

# INEX 2002 - 2006: Understanding XML Retrieval Evaluation

Mounia Lalmas and Anastasios Tombros

Queen Mary University of London,  
Mile End Road, London, UK  
{mounia, tassos}@dcs.qmul.ac.uk

**Abstract.** Evaluating the effectiveness of XML retrieval requires building test collections where the evaluation paradigms are provided according to criteria that take into account structural aspects. The Initiative for the Evaluation of XML retrieval (INEX) was set up in 2002, and aimed to establish an infrastructure and to provide means, in the form of large test collections and appropriate scoring methods, for evaluating the effectiveness of content-oriented XML retrieval. This paper describes the evaluation methodology developed in INEX, with particular focus on how evaluation metrics and the notion of relevance are treated.

## 1 Introduction

The continuous growth in XML information repositories has been matched by increasing efforts in the development of XML retrieval systems, in large part aiming at supporting content-oriented XML retrieval. These systems exploit the available structural information in documents, as marked up in XML, in order to implement a more *focussed* retrieval strategy and return document components – the so-called *XML elements* – instead of complete documents in response to a user query. This focussed retrieval approach is of particular benefit for information repositories containing long documents, or documents covering a wide variety of topics (e.g. books, user manuals, legal documents), where users' effort to locate relevant content can be reduced by directing them to the most relevant parts of these documents. As the number of XML retrieval systems increases, so does the need to evaluate their effectiveness.

The predominant approach to evaluate system retrieval effectiveness is with the use of test collections constructed specifically for that purpose. A test collection usually consists of a set of documents, user requests usually referred to as topics, and relevance assessments which specify the set of "right answers" for the user requests. Traditional IR test collections and methodology, however, cannot directly be applied to the evaluation of content-oriented XML retrieval as they do not consider structure. This is because they focus mainly on the evaluation of IR systems that treat documents as independent and well-distinguishable separate units of approximately equal size. Since content-oriented XML retrieval allows for document components to be retrieved, multiple elements from the same document can hardly be viewed as independent units. When allowing for the retrieval of arbitrary elements, we must also consider the overlap of elements; e.g. retrieving a complete section consisting of several paragraphs as

one element and then a paragraph within the section as a second element. This means that retrieved elements cannot always be regarded as separate units. Finally, the size of the retrieved elements should be considered, especially due to the task definition; e.g. retrieve minimum or maximum units answering the query, retrieve a component from which we can access, or browse to, a maximum number of units answering the query.

The evaluation of XML retrieval systems thus makes it necessary to build test collections where the evaluation paradigms are provided according to criteria that take into account the imposed structural aspects. The INitiative for the Evaluation of XML retrieval (INEX)<sup>1</sup>, which was set up in 2002, established an infrastructure and provided means, in the form of large test collections and appropriate scoring methods, for evaluating how effective content-oriented XML search systems are. This paper provides a detailed overview of the evaluation methodology developed in INEX, with particular focus on the treatment of the notion of relevance and on metrics for the evaluation of retrieval effectiveness.

## 2 The INEX Test-Beds

In traditional IR test collections, documents are considered as units of unstructured text, queries are generally treated as bags of terms or phrases, and relevance assessments provide judgments whether a document as a whole is relevant to a query or not. Although a test collection for XML IR consists of the same parts, each component is rather different from its traditional IR counterpart. XML documents organise their content into smaller, nested structural elements. Each of these elements in the document's hierarchy, along with the document itself (the root of the hierarchy), represent a retrievable unit. In addition, with the use of XML query languages, users of an XML IR system can express their information need as a combination of content and structural conditions. Consequently, relevance assessments for an XML collection must also consider the structural nature of the documents and provide assessments at different levels of the document hierarchy.

### 2.1 Document Collections

Up to 2004, the collection consisted 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, totalling 494 MB in size and 8 million in number of elements. On average, an article contains 1,532 XML nodes, where the average depth of the node is 6.9. In 2005, the collection was extended with further publications from the IEEE. A total of 4,712 new articles from the period of 2002-2004 were added, giving a total of 16,819 articles, and totalling 764MB in size and 11 million in number of elements.

INEX 2006 uses a different document collection, made from English documents from Wikipedia<sup>2</sup> [2]. The collection consists of the full-texts, marked-up in XML, of 659,388 articles of the Wikipedia project, and totaling more than 60 GB (4.6 GB without

---

<sup>1</sup> <http://inex.is.informatik.uni-duisburg.de/>

<sup>2</sup> <http://en.wikipedia.org>

images) and 52 million in number of elements. The collection has a structure similar to the IEEE collection. On average, an article contains 161.35 XML nodes, where the average depth of an element is 6.72.

## 2.2 Topics

Querying XML documents can be with respect to content and structure. Taking this into account, INEX identified two types of topics:

- *Content-only (CO)* topics are requests that ignore the document structure and are, in a sense, the traditional topics used in IR test collections. In XML retrieval, the retrieval results to such topics can be elements of various complexity, e.g. at different levels of the XML documents' structure.
- *Content-and-structure (CAS)* topics are requests that contain conditions referring both to content and structure of the sought elements. These conditions may refer to the content of specific elements (e.g. the elements to be returned must contain a section about a particular topic), or may specify the type of the requested answer elements (e.g. sections should be retrieved).

CO and CAS topics reflect two types of users with varying levels of knowledge about the structure of the searched collection. The first type simulates users who either do not have any knowledge of the document structure or who choose not to use such knowledge. This profile is likely to fit most users searching XML repositories. The second type of users aims to make use of any insight about the document structure that they may possess. CAS topics simulate users who do have some knowledge of the structure of the searched collection. They may then use this knowledge as a precision enhancing device in trying to make the information need more concrete. This user type is more likely to fit, e.g., librarians.

As in TREC, an INEX topic consists of the standard title, description and narrative fields. For CO topics, the title is a sequence of terms. For CAS topics, the title is expressed using the NEXI query language, which is a variant of XPATH defined for content-oriented XML retrieval evaluation - it is more focussed on querying content than many of the XML query languages [9].

In 2005, in an effort to investigate the usefulness of structural constraints, variants of the CO and CAS topics were developed. CO topics were extended into Content-Only + Structure (CO+S) topics. The aim was to enable the performance comparison of an XML system across two retrieval scenarios on the same topic, one when structural constraints are taken into account (+S) and the other when these are ignored (CO). The CO+S topics included an optional field called CAS title (<castitle>), which was a representation of the same information need contained in the <title> field of a CO topic but including additional knowledge in the form of structural constraint. CAS titles were expressed in the NEXI query language.

How to interpret the structural constraints (whether as part of CAS or CO+S topics) evolved over the years, since each structural constraint could be considered as a strict (must be matched exactly) or vague (does not need to be matched exactly) criterion. In the latter case, structural constraints were to be viewed as hints as to where to look

for relevant information. In 2002, the structural constraints of CAS topics were strictly interpreted. In 2003, both interpretations, strict and vague, were followed, whereas since 2004 only the latter was followed. As of today, INEX has a total of 401 topics.

### 2.3 Retrieval Tasks

The main INEX activity is the ad-hoc retrieval task. In IR literature, ad-hoc retrieval is described as a simulation of how a library might be used and involves the searching of a static set of documents using a new set of topics. Here, the collection consists of XML documents, composed of different granularity of nested XML elements, each of which represents a possible unit of retrieval. The user's query may also contain structural constraints, or hints, in addition to the content conditions.

A major departure from traditional IR is that XML retrieval systems need not only score elements with respect to their relevance to a query, but also determine the appropriate level of element granularity to return to users. In INEX, a relevant element is defined to be at the *right level of granularity* if it discusses all the topics requested in the user query – it is *exhaustive* to the query – and does not discuss other topics – it is *specific* to that query.

Up to 2004, ad-hoc retrieval was defined as the *general* task of returning, instead of whole documents, those XML elements that are most specific and exhaustive to the user's query. In other words, systems should return components that contain as much relevant information and as little irrelevant information as possible. Within this general task, several sub-tasks were defined, where the main difference was the treatment of the structural constraints.

The *CO sub-task* makes use of the CO topics, where an effective system is one that retrieves the most specific elements, and only those which are relevant to the topic of request. The *CAS sub-task* makes use of CAS topics, where an effective system is one that retrieves the most specific document components, which are relevant to the topic of request and match, either strictly or vaguely, the structural constraints specified in the query. In 2002, a strict interpretation of the CAS structural constraints was adopted, whereas in 2003, both, a strict and a vague interpretation was followed, leading to the *SCAS sub-task* (strict content-and-structure), defined as for the INEX 2002 CAS sub-task, and the *VCAS sub-task* (vague content-and-structure). In that last sub-task, the goal of an XML retrieval system was to return relevant elements that may not exactly conform to the structural conditions expressed within the user's query, but where the path specifications should be considered hints as to where to look. In 2004, the two sub-tasks investigated were the CO sub-task, and the VCAS sub-task. The SCAS sub-task was felt to be an unrealistic task because specifying an information need is not an easy task, in particular for semi-structured data with a wide variety of tag names.

However, within this general task, the actual relationship between retrieved elements was not considered, and many systems returned overlapping elements (e.g. nested elements). Indeed, the top 10 ranked systems for the CO sub-task in INEX 2004 contained between 70% to 80% overlapping elements. This had very strong implications with respect to measuring effectiveness (Section 3), where approaches that attempted to implement a more focussed approach (e.g., between two nested relevant elements, return the one most specific to the query) performed poorly. As a result, the *focussed sub-task*

was defined in 2005, intended for approaches aiming at targeting the appropriate level of granularity of relevant content that should be returned to the user for a given topic. The aim was for systems to find the most exhaustive and specific element on a path within a given document containing relevant information and return to the user only this most appropriate unit of retrieval. Returning overlapping elements was not permitted. The INEX ad-hoc general task, as carried out by most systems up to 2004, was renamed in 2005 as the *thorough sub-task*.

Within all the above sub-tasks, the output of XML retrieval systems was assumed to be a ranked list of XML elements, ordered by their presumed relevance to the query, whether overlapping elements were allowed or not. However, user studies [8] suggested that users were expecting to be returned elements grouped per document, and to have access to the overall context of an element. The *fetch & browse task* was introduced in 2005 for this reason. The aim was to first identify relevant documents (the fetching phase), and then to identify the most exhaustive and specific elements within the fetched documents (the browsing phase). In the fetching phase, documents had to be ranked according to how exhaustive and specific they were. In the browsing phase, ranking had to be done according to how exhaustive and specific the relevant elements in the document were, compared to other elements in the same document.

In 2006, the same task, renamed the *relevant in context sub-task*, required systems to return for each article an unranked set of non-overlapping elements, covering the relevant material in the document. In addition, a new task was introduced in 2006, the *best in context sub-task*, where the aim was to find the best-entry-point, here a single element, for starting to read articles with relevant information. This sub-task can be viewed as the extreme case of the fetch & browse approach, where only one element is returned per article.

## 2.4 Relevance

Most dictionaries define relevance as "pertinence to the matter at hand". In terms of IR, it is usually understood as the connection between a retrieved item and the user's query. In XML retrieval, the relationship between a retrieved item and the user's query is further complicated by the need to consider the structure in the documents. Since retrieved elements can be at any level of granularity, an element and one of its child elements can both be relevant to a given query, but the child element may be more focussed on the topic of the query than its parent element, which may contain additional irrelevant content. In this case, the child element is a better element to retrieve than its parent element, because not only it is relevant to the query, but it is also specific to the query. To accommodate the specificity aspect, INEX defined relevance along two dimensions:

- **Exhaustivity**, which measures how exhaustively an element discusses the topic of the user's request.
- **Specificity**, which measures the extent to which an element focuses on the topic of request (and not on other, irrelevant topics).

A multiple degree relevance scale was necessary to allow the explicit representation of how exhaustively a topic is discussed within an element with respect to its child

elements. For example, a section containing two paragraphs may be regarded more relevant than either of its paragraphs by themselves. Binary values of relevance cannot reflect this difference. INEX therefore adopted a four-point relevance scale [4]:

- Not exhaustive: The element does not contain any information about the topic.
- Marginally exhaustive: The element mentions the topic, but only in passing.
- Fairly exhaustive: The element discusses the topic, but not exhaustively.
- Highly exhaustive: The element discusses the topic exhaustively.

As for exhaustivity, a multiple degree scale was also necessary for the specificity dimension. This is to allow to reward retrieval systems that are able to retrieve the appropriate ("exact") sized elements. For example, a retrieval system that is able to locate the only relevant section in a book is more effective than one that returns a whole chapter. A four-point relevance scale was adopted:

- Not specific: the topic is not a theme discussed in the element.
- Marginally specific: the topic is a minor theme discussed in the element.
- Fairly specific: the topic is a major theme discussed in the element.
- Highly specific: the topic is the only theme discussed in the element.

Based on the combination of exhaustivity and specificity, it becomes possible to identify those relevant elements which are both exhaustive and specific to the topic of request and hence represent the most appropriate unit to return to the user. In the evaluation we can then reward systems that are able to retrieve these elements.

Obtaining relevance assessments is a very tedious and costly task [6]. An observation made in [1] was that the assessment process could be simplified if first, relevant passages of text were identified by highlighting, and then the elements within these passages were assessed. As a consequence, at INEX 2005, the assessment method was changed, leading to the redefinition of the scales for specificity. The procedure was a two-phase process. In the first phase, assessors highlighted text fragments containing only relevant information. The specificity dimension was then automatically measured on a continuous scale [0,1], by calculating the ratio of the relevant content of an XML element: a completely highlighted element had a specificity value of 1, whereas a non-highlighted element had a specificity value of 0. For all other elements, the specificity value was defined as the ratio (in characters) of the highlighted text (i.e. relevant information) to the element size. For example, an element with specificity of 0.72 has 72% of its content highlighted.

In the second phase, for all elements within highlighted passages (and parent elements of those), assessors were asked to assess their exhaustivity. Following the outcomes of extensive statistical analysis of the INEX 2004 results [5] - which showed that in terms of comparing retrieval effectiveness the same conclusions could be drawn using a smaller number of grades for the exhaustivity dimension<sup>3</sup> - INEX 2005 adopted the following 3 + 1 exhaustivity values:

<sup>3</sup> The same observation was reached for the specificity dimension, but as the assessment procedure was changed in INEX 2005, the new highlighting process allowed for a continuous scale of specificity to be calculated automatically.

- Highly exhaustive (2): the element discussed most, or all, aspects of the topic.
- Partly exhaustive (1): the element discussed only few aspects of the topic.
- Not exhaustive (0): the element did not discuss the topic.
- Too Small (?): the element contains relevant material but is too small to be relevant on its own.

The category of "too small" was introduced to allow assessors to label elements which, although contained relevant information, were too small to be able to sensibly reason about their level of exhaustivity. An extensive statistical analysis was performed on the INEX 2005 results [5], which showed that in terms of comparing retrieval performance, not using the exhaustivity dimension led to similar results. As a result, INEX 2006 dropped the exhaustivity dimension, and relevance was defined only along the specificity dimension.

### 3 Metrics

Measures of XML retrieval effectiveness must consider the dependency between elements. Unlike traditional IR, users in XML retrieval have access to other, structurally related elements from returned result elements. They may hence locate additional relevant information by browsing or scrolling. This motivates the need to consider so-called *near-misses*, which are elements from where users can access relevant content, within the evaluation. In this section, we restrict ourselves to the metrics used to evaluate the thorough and focussed sub-tasks, as the evaluation of the other sub-tasks is still an ongoing research issue.

The effectiveness of most ad-hoc retrieval tasks is measured by the established and widely used precision and recall metrics, or their variants. When using this family of measures, if we consider near-misses when evaluating retrieval effectiveness, then systems that return *overlapping* elements (e.g. both a paragraph and its enclosing section) will be evaluated as more effective than those that do not return overlapping elements (e.g. either the paragraph or its enclosing section). If both the paragraph and its enclosing section are relevant, then this family of effectiveness measures will count both these nested elements as separate relevant components that increase the count of relevant and retrieved elements. Therefore, despite not retrieving entirely new relevant information, systems that favour the retrieval of overlapping components would receive higher effectiveness scores. To address this problem, INEX used the XCG measures, which are an extension of the Cumulative Gain (CG) based measures [3]. These measures are not based on a counting mechanism, but on cumulative gains associated with returned results.

For each returned element, a gain value  $xG[.]$  is calculated, which is a value in the interval  $[0, 1]$ . A value of 0 reflects no gain, 1 is the highest gain value, and values between 0 and 1 represent various gain levels. The gain value depends on the element's exhaustivity and specificity. Given that INEX employs two relevance dimensions, the gain value is calculated as a combination of these dimensions, thus reflecting the worth of a retrieved element. INEX uses *quantisation functions* to provide a relative ordering of the various combinations of exhaustivity and specificity values and a mapping

of these to a single relevance scale in  $[0, 1]$ . Various quantisation functions have been used over the years as a means to model assumptions regarding the worth of retrieved elements to users or scenarios. For example, INEX 2003 used the quantisations defined below, where  $e$  and  $s$  stand, respectively, for exhaustivity and specificity.

$$quant_{strict}(e, s) := \begin{cases} 1 & \text{if } e = 3 \text{ and } s = 3, \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

The strict function is used to evaluate XML retrieval methods with respect to their capability of retrieving highly exhaustive and highly specific components.

$$quant_{gen}(e, s) := \begin{cases} 1 & \text{if } (e, s) = (3, 3), \\ 0.75 & \text{if } (e, s) \in \{(2, 3), (3, \{2, 1\})\}, \\ 0.5 & \text{if } (e, s) \in \{(1, 3), (2, \{2, 1\})\}, \\ 0.25 & \text{if } (e, s) \in \{(1, 2), (1, 1)\}, \\ 0 & \text{if } (e, s) = (0, 0). \end{cases} \quad (2)$$

The generalised function allows the reward of fairly and marginally relevant elements in the results. Other quantisations were introduced in subsequent years of INEX, emphasising specificity or exhaustivity. In [5], however, it is shown that, although quantisation functions express different user preferences, many of them behave similarly when ranking systems. As a consequence, one form of strict and one form of general quantisation functions have been used since 2005, and were modified to adapt to the new scale used in INEX 2005. In INEX 2006, as the exhaustivity dimension was dropped, the quantisation function simply maps an element to its specificity value.

Given a ranked list of elements  $e_j$ , each with their calculated gain value  $xC[e_j] = quant(e_j)$  where  $quant$  is a chosen quantisation function, the cumulative gain at rank  $i$ , denoted as  $xCG[i]$ , is computed as the sum of the relevance scores up to that rank:

$$xCG[i] := \sum_{j=1}^i xC(e_j) \quad (3)$$

For each query, an ideal gain vector,  $xCI$ , is derived by filling the rank positions with  $xG(c'_j)$  in decreasing order for all assessed elements  $c'_j$ . A retrieval run's  $xCG$  vector is compared to this ideal ranking by plotting both the actual and ideal cumulative gain functions against the rank position. Normalised  $xCG$  ( $nxCG$ ) is:

$$nxCG[i] := \frac{xCG[i]}{xCI[i]} \quad (4)$$

For a given rank  $i$ ,  $nxCG[i]$  reflects the relative gain the user has accumulated up to that rank, compared to the gain he/she could have attained if the system would have produced the optimum best ranking, where 1 represents ideal performance.

XCG also defines effort-precision/gain-recall ( $MAep$ ). The effort-precision  $ep$  at a given gain-recall value  $gr$  is defined as the number of visited ranks required to reach a given level of gain relative to the total gain that can be obtained, and is defined as:

$$ep[r] := \frac{i_{ideal}}{i_{run}} \quad (5)$$

where  $i_{ideal}$  is the rank position at which the cumulative gain of  $r$  is reached by the ideal system and  $i_{run}$  is the rank position at which the cumulative gain of  $r$  is reached by the system run. A score of 1 reflects ideal performance, i.e. when the user needs to spend the minimum necessary effort to reach a given level of gain. The gain-recall  $gr$  is calculated as:

$$gr[i] := \frac{xCG[i]}{xCI[n]} = \frac{\sum_{j=1}^i xG[j]}{\sum_{j=1}^n xI[j]} \quad (6)$$

where  $n$  is the number of elements  $c$  where  $xC[c] > 0$ . This method follows the same viewpoint as standard precision/recall, where recall is the control variable and precision the dependent variable. Here, the gain-recall is the control variable and effort-precision is the dependent variable. As with precision/recall, interpolation techniques are used to estimate effort-precision values at non-natural gain-recall points, e.g. when calculating effort-precision at standard recall points of  $[0.1, 1]$ , denoted as e.g.  $ep@0.1$ . For this purpose, a simple linear interpolation method is used. Also, the *non-interpolated mean average effort-precision*, denoted as *MAep*, is calculated by averaging the effort-precision values obtained for each rank where a relevant document is returned.

In the case of the thorough sub-task (when overlap is not an issue), the full recall-base is used to derive both the ideal gain vector  $xCI$  and the gain vectors,  $xCG$ .

For the focussed retrieval task, the elements in the ideal recall-base represent the desired target elements that should be retrieved, while all other elements in the full recall-base may be awarded partial scores. In this case, the ideal gain vector  $xCI$  is derived from the ideal recall-base, whereas the gain vectors,  $xCG$ , for the retrieval approaches under evaluation are based on the full recall-base to enable the scoring of near-miss elements. As any relevant elements of the full recall-base not included in the ideal recall-base are considered as near-misses, this strategy allows to support the evaluation viewpoint whereby elements in the ideal recall-base *should* be retrieved, whereas the retrieval of near-misses *could* be rewarded as partial success.

The construction of the ideal recall-base requires a preference function among exhaustivity and specificity value pairs. Quantisation functions are used for this purpose as these reflect the worth of retrieved elements. Given a chosen quantisation function, it is possible to quantify the value, or worth, of an element and identify the "best" components within an XML document as those elements with the highest quantised score. Also needed is a methodology for traversing an XML document (its tree structure) and selecting ideal elements based on their relative preference relations to their structurally related elements. The approach adopted in INEX, is to traverse the XML tree of a document bottom-up and to select the element with the highest quantised score. In the case where two elements have an equal score, the one higher in the XML structure is selected.

## 4 Conclusion

INEX has focused on developing an infrastructure, test collections, and appropriate scoring methods for evaluating the effectiveness of content-oriented XML retrieval. The initiative is now entering its sixth year, with INEX 2007 set to begin in April 2007. The major achievements in XML retrieval evaluation can be summarised as follows:

- A larger and more realistic test collection has been achieved with the addition of the Wikipedia documents. The content of the Wikipedia collection can also appeal to users with backgrounds other than computer science, making the carrying out of user studies with this collection more appropriate.
- A better understanding of information needs and retrieval scenarios. The set of retrieval tasks that were used at INEX 2006 is considered as a good representation of actual retrieval tasks that users of an XML retrieval system may wish to perform.
- A better understanding of how to measure the effectiveness of different retrieval systems by using appropriate metrics. In particular, we now have an understanding of how to deal with near-misses and overlapping elements, and which metrics to use under which retrieval assumptions.

In addition, INEX has been expanding in scope with the addition of a number of additional research tracks that tackle other IR problems related to XML documents. The additional tracks deal with issues such as retrieval of multimedia items, user interaction, retrieval from heterogeneous collections of documents, classification and clustering, etc. As an ongoing effort, empirical data about user behaviour for validating the effectiveness metrics are being considered. The current emphasis in INEX is to identify who the real users of XML retrieval systems are, how they might use retrieval systems and for which realistic tasks. A new research track, the user case studies track, is currently investigating this issue.

**Acknowledgments.** The authors would like to acknowledge the INEX organisers and participants for their valuable contributions throughout the various INEX campaigns.

## References

1. Clarke, C.: Range results in XML retrieval. In: Proceedings of the INEX Workshop on Element Retrieval Methodology (2005)
2. Denoyer, L., Gallinari, P.: The Wikipedia XML Corpus. SIGIR Forum 40(1) (2006)
3. Järvelin, K., Kekäläinen, J.: Cumulated gain-based evaluation of IR techniques. ACM TOIS 20(4), 422–446 (2002)
4. Kekäläinen, J., Järvelin, K.: Using graded relevance assessments in IR evaluation. JASIST 53(13), 1120–1129 (2002)
5. Ogilvie, P., Lalmas, M.: Investigating the exhaustivity dimension in content-oriented XML element retrieval evaluation. In: Proceedings of ACM CIKM, pp. 84–93 (2006)
6. Piwowarski, B., Lalmas, M.: Providing consistent and exhaustive relevance assessments for XML retrieval evaluation. In: Proceedings of ACM CIKM, pp. 361–370 (2004)
7. Kazai, G., Lalmas, M.: eXtended Cumulated Gain Measures for the Evaluation of Content-oriented XML Retrieval. ACM TOIS 24(4), 503–542 (2006)
8. Tombros, A., Malik, S., Larsen, B.: Report on the INEX 2004 interactive track. ACM SIGIR Forum 39(1) (2005)
9. Trotman, A., Sigurbjornsson, B.: Narrowed extended XPATH I (NEXI). In: Proceedings of the INEX Workshop on Element Retrieval Methodology (2004)