster | Classify | Bookmark | RSS

Query for Summary: phone [X] [change]

[Save Summary] [Show My Summaries]

1.0 the government's digital britain report has proposed that every citizen with a fixed line phone will pay 50p per month to pay for the roll out of faster next generation networks .

Search | Summary | Cluster | Classify | Bookmark | RSS

Select a cluster option and click on cluster tab after selection in order to get new cluster tree.

⊙ All documents on the first page (titel,snipped)
○ Selected documents (content)

Euroopan Parlamentti
Pahin
(Other)
  EP
  Murrosten
  Buzek Tapasi Tukholmassa
  Analysoi
  Afrikan
  Kiinan Hallituksen
  EU

TRIM
TAMPERE RESEARCH CENTER
FOR INFORMATION AND MEDIA

# WebExplorer Interface *(3)*

# WebExplorer System Architecture



**Client-JavaScript**

search

c
→ c
  → c
    → c
    → c
    → c
    → c
→ c
  → c
  → c

Results
1. snippet
2. snippet
...
n. snippet

**Server WebExplorer**

| Search Servlet | Expansion Engine |
| Summarizer Servlet | Summarizer Engine |
| Classifier Servlet | Classifier Engine |
| Clustering Servlet | Clustering Engine |
| Bookmarking Servlet | RSS Servlet |
| Ontology Servlet | Info Servlet |

**Document servers**

Document DB TRIP

Document DB Lemur

**Web** Google

**Postgress RDBMS**

Ontology KB

WebExplorer  System Architecture Diagram

# *Formative Evaluation*

## *To test the feasibility of WebExplorer*

|  | Group 1 | Group 2 |
|---|---|---|
| Ontology IR | +2 | +1 |
| CL-IR | +2 | +1.5 |
| Summarization | 0 | +1 |
| Classification | +0.5 | -- |
| Bookmarking | +1.5 | 0.5 |
| RSS Feeds | +0.5 | -- |

Scale: -2 to +2

-2:No benefit at all

-1:Only little benefit

 0:I can't say

+1:some benefit

+2:much benefit

School grade (4-10) for WebExplorer:7

Median values of the groups' opinions

TRIM
TAMPERE RESEARCH CENTER
FOR INFORMATION AND MEDIA

# *Conclusion*

- WebExplorer is an answer to problems in semantic information access
  - **light-weight** disposable **ontologies** for direct content access
  - **independencies** of:
    - **collections** (partially), indexing ways,
    - availability of metadata / **annotations**
    - changes of needs, **variability** of "world models"
    - **search engines**, **query languages**
    - vocabulary variation and **natural languages**
  - annotation - media analysis
  - a compromise different from semantic annotation or library indexing with control at the user end
    - costs and benefits distributed in a new way

# *Conclusion*

Document Analysis functions offer

- – Online query biased **summarization** of individual documents and document sets

- – **Clustering** of documents

- – Ontology based cross-lingual **classification**

- – **Bookmark** managemet

- – Displaying search results on **map** and **timeline**