# Interactive information retrieval in XML documents

## Nils Pharo, Oslo University College

# Outline

- XML retrieval
- INEX
- INEX interactive track
- Recent study:
  - research questions
  - results
  - conclusion

# XML retrieval

```
<?xml version="1.0" encoding="UTF-8 "?>
<book class="H.3.3">
<author>John Smith</author>
<title>XML Retrieval</title>
<chapter>
<heading>Introduction</heading> This text explains all
    about XML and IR. </chapter>
<chapter>
<heading> XML Query Language XQL</heading>
<section> <heading>Examples</heading> </section>
<section> <heading>Syntax</heading> Now we describe
    the XQL syntax. </section>
</chapter>
</book>
```

# XML-markuped documents

Facilitates

- possibilities to exploit weighting
- the retrieval of document parts

# INEX – the Initiative for the Evaluation of XML Retrieval

- International community whose main goal is "to promote the evaluation of focused retrieval by providing large test collections of structured documents, uniform evaluation measures, and a forum for organizations to compare their results"
- Since 2002
- Partly financed by DELOS $\rightarrow$ 2007
- http://www.inex.otago.ac.nz/

# INEX organization

- Test collection
  - Journal articles from IEEE ($\rightarrow$2005)
  - Wikipedia articles (2006$\rightarrow$)
  - Topics created by participants
  - Document assessed by topic creators
- 186 participating institutions (in 2008)

# INEX tracks

- Ad Hoc
- Book
- Efficiency
- Entity Ranking
- *Interactive (iTrack)*
- Question Answering
- Link-the-Wiki
- XML-Mining

# INEX Interactive track

- investigate end-users interaction with elements of XML documents
- develop approaches for element retrieval which are effective in user-based environments

# Organisation of INEX iTrack

Distributed setup:

- a common subject recruiting procedure

- a common set of user tasks and data collection instruments such as interview guides and questionnaires

- a common logging procedure for user/system interaction

# 2008: two task categories

**Fact finding**: "In the recent Olympics there was a controversy over the age of some of the female gymnasts. You want to know the minimum age for Olympic competitors in gymnastics."

**Research**: "Every year there are several ranking lists over the best universities in the world. These lists are seldom similar. You are writing an article discussing and comparing the different ranking systems and need information about the different lists and what criteria and factors they use in their ranking."

# Method

- The searchers perform simulated tasks

- Documents and elements can be relevance judged

- Questionnaires to collect background data such as; demographic data, search experience, topic knowledge

# System

# Document view

# Example study

- Goal: identify searchers preferred level of granularity

- Compared two genre: journal articles (2005 data) and encyclopedic articles (2006 data)

- Log analysis using Excel

# Result list interaction

| Granularity level | IEEE | Wikipedia |
|---|---|---|
| Level 1/2  (article or metadata) | 71 % | 73.5 % |
| Level 3 (section) | 14.5 % | 15 % |
| Level 4/5  (subsections) | 14.5 % | 11.5 % |
| Total | 100 % | 100 |

# Within documents retrieval

| Granularity level | IEEE | Wikipedia |
|---|---|---|
| Level 1 (article) | 8 % | 9 % |
| Level 2 (metadata) | 6.5 % | |
| Level 3 (section) | 55.5 % | 49 % |
| Level 4/5 (subsections) | 30 % | 42 % |
| Total | 100 % | 100 % |

# Relevance assessments

| Granularity level | IEEE | Wikipedia |
|---|---|---|
| Level 1 (article) | 7 % | 40 % |
| Level 2 | 14 % | |
| Level 3 | 50 % | 33.5 % |
| Level 4/5 | 29 % | 26.5 % |
| Total | 100 % | |

# Two possible reasons for difference in relevance assessments

- The Size factor: Wikipedia articles are smaller than scientific articles
- The Genre factor: the structure of scientific articles and encyclopaedic articles differ

# Conclusions

- searchers prefer to use article as initial entry point
- within documents searchers look at elements relatively proportional to their distribution
- genre has a strong influence on what items are assessed for relevance