

# Evaluation of Effective XML Information Retrieval

A thesis submitted for the degree of  
Doctor of Philosophy

Jovan Pehcevski, Grad. Electro-Technical Eng.  
(University “St. Cyril and Methodius”, Skopje, Republic of Macedonia)  
School of Computer Science and Information Technology,  
Science, Engineering, and Technology Portfolio,  
RMIT University,  
Melbourne, Victoria, Australia.

28th August, 2006



This thesis is dedicated to my father, Dimitar Pehcevski,  
who never stopped believing

## **Declaration**

I certify that except where due acknowledgement has been made, the work is that of the author alone; the work has not been submitted previously, in whole or in part, to qualify for any other academic award; the content of the thesis is the result of work which has been carried out since the official commencement date of the approved research program; and, any editorial work, paid or unpaid, carried out by a third party is acknowledged.

Jovan Pehcevski

School of Computer Science and Information Technology

RMIT University

28th August, 2006

## Acknowledgments

*I've come up with a set of rules that describe our reactions to technologies:*

- 1. Anything that is in the world when you're born is normal and ordinary and is just a natural part of the way the world works.*
- 2. Anything that's invented between when you're fifteen and thirty-five is new and exciting and revolutionary and you can probably get a career in it.*
- 3. Anything invented after you're thirty-five is against the natural order of things.*

— Douglas Adams, *The Salmon of Doubt*, Harmony Books, 2002.

There is a common wisdom that every saga has a beginning. There is another not so common wisdom that the road towards a PhD degree does not always have an end (or at least not a satisfactory one). I believe I am one of the few blessed people in this world to have a companion like Irena, my real life partner and my immortal beloved, in travelling the road less travelled. I thank her for tolerating my (at times) peculiar research behaviour, my weird passion towards science-fiction books and movies, all my sleepless nights, and for everything that has happened in between during the last four crazy years. I do hope that, one bright and sunny day, I will be able to make all of this up to her. I also want to express my deepest gratitudes to members of our closest families: my mother Jordanka, my brother Dragan and his lovely wife Lidija, my grandfather Jovan and my other two grandparents, my uncle Venko with his family, Irena's father Metodi, Irena's mother Jelica, and Irena's brother Vlado. I thank them for their never-ending love and support, and for believing in our goals.

On my academic side, my deepest gratitudes go to my first thesis adviser, Dr. James Thom, for his thorough guidance during my PhD candidature. It was Jamie who introduced me to the academic research, who opened the doors to its endless opportunities, and who to this day is probably still wondering why I was writing papers while running out of time to write my thesis. His research experience and continuous feedback were invaluable for the achieved quality of my research work. I thank Dr. Seyed Mohammad Mehdi (Saied) Tahaghoghi, my second thesis adviser and my favourite lecturer, for his many useful writing suggestions and technical tips, and for his generous academic and personal support. Saied introduced me to the world of academic teaching, and assisted me to learn a great deal about

effective teaching techniques. I am also indebted to Dr. Anne-Marie Vercoustre, my valuable consultant and a former second thesis adviser, for her encouragement towards achieving my PhD goal. I am personally looking forward to continue our collaboration with Anne-Marie at INRIA later this year. I also want to express my gratitude to Michael Harris and Donal Ellis, two great Aussie mates, who taught me many useful course management tips and techniques. Finally, I thank Prof. Justin Zobel for managing our vibrant research group, and for organising the “Theory of IR” meetings where I learned a lot about the activities surrounding the exciting field of information retrieval.

On my socialising side, I thank Tome Jovanovski and Jasmina Karevski for their honest friendship during the last three years. They shared many of the joyful and most memorable moments with us here in Australia, and for that Irena and I are most grateful. There are some past and present members of my RMIT research group that I especially want to thank: Timo Volkmer, for opening the door into the world of Douglas Adams; Abhijit Chattaraj, for discussing and sharing with me the intricacies of the world of Douglas Adams; Falk Scholer, for our joint admiration of the world of elves, dwarfs, and wizards, which does not have anything to do with the world of Douglas Adams; Martin Plowman, for reminding me that I could well be one of the legal aliens that he is so desperately seeking to interview for his PhD thesis; Vaughan Shanks, for his many live demonstrations of the great Aussie spirit; and Ranjan Sinha, Bodo Billerbeck, Halil Ali, and Steven Garcia, for just being there.

Last, but definitely not least, I want to thank some of the academics being part of the wonderful INEX research community; indeed, without INEX, the research work presented in this thesis would not have been possible. I thank Prof. Mounia Lalmas, for her ongoing encouragement and support; Andrew Trotman, for his friendship and for our very many useful research discussions; Shlomo Geva, for his useful insights and always intriguing e-mails; and Benjamin Piwowarski, Birger Larsen, Tassos Tombros, Jaap Kamps, Gabriella Kazai, Arjen de Vries, Djoerd Hiemstra, and Miro Lehtonen, for engaging with me into many interesting and helpful e-mail discussions. I sincerely thank all of them for their help and support, and I truly hope we will continue our successful collaboration in future.

In the meantime, I continue to be very excited and enthusiastic about all things concerning *life*, the *universe*, and *everything*. But that’s most likely because I am still only thirty-four.

— Jovan Pehcevski, Melbourne, Australia, August 2006

## Credits

Portions of the material in this thesis have previously appeared in the following publications:

- “Hybrid XML Retrieval: Combining Information Retrieval and a Native XML Database”, *Information Retrieval* [Pehcevski et al., 2005b] (Chapter 3).
- “Relevance in XML Retrieval: The User Perspective”, *Proceedings of the SIGIR 2006 Workshop on XML Element Retrieval Methodology* [Pehcevski, 2006] (Chapter 4).
- “HiXEval: Highlighting XML Retrieval Evaluation”, *Advances in XML Information Retrieval and Evaluation: Fourth Workshop of the INitiative for the Evaluation of XML Retrieval (INEX 2005)*, LNCS 3977 [Pehcevski and Thom, 2006] (Chapters 4 and 5).
- “RMIT University at INEX 2005: Ad hoc Track”, *Advances in XML Information Retrieval and Evaluation: Fourth Workshop of the INitiative for the Evaluation of XML Retrieval (INEX 2005)*, LNCS 3977 [Pehcevski et al., 2006] (Chapter 6).
- “Combining Image and Structured Text Retrieval”, *Advances in XML Information Retrieval and Evaluation: Fourth Workshop of the INitiative for the Evaluation of XML Retrieval (INEX 2005)*, LNCS 3977 [Awang Iskandar et al., 2006] (Chapter 6 and Appendix C).
- “Users and Assessors in the Context of INEX: Are Relevance Dimensions Relevant?”, *Proceedings of the INEX 2005 Workshop on Element Retrieval Methodology* [Pehcevski et al., 2005c] (Chapter 4 and Appendix B).
- “Hybrid XML Retrieval Revisited”, *Advances in XML Information Retrieval: Third International Workshop of the INitiative for the Evaluation of XML Retrieval (INEX 2004)*, LNCS 3493 [Pehcevski et al., 2005a] (Chapter 3).
- “Enhancing Content-And-Structure Information Retrieval using a Native XML Database”, *Proceedings of the first Twente Data Management Workshop (TDM’04) on XML Databases and Information Retrieval* [Pehcevski et al., 2004b] (Chapter 3).
- “RMIT INEX Experiments: XML Retrieval using Lucy/eXist”, *Proceedings of the INEX 2003 workshop* [Pehcevski et al., 2004a] (Chapter 3).

- “XML-Search Query Language: Needs and Requirements”, *Proceedings of the AUSWeb 2003 conference* [Pehcevski et al., 2003] (Chapters 1 and 2).

This work was supported by the International Postgraduate Research Scholarship (IPRS), the Australian Research Council and the Search Engine Group at RMIT University.

The thesis was written using the `KWrite` editor on Mandrake GNU/Linux, and typeset using the  $\text{\LaTeX} 2_{\epsilon}$  document preparation system.

All trademarks are the property of their respective owners.

### **Note**

Unless otherwise stated, all fractional results have been rounded to the displayed number of decimal figures.

# Contents

<b>Abstract</b>	<b>1</b>
<b>1 Introduction</b>	<b>3</b>
1.1 Basic XML retrieval concepts . . . . .	4
1.2 INEX . . . . .	8
1.3 Challenges in XML retrieval . . . . .	11
1.3.1 How can <i>information retrieval</i> and <i>database</i> techniques be combined for effective XML retrieval? . . . . .	12
1.3.2 What does <i>user and assessor experience</i> suggest about how <i>relevance</i> should be defined in XML retrieval? . . . . .	13
1.3.3 How should the effectiveness of XML retrieval be <i>evaluated</i> ? . . . . .	14
1.3.4 How effective is XML retrieval in different <i>application scenarios</i> ? . . . .	15
1.4 Thesis structure . . . . .	18
<b>2 XML Information Retrieval</b>	<b>21</b>
2.1 XML retrieval approaches . . . . .	22
2.1.1 Query languages . . . . .	22
2.1.2 Full-text information retrieval approaches . . . . .	23
2.1.3 Native XML database approaches . . . . .	27
2.1.4 Scoring approaches . . . . .	28
2.2 Relevance in information retrieval . . . . .	36
2.2.1 Definitions and dimensions . . . . .	36
2.2.2 INEX relevance . . . . .	37
2.3 Evaluation approaches . . . . .	40
2.3.1 Assumptions . . . . .	40

2.3.2	Metrics and measures . . . . .	43
2.3.3	Significance, fidelity, and reliability . . . . .	56
2.4	Methodology of XML element retrieval . . . . .	58
2.5	Summary . . . . .	62
<b>3</b>	<b>Hybrid XML Retrieval</b>	<b>65</b>
3.1	Technological aspects . . . . .	66
3.1.1	A full-text information retrieval approach . . . . .	66
3.1.2	A native XML database approach . . . . .	68
3.1.3	A hybrid approach to XML retrieval . . . . .	69
3.2	Retrieval modelling aspects . . . . .	70
3.2.1	Identifying the appropriate answer granularity . . . . .	70
3.2.2	Ranking the final answers . . . . .	72
3.2.3	Tuning the retrieval parameters . . . . .	74
3.3	Experiments on INEX 2003 and 2004 test collections . . . . .	79
3.3.1	Evaluation methodology . . . . .	79
3.3.2	INEX 2003 experiments . . . . .	86
3.3.3	INEX 2004 experiments . . . . .	97
3.4	Summary . . . . .	105
<b>4</b>	<b>Relevance in XML Retrieval</b>	<b>107</b>
4.1	Analysis of INEX 2004 relevance . . . . .	108
4.1.1	INEX 2004 relevance dimensions . . . . .	108
4.1.2	Methodology . . . . .	108
4.1.3	Assessor behaviour analysis for INEX 2004 CO topics . . . . .	112
4.1.4	User behaviour analysis for INEX 2004 Interactive topics . . . . .	115
4.1.5	Analysis of the level of agreement . . . . .	118
4.1.6	Concluding remarks on INEX 2004 relevance . . . . .	122
4.2	Analysis of INEX 2005 relevance . . . . .	124
4.2.1	INEX 2005 relevance dimensions . . . . .	124
4.2.2	Assessor behaviour analysis for INEX 2005 topics . . . . .	124
4.2.3	Analysis of the level of agreement . . . . .	127
4.2.4	Concluding remarks on INEX 2005 relevance . . . . .	131
4.3	A topical-hierarchical relevance definition . . . . .	132

4.3.1	Relevance dimensions . . . . .	132
4.3.2	Relevance scale . . . . .	134
4.3.3	User satisfaction . . . . .	135
4.4	Experiments with the new relevance definition . . . . .	136
4.4.1	Comparison to the INEX 2004 relevance definition . . . . .	136
4.4.2	Comparison to the INEX 2005 relevance definition . . . . .	139
4.5	Summary . . . . .	144
<b>5</b>	<b>Evaluation of XML Retrieval</b>	<b>149</b>
5.1	HiXEval: Highlighting XML retrieval evaluation . . . . .	150
5.1.1	Assumptions . . . . .	150
5.1.2	Formal definition . . . . .	151
5.2	Fidelity tests . . . . .	155
5.2.1	Simulated runs . . . . .	156
5.2.2	Expected rankings . . . . .	160
5.2.3	System-oriented task . . . . .	161
5.2.4	User-oriented task . . . . .	166
5.2.5	Further analysis of evaluation behaviour . . . . .	169
5.2.6	Concluding comments on HiXEval fidelity . . . . .	171
5.3	HiXEval versus XCG in XML retrieval experiments . . . . .	173
5.3.1	Comparison of run orderings . . . . .	173
5.3.2	Reliability tests . . . . .	179
5.4	Summary . . . . .	183
<b>6</b>	<b>Scenarios of XML Retrieval</b>	<b>185</b>
6.1	Ad-hoc retrieval scenario . . . . .	186
6.1.1	XML retrieval approach . . . . .	186
6.1.2	INEX 2005 CO and +S sub-tasks . . . . .	192
6.1.3	INEX 2005 CAS sub-task . . . . .	197
6.2	Multimedia retrieval scenario . . . . .	199
6.2.1	XML retrieval approach . . . . .	202
6.2.2	INEX 2005 MM task . . . . .	206
6.3	Summary . . . . .	210

<b>7</b>	<b>Conclusions and Future Work</b>	<b>213</b>
7.1	Hybrid approach for effective XML retrieval . . . . .	213
7.2	New relevance definition for XML retrieval . . . . .	215
7.3	Highlighting XML retrieval evaluation . . . . .	217
7.4	Application scenarios of XML retrieval . . . . .	219
7.5	Conclusion summary . . . . .	221
<b>A</b>	<b>A similarity framework for XML retrieval</b>	<b>223</b>
A.1	Motivation . . . . .	223
A.2	TF/IDF ranking model . . . . .	224
A.3	Similarity framework . . . . .	228
A.4	Modelling scoring approaches . . . . .	239
<b>B</b>	<b>Measuring overlap</b>	<b>241</b>
B.1	Level of overlap for background topic B1 . . . . .	244
B.2	Level of overlap for comparison topic C2 . . . . .	245
<b>C</b>	<b>INEX 2005 MM track experiments</b>	<b>247</b>
C.1	TRECEval analysis . . . . .	248
C.2	HiXEval analysis . . . . .	248
C.3	Comparison of run orderings . . . . .	252
	<b>Bibliography</b>	<b>255</b>

# Glossary of acronyms

A-overlap	Ascendants overlap
AP	Average Precision
BA	Broad Answer
BEP	Best Entry Point
CAS	Content And Structure
CBEP	Combined BEP
CBIR	Content-Based Image Retrieval
CG	Cumulated Gain
CO	Content Only
CO+S	Content Only plus Structure
CRE	Coherent Retrieval Element
D-overlap	Descendants overlap
DFR	Divergence From Randomness
EA	Exact Answer
ep/gr	effort precision / gain recall
EPRUM	Expected Precision-Recall with User Modelling
Ex	Exhaustivity
F@r	F-measure (harmonic mean between Precision and Recall) at rank r
FS	Fully Seen
FullRB	Full Recall-Base
GIFT	GNU Image Finding Tool
HSV	Hue Saturation Value
HiXEval	Highlighting XML Retrieval Evaluation
HyREX	Hypermedia Retrieval Engine for XML

IDF	Inverse Document Frequency
IDL	Inverse Document Length
IEF	Inverse Element Frequency
iMAP	interpolated Mean Average Precision
INEX	INitiative for the Evaluation of XML Retrieval
IR	Information Retrieval
LCA	Lowest Common Ancestor
MA	Mutually Agreed
MAP	Mean Average Precision
MAep	Mean Average effort precision
MM	Multimedia
MpE	Ranking heuristic: more matching elements (M), shorter XPath length (p), nearer to end (E)
NA	Narrow Answer
nCG	normalised Cumulated Gain
nCRE	CRE that represents either LCA of at least two matching elements, or a matching element whose parent is not recognised as LCA
NEXI	Narrowed Extended XPath I
NR	Not Relevant
NS	Not Seen
nxCG	normalised extended Cumulated Gain
oCRE	CRE that represents LCA of at least two matching elements
O-overlap	Overall overlap
P-overlap	Probabilistic overlap
P@r	Precision at rank r
PA	Partial Answer
PBEP	Parent BEP
PME	Ranking heuristic: longer XPath length (P), more matching elements (M), nearer to end (E)
PS	Partially Seen
PTF	Ranking heuristic: longer XPath length (P), more distinct query terms (T), more frequent query term occurrences (F)
Prec	Precision

R-prec	Mean Recall-precision
R@r	Recall at rank r
RMIT	Royal Melbourne Institute of Technology
RP	Recall-precision
RSV	Retrieval Status Value
Rec	Recall
SBEP	Start-reading-here BEP
SCAS	Strict Content And Structure
SS	CAS query interpretation: strict target and strict support elements
SV	CAS query interpretation: strict target and vague support elements
Sp	Specificity
T2I	Tolerance to Irrelevance
TF	Term Frequency
TPF	Ranking heuristic: more distinct query terms (T), longer XPath length (P), more frequent query term occurrences (F)
TREC	Text REtrieval Conference
Trel	Total amount of relevant information
VCAS	Vague Content And Structure
VS	CAS query interpretation: vague target and strict support elements
VV	CAS query interpretation: vague target and vague support elements
WWW	World Wide Web
W3C	World Wide Web Consortium
XCG	Extended Cumulated Gain
XML	eXtensible Markup Language
XPath	XML Path language
XQuery	XML Query language
XSL	eXtensible Stylesheet Language

# List of Figures

1.1	A very small XML document collection containing just two XML documents.	5
1.2	A Document Type Definition used to validate XML documents . . . . .	6
1.3	INEX 2005 CO+S topic 203 . . . . .	10
2.1	A sample of relevance assessments for INEX 2005 CO+S topic 203 . . . . .	41
2.2	Identifying ideal elements from a recall-base . . . . .	53
3.1	INEX 2003 CO Topic 99 . . . . .	66
3.2	Three approaches to XML retrieval: a full-text information retrieval approach, a native XML database approach, and a hybrid XML retrieval approach . . .	67
3.3	Identifying matching, oCRE, and nCRE retrieval elements . . . . .	72
3.4	Tuning retrieval parameters in three Zettair similarity measures using CO topics of the INEX 2002 test collection . . . . .	75
3.5	Performance results for INEX 2002 runs using MpE ranking heuristic and eight distinct cases of retrieved elements per document, obtained with MAP using strict quantisation in <code>inex_eval</code> . . . . .	78
3.6	A sample of relevance assessments for INEX 2003 CO topic 99 . . . . .	81
3.7	Identifying General and Specific highly relevant elements . . . . .	82
3.8	Distribution of highly relevant elements across INEX 2003 and 2004 CO topics, using three cases of relevance assessments . . . . .	83
3.9	Categories of INEX 2003 and 2004 CO topics using General relevance assessments	84
3.10	Evaluation of the overall performance of three INEX 2003 CO runs, using strict quantisation in <code>inex_eval</code> . . . . .	89
3.11	Evaluation of the overall performance of three INEX 2003 CAS runs, using strict quantisation in <code>inex_eval</code> . . . . .	96

3.12	Evaluation of the overall performance of three INEX 2004 CO runs, using strict quantisation in <code>inex_eval</code> . . . . .	99
3.13	Evaluation of the overall performance of four INEX 2004 CAS runs, using strict quantisation in <code>inex_eval</code> . . . . .	103
3.14	Distribution of highly relevant elements across INEX 2004 CAS topics using three cases of relevance assessments, and two categories of INEX 2004 CAS topics identified in the case of General relevance assessments . . . . .	104
4.1	Background topic B1 used in INEX 2004 Interactive track . . . . .	110
4.2	Comparison topic C2 used in INEX 2004 Interactive track . . . . .	111
4.3	Distribution of relevant elements across nine relevance points for INEX 2004 CO topics, as judged by 34 assessors . . . . .	113
4.4	Distribution of relevant elements across nine relevance points for INEX 2004 Interactive topics, as judged by 88 users . . . . .	116
4.5	Identifying Exact, Partial, Broad, and Narrow relevant elements . . . . .	143
5.1	Evaluation of the overall performance of simulated runs for a system-oriented task, using <code>HiXEval</code> . . . . .	165
5.2	Evaluation of the overall performance of simulated runs for a user-oriented task, using <code>HiXEval</code> . . . . .	168
6.1	Evaluation of the overall performance of three CO runs submitted in INEX 2005 <code>FetchBrowse</code> article-level retrieval strategy, using <code>HiXEval</code> . . . . .	196
6.2	INEX 2005 MM topic 6, with image <code>BN7386.10.jpg</code> in the target element . . . . .	200
6.3	Querying image <code>BN7386.10.jpg</code> using <code>GIFT</code> . . . . .	204
6.4	<code>GIFT</code> results for a sample image query . . . . .	205
6.5	Evaluation of the overall performance of six RMIT runs officially submitted in INEX 2005 MM track, using <code>HiXEval</code> . . . . .	208
6.6	Performance results for additional RMIT runs obtained with Precision at rank cut-offs 1, 5 and 10 in <code>HiXEval</code> , as parameter $\beta$ varies from 0.0 to 1.0 . . . . .	209
6.7	Performance results for additional RMIT runs obtained with <code>MAP</code> and <code>R-prec</code> in <code>HiXEval</code> , as parameter $\beta$ varies from 0.0 to 1.0 . . . . .	210
B.1	Two sets of highly relevant elements drawn from relevance assessments for INEX 2004 CO topics 192 and 198 . . . . .	242

C.1	Evaluation of the overall performance of best performing runs submitted by INEX 2005 MM track participants, using <code>TRECEval</code> . . . . .	251
C.2	Evaluation of the overall performance of best performing runs submitted by INEX 2005 MM track participants, using <code>HiXEval</code> . . . . .	252

# List of Tables

2.1	List of documents obtained for a ranked query using a full-text information retrieval approach . . . . .	27
2.2	List of matching and LCA elements for two Boolean queries using a native XML database approach . . . . .	29
2.3	The 10-point relevance scale, used in INEX 2003 and 2004 . . . . .	39
2.4	Quantisation functions used in INEX since 2002 . . . . .	46
3.1	eXist list of matching elements for INEX 2003 CO topic 99 . . . . .	69
3.2	Ranked list of CREs for INEX 2003 CO topic 99, obtained with MpE heuristic	73
3.3	Ranked list of CREs for INEX 2003 CO topic 99, obtained with PME heuristic	74
3.4	Performance results for INEX 2002 runs using 16 CRE ranking heuristics and two types of answer elements, obtained with MAP using strict quantisation in <code>inex_eval</code> . . . . .	76
3.5	Performance results for INEX 2002 runs using different combinations of ranking heuristics and answer elements, obtained with MAP using five quantisations in <code>inex_eval</code> . . . . .	77
3.6	List of INEX 2003 and 2004 CO and CAS runs . . . . .	80
3.7	Performance results for five INEX 2003 CO runs and 20 distinct cases of retrieved elements per document, obtained with MAP using strict quantisation in <code>inex_eval</code> . . . . .	87
3.8	Performance results for five INEX 2003 CO runs and three distinct cases of retrieved elements per document, calculated under three CO topic categories and the case of Original relevance assessments . . . . .	90

3.9	Performance results for five INEX 2003 CO runs and three distinct cases of retrieved elements per document, calculated under three CO topic categories and the case of General relevance assessments . . . . .	91
3.10	Performance results for five INEX 2003 CO runs and three distinct cases of retrieved elements per document, calculated under three CO topic categories and the case of Specific relevance assessments . . . . .	92
3.11	Performance results for six INEX 2003 CO runs and three distinct cases of retrieved elements per document, obtained with MAP using strict quantisation in <code>inex_eval_ng</code> . . . . .	93
3.12	Performance results for INEX 2003 CAS runs when all elements are retrieved per document, obtained with MAP using strict quantisation in <code>inex_eval</code> . . .	95
3.13	Performance results for three INEX 2003 CAS runs when all elements are retrieved per document, calculated under three SCAS topic categories . . . .	97
3.14	Performance results for INEX 2004 CO runs when all elements are retrieved per document, obtained with MAP using strict quantisation in <code>inex_eval</code> . . .	98
3.15	Performance results for three INEX 2004 CO runs when all elements are retrieved per document, calculated under three CO topic categories and the case of General relevance assessments . . . . .	100
3.16	Performance results for INEX 2004 CAS runs when all elements are retrieved per document, obtained with MAP using strict quantisation in <code>inex_eval</code> . . .	102
3.17	Performance results for four INEX 2004 CAS runs when all elements are retrieved per document, calculated under three VCAS topic categories and the case of General relevance assessments . . . . .	105
4.1	Distribution of relevant elements with four element names across INEX 2004 CO topics, as judged by 34 assessors . . . . .	114
4.2	Correlation between grades of the two INEX 2004 relevance dimensions, as judged by 34 assessors . . . . .	115
4.3	Distribution of relevant elements with four element names across INEX 2004 Interactive topics, as judged by 88 users . . . . .	117
4.4	Correlation between grades of the two INEX 2004 relevance dimensions, as judged by 88 users . . . . .	117
4.5	Level of assessor and user agreement for Background topic B1 . . . . .	119
4.6	Level of assessor and user agreement for Comparison topic C2 . . . . .	121

4.7	Statistical analysis of distribution of E? (too-small), E1, and E2 relevant elements across the 29 CO+S and 34 VVCAS INEX 2005 topics . . . . .	126
4.8	Document-level and element-level assessor agreement for five topics double-judged at INEX 2005 . . . . .	128
4.9	Fine-grained element-level assessor agreement for five topics double-judged at INEX 2005 . . . . .	130
4.10	Statistical analysis of user responses on questions Q4.5 and Q4.6, gathered from 29 users that participated in Task C of INEX 2005 Interactive track . .	135
4.11	Number of users that chose a combination of responses on questions Q4.5 and Q4.6, used in Task C of INEX 2005 Interactive track . . . . .	136
4.12	Level of assessor and user agreement for topics B1 and C2, calculated for 20 mappings between the INEX 2004 10-point relevance scale and the new five-point relevance scale . . . . .	138
4.13	Statistical analysis of overall distribution of user and assessor judgements, calculated across two General and two Challenging topics used in Task A of INEX 2005 Interactive track . . . . .	140
4.14	Statistical analysis of overall distribution of user and assessor judgements, calculated across four topics used in Task C of INEX 2005 Interactive track .	142
4.15	Statistical analysis of distribution of Exact, Partial, Broad, and Narrow relevant elements across the 29 CO+S and 34 VVCAS INEX 2005 topics . . . . .	145
4.16	Statistical analysis of distribution of three Exhaustivity values across Exact, Partial, Broad, and Narrow relevant elements found for the 29 CO+S and 34 VVCAS INEX 2005 topics . . . . .	146
5.1	Statistical analysis of distribution of passages, best entry points, and all relevant elements across the 29 INEX 2005 CO+S topics . . . . .	158
5.2	Simulated runs created from a recall-base for INEX 2005 CO+S topic 203 . .	159
5.3	Performance results for simulated runs using a system-oriented task, obtained with measures in HiXEval . . . . .	162
5.4	Performance results for simulated runs using a user-oriented task, obtained with measures in HiXEval . . . . .	167
5.5	Performance results for simulated runs using two evaluation scenarios, obtained with three rank cutoff measures in HiXEval . . . . .	170

5.6	Comparing run orderings obtained with pairs of evaluation measures from two XCG metrics and HiXEval, using 55 submitted runs in <i>Thorough</i> retrieval strategy	174
5.7	Comparing run orderings obtained with pairs of evaluation measures from two XCG metrics and HiXEval, using 44 submitted runs in <i>Focussed</i> retrieval strategy	175
5.8	Comparing run orderings obtained with pairs of evaluation measures from two XCG metrics and HiXEval, using 31 submitted runs in <i>FetchBrowse</i> article-level and element-level retrieval strategies . . . . .	177
5.9	Comparing run orderings obtained with pairs of evaluation measures from two XCG metrics and HiXEval, using 25 submitted runs in SSCAS, 23 runs in SVCAS and VSCAS, and 28 runs in VVCAS retrieval strategies . . . . .	178
5.10	Reliability tests for measures from two XCG metrics and HiXEval, using 55 runs submitted in <i>Thorough</i> retrieval strategy . . . . .	180
5.11	Reliability tests for measures from two XCG metrics and HiXEval, using 44 runs submitted in <i>Focussed</i> retrieval strategy . . . . .	181
5.12	Reliability tests for measures from two XCG metrics and HiXEval, using 31 runs submitted in <i>FetchBrowse</i> element-level retrieval strategy . . . . .	182
6.1	Ranked list of CREs for INEX 2005 CO+S topic 203, obtained with TPF ranking heuristic . . . . .	188
6.2	Ranked list of CREs for INEX 2005 CO+S topic 203, obtained with PTF ranking heuristic . . . . .	189
6.3	List of nine CO, nine +S, and eight CAS runs submitted in different retrieval strategies of INEX 2005 Ad-hoc track . . . . .	191
6.4	Performance results for three CO and three +S runs submitted in INEX 2005 <i>Thorough</i> retrieval strategy, obtained with measures in HiXEval . . . . .	192
6.5	Performance results for three CO and three +S runs submitted in INEX 2005 <i>Focussed</i> retrieval strategy, obtained with measures in HiXEval . . . . .	194
6.6	Performance results for three CO and three +S runs submitted in INEX 2005 <i>FetchBrowse</i> article-level strategy, obtained with measures in HiXEval . . . . .	195
6.7	Performance results for three CO and three +S runs submitted in INEX 2005 <i>FetchBrowse</i> element-level strategy, obtained with measures in HiXEval . . . . .	197
6.8	Performance results for CAS runs submitted in INEX 2005 SS, SV, VS, and VV retrieval strategies, obtained with measures in HiXEval . . . . .	198
6.9	Distribution of two types of INEX 2005 MM topics across three categories . . . . .	201

6.10	List of six RMIT CAS runs officially submitted in the INEX 2005 MM track	206
6.11	Performance results for six RMIT runs officially submitted in the INEX 2005 MM track, obtained with measures in <code>HiXEval</code>	207
A.1	Document ( $S_{q,d}$ ), text-node ( $S_{q,p}$ ), and element-node ( $S_{q,e}$ ) combining similarity functions	230
A.2	Term weights $w_t$	231
A.3	Document ( $w_{d,t}$ ), text-node ( $w_{p,t}$ ), element-node ( $w_{e,t}$ ), and query ( $w_{q,t}$ ) term weights	232
A.4	Relative term frequencies used in document ( $r_{d,t}$ ), text-node ( $r_{p,t}$ ), element-node ( $r_{e,t}$ ), and query ( $r_{q,t}$ ) context	233
A.5	Document ( $W_d$ ), text-node ( $W_p$ ), element-node ( $W_e$ ), and query ( $W_q$ ) lengths	234
A.6	Q-expression <code>D[BB-ABB-BAA]</code> used in document context, representing the cosine similarity measure	235
A.7	Q-expression <code>E[DM-AFI-BAA]</code> used in element-node context, representing the multinomial language model	237
A.8	Text-node propagation functions $\mathcal{F}(e, p)$	239
A.9	List of scoring approaches modelled by the similarity framework	240
B.1	Overlap values for four element sets, obtained with four overlap measures	243
B.2	Statistical analysis of the level of overlap among highly relevant elements found in relevance judgements for topic B1, as judged by 41 users	244
B.3	Statistical analysis of the level of overlap among highly relevant elements found in relevance judgements for topic C2, as judged by 36 users	245
C.1	Evaluation of the overall performance of runs officially submitted by INEX 2005 MM track participants, obtained with <code>P@r</code> (values 1, 5, and 10), <code>MAP</code> and <code>R-Prec</code> measures in <code>TRECEval</code>	249
C.2	Evaluation of the overall performance of runs officially submitted by INEX 2005 MM track participants, obtained with <code>P@r</code> (values 1, 5, and 10), <code>MAP</code> and <code>R-Prec</code> measures in <code>HiXEval</code>	250
C.3	Comparing run orderings obtained with pairs of evaluation measures from <code>TRECEval</code> and <code>HiXEval</code> , using 25 submitted runs in INEX 2005 MM track	253



# Abstract

XML is being adopted as a common storage format in scientific data repositories, digital libraries, and on the World Wide Web. Accordingly, there is a need for content-oriented XML retrieval systems that can efficiently and effectively store, search and retrieve information from XML document collections. Unlike traditional information retrieval systems where whole documents are usually indexed and retrieved as information units, XML retrieval systems typically index and retrieve document components of varying granularity. To evaluate the effectiveness of such systems, test collections — where relevance assessments are provided according to an XML-specific definition of relevance — are necessary. Such test collections have been built during four rounds of the Initiative for the Evaluation of XML Retrieval (INEX).

There are many different approaches to XML retrieval; most approaches either extend full-text information retrieval systems to handle XML retrieval, or use database technologies that incorporate existing XML standards to handle both XML presentation and retrieval. We present a *hybrid approach* to XML retrieval that combines text information retrieval features with XML-specific features found in a native XML database. Results from our experiments on the INEX 2003 and 2004 test collections demonstrate the usefulness of applying our hybrid approach to different XML retrieval tasks.

A realistic definition of *relevance* is necessary for meaningful comparison of alternative XML retrieval approaches. The three relevance definitions used by INEX since 2002 comprise two relevance dimensions, each based on topical relevance. We perform an extensive analysis of the two INEX 2004 and 2005 relevance definitions, and show that assessors and users find them difficult to understand. We propose a new definition of relevance for XML retrieval, and demonstrate that a relevance scale based on this definition is useful for XML retrieval experiments.

Finding the appropriate approach to evaluate XML retrieval effectiveness is the subject of ongoing debate within the XML information retrieval research community. We present an overview of the evaluation methodologies implemented in the current INEX metrics, which reveals that the metrics follow different assumptions and measure different XML retrieval behaviours. We propose a new *evaluation metric* for XML retrieval and conduct an extensive analysis of the retrieval performance of simulated runs to show what is measured. We compare the evaluation behaviour obtained with the new metric to the behaviours obtained with two of the official INEX 2005 metrics, and demonstrate that the new metric can be used to reliably evaluate XML retrieval effectiveness.

To analyse the effectiveness of XML retrieval in different *application scenarios*, we use evaluation measures in our new metric to investigate the behaviour of XML retrieval approaches under the following two scenarios: the ad-hoc retrieval scenario, exploring the activities carried out as part of the INEX 2005 Ad-hoc track; and the multimedia retrieval scenario, exploring the activities carried out as part of the INEX 2005 Multimedia track. For both application scenarios we show that, although different values for retrieval parameters are needed to achieve the optimal performance, the desired textual or multimedia information can be effectively located using a combination of XML retrieval approaches.