Norbert Fuhr
Mounia Lalmas
Saadia Malik
Zoltán Szlávik (Eds.)

# Advances in XML Information Retrieval

Third International Workshop of the Initiative
for the Evaluation of XML Retrieval, INEX 2004
Dagstuhl Castle, Germany, December 2004
Revised Selected Papers

Norbert Fuhr   Mounia Lalmas
Saadia Malik   Zoltán Szlávik (Eds.)

# Advances in XML Information Retrieval

Third International Workshop of the Initiative
for the Evaluation of XML Retrieval, INEX 2004
Dagstuhl Castle, Germany, December 6-8, 2004
Revised Selected Papers

Volume Editors

Norbert Fuhr
University of Duisburg-Essen
Faculty of Engineering Sciences, Information Systems
47048 Duisburg, Germany
E-mail: fuhr@uni-duisburg.de

Mounia Lalmas
Queen Mary University of London, Department of Computer Science
London E1 4NS, England, United Kingdom
E-mail: mounia@dcs.qmul.ac.uk

Saadia Malik
University of Duisburg-Essen
Faculty of Engineering Sciences, Information Systems
47048 Duisburg, Germany
E-mail: malik@is.informatik.uni-duisburg.de

Zoltán Szlávik
Queen Mary University of London, Department of Computer Science
London E1 4NS, England, United Kingdom
E-mail: zolley@dcs.qmul.ac.uk

# Lecture Notes in Computer Science 3493

*Commenced Publication in 1973*
Founding and Former Series Editors:
Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

# Preface

The ultimate goal of many information access systems (e.g., digital libraries, the Web, intranets) is to provide the right content to their end-users. This content is increasingly a mixture of text, multimedia, and metadata, and is formatted according to the adopted –W3C standard for information repositories, the so-called eXtensible Markup Language (XML). Whereas many of today's information access systems still treat documents as single large (text) blocks, XML offers the opportunity to exploit the internal structure of documents in order to allow for more precise access thus providing more specific answers to user requests. Providing effective access to XML-based content is therefore a key issue for the success of these systems.

The aim of the INEX campaign (Initiative for the Evaluation of XML Retrieval), which was set up at the beginning of 2002, is to establish infrastructures, XML test suites, and appropriate measurements for evaluating the performance of information retrieval systems that aim at giving effective access to XML content. More precisely, the goal of the INEX initiative is to provide means, in the form of a large XML test collection and appropriate scoring methods, for the evaluation of content-oriented XML retrieval systems.

INEX 2004 was responsible for a range of evaluation activities in the field of XML information retrieval, with five tracks: (1) *Ad Hoc Retrieval Track*, the main track, which can be regarded as a simulation of how a digital library might be used, where a static set of XML documents and their components is searched using a new set of queries (topics) containing both content and structural conditions; (2) *Interactive Track*, which aimed to investigate the behavior of users when interacting with components of XML documents; (3) *Heterogeneous Collection Track*, where retrieval is based on a collection comprising various XML subcollections from different digital libraries, as well as material from other resources; (4) *Relevance Feedback Track*, dealing with relevance feedback methods for XML; and (5) *Natural Language Track*, where natural language formulations of structural conditions of queries have to be answered.

The INEX 2004 workshop, held at Schloss Dagstuhl (Germany), 6–8 December 2004, brought together researchers in the field of XML retrieval who participated in the INEX 2004 evaluation campaign. Participants were able to present and discuss their approaches to XML retrieval. These proceedings contain revised papers describing work carried out during INEX 2004 in the various tracks by the participants.

assessment tool), and Gabriella Kazai and Arjen de Vries (metrics). The organizers of the various tracks did a great job and their work is greatly appreciated: Anastasios Tombros, Birger Larsen, Thomas Rölleke, Carolyn Crouch, Shlomo Geva and Tony Sahama. Finally, we would like to thank the participating organizations and people for their participation in INEX 2004.

March 2005                                                          Norbert Fuhr
                                                                    Mounia Lalmas
                                                                    Saadia Malik
                                                                    Zoltán Szlávik

# Organizers

## Organizers

### Project Leaders

Norbert Fuhr, University of Duisburg-Essen
Mounia Lalmas, Queen Mary University of London

### Contact Person

Saadia Malik, University of Duisburg-Essen

### Topic Format Specification

Börkur Sigurbjörnsson, University of Amsterdam
Andrew Trotman, University of Otago

### Online Relevance Assessment Tool

Benjamin Piwowarski, University of Chile

### Metrics

Gabriella Kazai, Queen Mary University of London
Arjen P. de Vries, Centre for Mathematics and Computer Science

### Interactive Track

Birger Larsen, Royal School of Library and Information Science
Saadia Malik, University of Duisburg-Essen
Anastasios Tombros, Queen Mary University of London

### Relevance Feedback Track

Carolyn Crouch, University of Minnesota-Duluth
Mounia Lalmas, Queen Mary University of London

### Heterogeneous Collection Track

Thomas Rölleke, Queen Mary University of London
Zoltán Szlávik, Queen Mary University of London

### Natural Language Processing

Shlomo Geva, Queensland University of Technology
Tony Sahama, Queensland University of Technology

# Lecture Notes in Computer Science

For information about Vols. 1–3411

please contact your bookseller or Springer

# Table of Contents

**Ad Hoc Retrieval and Relevance Feedback**

**Relevance Feedback**

## Ad Hoc Retrieval and Heterogeneous Document Collection

## Heterogeneous Document Collection

## Natural Language Processing of Topics

## Interactive Studies

# Overview of INEX 2004

Saadia Malik[1], Mounia Lalmas[2], and Norbert Fuhr[3]

[1] Information Systems, University of Duisburg-Essen, Duisburg, Germany
`malik@is.informatik.uni-duisburg.de`
[2] Department of Computer Science, Queen Mary University of London, London, UK
`mounia@dcs.qmul.ac.uk`
[3] Information Systems, University of Duisburg-Essen, Duisburg, Germany
`fuhr@uni-duisburg.de`

## 1 Introduction

The widespread use of the eXtensible Markup Language (XML) in scientific data repositories, digital libraries and on the web, brought about an explosion in the development of XML retrieval systems. These systems exploit the logical structure of documents, which is explicitly represented by the XML markup: instead of whole documents, only components thereof (the so-called XML elements) are retrieved in response to a user query. This means that an XML retrieval system needs not only to find relevant information in the XML documents, but also determine the appropriate level of granularity to return to the user, and this with respect to both content and structural conditions.

Evaluating the effectiveness of XML retrieval systems requires a test collection (XML documents, tasks/topics, and relevance judgements) where the relevance assessments are provided according to a relevance criterion that takes into account the imposed structural aspects. A test collection as such has been built as a result of three rounds of the Initiative for the Evaluation of XML Retrieval[1] (INEX 2002, INEX 2003 and INEX 2004). The aim of this initiative is to provide means, in the form of large testbeds and appropriate scoring methods, for the evaluation of content-oriented retrieval of XML documents.

This paper presents an overview of INEX 2004. In section 2, we give a brief summary of the participants. Section 3 provides an overview of the test collection along with the description of how the collection was constructed. Section 4 outlines the retrieval tasks in the main track, which is concerned with the ad hoc retrieval of XML documents. Section 5 briefly reports on the submission runs for the retrieval tasks, and Section 6 describes the relevance assessment phase. The different metrics used are discussed in Section 7, followed by a summary of the evaluation results in Section 8. Section 9 presents a short description of four new tracks that started in INEX 2004, namely the heterogenous collection track, the relevance feedback track, the natural language processing track and the interactive track. The paper finishes with some conclusions and an outlook for INEX 2005.

---

[1] http://inex.is.informatik.uni-duisburg.de/

## 2     Participating Organisations

In response to the call for participation issued in March 2004, around 55 organisations registered from 20 different countries within six weeks. Throughout the year, the number of participants decreased due to insufficient contribution while a number of new groups joined later at the assessment phase. The active participants are listed in Table 1.

## 3     The Test Collection

The INEX test collection, as for any IR test collection aiming at evaluating retrieval effectiveness, is composed of three parts: the set of documents, the set of topics, and the relevance assessments (these are described in Section 6).

### 3.1     Documents

The document collection was donated by the IEEE Computer Society. It consists of the full-text of 12,107 articles, marked up in XML, from 12 magazines and 6 transactions of the IEEE Computer Society's publications, covering the period of 1995-2002, and totalling 494 MB in size, and 8 millions in number of elements. The collection contains scientific articles of varying length. On average, an article contains 1,532 XML nodes, where the average depth of the node is 6.9. More details can be found in [3].

### 3.2     Topics

As in previous years, in INEX 2004 we distinguish two types of topics, reflecting two user profiles, where the users differ in the amount of knowledge they have about the structure of the collection:

- – **Content-only (CO) topics** are requests that ignore the document structure and contain only content related conditions, e.g. only specify what a document/component should be about (without specifying what that component is). The CO topics simulate users who do not (want to) know, or do not want to use, the actual structure of the XML documents. This profile is likely to fit most users searching XML digital libraries.
- – **Content-and-structure (CAS) topics** are topic statements that contain explicit references to the XML structure, and explicitly specify the contexts of the user's interest (e.g. target elements) and/or the contexts of certain search concepts (e.g. containment conditions). The CAS topics simulate users that have some knowledge of the structure of the XML. Those users might want to use this knowledge to try to make their topics more concrete, by adding structural constraints. This user profile could fit librarians that have some knowledge of the collection structure.

The topic format and guidelines were based on TREC guidelines, but were modified to accommodate the two types of topics used in INEX. Both CO and CAS topics are made up of four parts. The parts explain the same information need, but for different purposes.

Table 1. List of active INEX 2004 participants

| Organisations | Assessed topics | no of runs submitted |
|---|---|---|
| University of Amsterdam | 2 | 6 |
| University of California, Berkeley | 2 | 5 |
| VSB-Technical University of Ostrava | 2 | 0 |
| RMIT University | 1 | 6 |
| University of Otago | 2 | 1 |
| IBM Haifa Research Lab | 2 | 6 |
| University of Illinois at Urbana-Champaign | 2 | 0 |
| Nara Institute of Science and Technology | 2 | 4 |
| University of Wollongong in Dubai | 2 | 0 |
| Fondazione Ugo Bordoni | 2 | 3 |
| IRIT | 2 | 6 |
| Ecoles des Mines de Saint-Etienne, France | 1 | 2 |
| University of Munich (LMU) | 2 | 5 |
| Queen Mary University of London | 2 | 0 |
| Royal School of LIS | 1 | 3 |
| LIP6 | 1 | 6 |
| University of Tampere | 2 | 6 |
| University of Helsinki | 2 | 3 |
| Carnegie Mellon University | 1 | 6 |
| Cirquid project (CWI and University of Twente) | 2 | 6 |
| The Selim and Rachel Benin School of Engineering and Computer Science | 2 | 0 |
| University of Minnesota Duluth | 1 | 2 |
| Bamberg University | 2 | 0 |
| UCLA | 2 | 6 |
| Max-Planck-Institut fuer Informatik | 4 | 4 |
| Kyungpook National University | 2 | 4 |
| Utrecht University | 1 | 6 |
| The Robert Gordon University | 0 | 2 |
| University of Milano | 2 | 0 |
| Oslo University College | 2 | 6 |
| Cornell University | 2 | 0 |
| Universität Rostock | 2 | 0 |
| Universidade Estadual de Montos Claros | 2 | 6 |
| INRIA | 2 | 0 |
| LIMSI/CNRS | 2 | 0 |
| University of Waterloo | 2 | 3 |
| Queensland University of Technology | 2 | 6 |
| Indiana University | 2 | 2 |
| University of Granada | 2 | 0 |
| University of Kaiserslautern | 2 | 0 |
| The University of Iowa | 2 | 0 |
| Rutgers University | 2 | 0 |
| IIT Information Retrieval Lab | 2 | 0 |

- **Title:** a short explanation of the information need. It serves as a summary of both the content and, in the case of CAS topics, also the structural requirements of the user's information need. For the expression of these constraints the Narrowed Extended XPath I (NEXI) query syntax is used [8].
- **Description:** a one or two sentence natural language definition of the information need.
- **Narrative:** a detailed explanation of the information need and the description of what makes a document/component relevant or not. The narrative was there to explain not only what information is being sought for, but also the context and motivation of the information need, i.e., why the information is being sought and what work task it might help to solve. The latter was required for the interactive track (see Section 9.4).
- **Keywords:** a set of comma-separated scan terms that were used in the collection exploration phase of the topic development process (see later) to retrieve relevant documents/ components. Scan terms may be single words or phrases and may include synonyms, and terms that are broader or narrower terms than those listed in the topic description or title.

The title and the description must be interchangeable, which was required for the natural language processing track (see Section 9.3). The DTD of the topics is shown in Figure 1.

```
<!ELEMENT inex_topic  (title,description,narrative,keywords)>
<!ATTLIST inex_topic
  topic_id   CDATA  #REQUIRED
  query_type CDATA  #REQUIRED
  ct_no      CDATA  #REQUIRED
>
<!ELEMENT title (#PCDATA)>
<!ELEMENT description  (#PCDATA)>
<!ELEMENT narrative    (#PCDATA)>
<!ELEMENT keywords     (#PCDATA)>
```

**Fig. 1.** Topic DTD

The attributes of the topic are: topic_id (which ranges from 127 to 201), query_type (with value CAS or CO) and ct_no, which refers to the candidate topic number (which ranges from 1 to 198). Examples of both types of topics can be seen in Figure 2 and Figure 3.

The topics were created by participating groups. Each participant was asked to submit up to 6 candidate topics (3 CO and 3 CAS). A detailed guideline was provided to the participants for the topic creation. Four steps were identified for this process: 1) Initial Topic Statement creation 2) Collection Exploration 3) Topic Refinement and 4) Topic Selection. The first three steps were performed by the participants themselves while the selection of topics was decided by the organisers.

During the first step, participants created their initial topic statement. These were treated as a user's description of his/her information need and were formed without

```
<inex_topic topic_id="127" query_type="CAS" ct_no="13">
<title>//sec//(p| fgc)[about( ., Godel Lukasiewicz and other
fuzzy implication definitions)]</title>
<description>Find paragraphs or figure-captions containing the
definition of Godel, Lukasiewicz or other fuzzy-logic
implications</description>
<narrative>Any relevant element of a section must contain the
definition of a fuzzy-logic implication operator or a pointer to
the element of the same article where the definition can be
found. Elements containing criteria for identifying or comparing
fuzzy implications are also of interest. Elements which discuss
or introduce non-implication fuzzy operators are not relevant.
</narrative>
<keywords>Godel implication, Lukasiewicz implication, fuzzy
implications, fuzzy-logic implication </keywords>
</inex_topic>
```

**Fig. 2.** A CAS topic from the INEX 2004 test collection

```
<inex_topic topic_id="162" query_type="CO" ct_no="1">
<title> Text and Index Compression Algorithms </title>
<description>Any type of coding algorithm for text and index
compression</description>
<narrative>We have developed an information retrieval system
implementing compression techniques for indexing documents. We
are interested in improving the compression rate of the system
preserving a fast access and decoding of the data. A relevant
document/component should introduce new algorithms or compares
the performance of existing text-coding techniques for text and
index compression. A document/component discussing the cost of
text compression for text coding and decoding is highly relevant.
Strategies for dictionary compression are not
relevant.</narrative>
<keywords>text compression, text coding, index compression
algorithm</keywords>
</inex_topic>
```

**Fig. 3.** A CO topic from the INEX 2004 test collection

regard to system capabilities or collection peculiarities to avoid artificial or collection biased queries. During the collection exploration phase, participants estimated the number of relevant documents/components to their candidate topics. The HyREX retrieval system [1] was provided to participants to perform this task. Participants had to judge the top retrieved results and were asked to record the relevant document/component (XPath) paths in the top 25 retrieved components/documents. Those topics having at least 2 relevant documents/components and less than 20 documents/components in the top 25 retrieved elements could be submitted as candidate topics. In the topic refinement stage, the topics were finalised ensuring coherency and that each part of the topic could be used in stand-alone fashion.