

# Understanding Content-and-Structure

Jaap Kamps

Maarten Marx

Maarten de Rijke

Börkur Sigurbjörnsson

Informatics Institute

University of Amsterdam

{kamps,marx,mdr,borkur}@science.uva.nl

*INEX 2005 Methodology Workshop*

Glasgow, July, 2005

## Overview

- ▶ Expressive power of query languages
- ▶ What do users mean by a target constraints?
- ▶ What sort of query trees are popular?

## Who is our user?

- ▶ Our ideal user
  - ◇ Experienced searcher
  - ◇ Not necessarily technically literate
  - ◇ Example: information professional, librarian■
- ▶ Ignorant users
  - ◇ **Some** knowledge of the **tag-names**
  - ◇ **No** knowledge of the **hierarchy** of tag-names
- ▶ Semi-ignorant users
  - ◇ **Some** knowledge of the **tag-names**
  - ◇ **Some** knowledge of the **hierarchy** of tag-names

## Expressive power of query languages

- ▶ Assign users a query language that matches their knowledge
  - ◇ The expressive power of the query language should fit the expressive power of the user
- ▶ A query language should be safe and complete
  - ◇ **Safety**: If the user cannot distinguish A and B, then no query can
  - ◇ **Completeness**: If the user can distinguish A and B, then a query can■
- ▶ A note on XPath
  - ◇ Expressiveness of XPath fragments is a hot research topic
  - ◇ It's ideal research purposes for talking about trees
  - ◇ It's not meant as an end-user language

## Overview

- ▶ Expressive power of query languages
- ▶ What do users mean by a target constraints?
- ▶ What sort of trees are popular?

## User wants a specific granularity?

- ▶ E3S3 judgements from 2003 and 2004

	article+	sec+	p+	abs	vt
article (11)	<b>42.89%</b>	25.62%	18.79%	0.19%	0.76%
sec (20)	10.17%	<b>38.60%</b>	25.92%	1.36%	0.20%
p (5)	11.26%	21.13%	<b>49.30%</b>	3.52%	–
abs (4)	14.40%	37.29%	22.88%	<b>11.86%</b>	–
vt (2)	–	–	42.86%	–	<b>53.13%</b>

- ▶ Users seem not to take granularity constraint seriously
- ▶ Users seem to have a bias toward wanting what they ask
  - ◇ Information need controls granularity?
  - ◇ Target constraint controls assessments? ■
- ▶ But the assessment guidelines said ... (to be continued)

## Can we have structured information needs?

- ▶ “But the assessment guidelines said that structural constraints should be ignored”
  - ◇ I’d say, rightfully so
- ▶ Why should we consider structural constraints strictly?■
- ▶ What is an abstract?
  - ◇ Text inside an `<abs>` tag?
  - ◇ “A statement summarizing the important points of a text”? (dictionary.com)
  - ◇ Wouldn’t an **introduction** section be satisfactory?
  - ◇ Wouldn’t a **conclusions** section be satisfactory?■
- ▶ Claim: Structure is never an inherent part of an information need

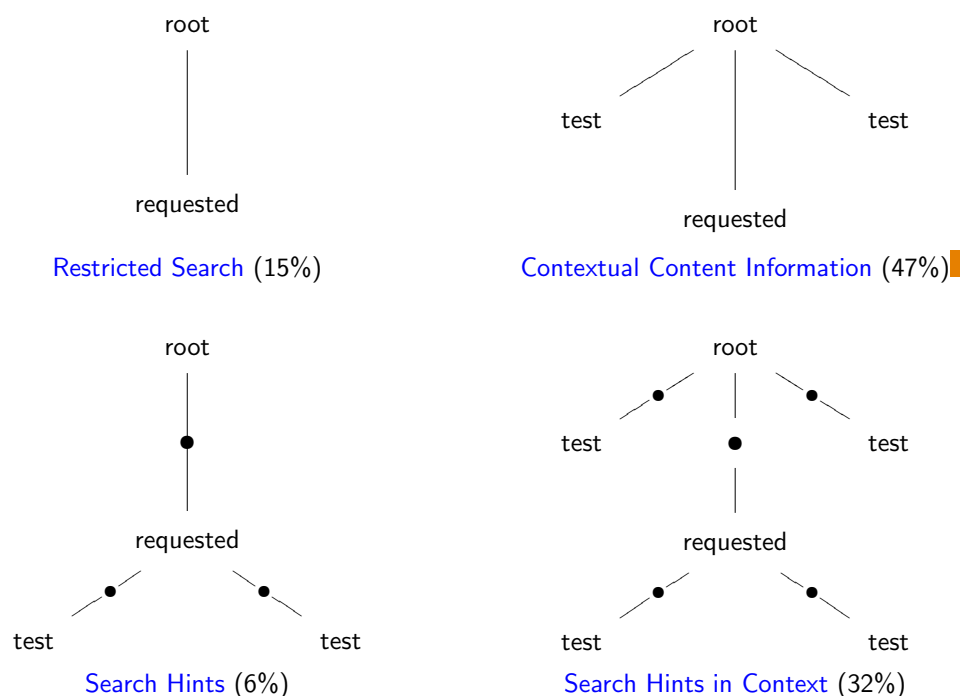
## Overview

- ▶ Expressive power of query languages
- ▶ What do users mean by a target constraints?
- ▶ What sort of trees are popular in the INEX collection?

## Classification of queries

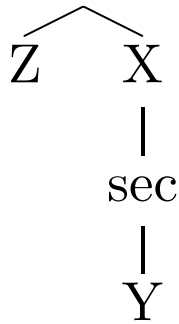
- ▶ We classify the INEX queries based on whether or not they express **hierarchical relationships** between elements
- ▶ We do not consider hierarchical relationships with the `<article>` tag since that relationship is trivial
- ▶ Examples of non-hierarchical queries:
  - ◇ `//sec[about(., java thread implementation)]`
  - ◇ `//article[about(../abs, java)]//sec[about(., thread implementation)]`
- ▶ Examples of hierarchical queries:
  - ◇ `//article[about(../fm//abs, java)]//bdy//sec[about(., thread implementation)]`

## Shapes of the queries



## Popular trees

article



- ▶ Most popular query template

- ▶ Appears 14 times among the 2004 queries

- ▶ Popular replacements for  $X$

◇  $\epsilon$  (10); body (4)

- ▶ Popular replacements for  $Z$

◇  $\alpha$  (8);  $\alpha$  (3);  $\alpha$  (2);  $\alpha$  (1);  $\alpha$  (1)

- ▶ Popular replacements for  $Y$

◇  $\alpha$  (12);  $\alpha$  (1);  $\alpha$  (1)

## Summary

- ▶ Expressive power of query languages
  - ◇ Choose a suitably expressive query language for users
  - ◇ This will minimize change of “semantic mistakes”
- ▶ What do users mean by structural constraints?
  - ◇ It's only a hint
  - ◇ Structure is never an inherent part of an information need
- ▶ What sort of query trees are popular?
  - ◇ Non-hierarchical
  - ◇ Including constraints on context

## Micro average

		article	bdy	sec+	p+	abs	vt	bib+	bb	fig	figc	fm
article	12	<b>18.5</b>	<b>10.5</b>	<b>15.2</b>	<b>11.8</b>	0.1	0.5	–	10.5	–	0.5	0.1
sec	20	<b>5.2</b>	<b>5.0</b>	<b>38.6</b>	<b>25.9</b>	1.4	0.2	0.9	0.4	0.4	0.5	0.9
p	5	<b>5.6</b>	<b>5.6</b>	<b>21.1</b>	<b>49.3</b>	3.5	–	–	–	–	2.1	2.8
'*'	2	<b>2.6</b>	<b>2.6</b>	<b>76.9</b>	<b>17.9</b>	–	–	–	–	–	–	–
abs	4	9.3	5.1	37.3	22.9	<b>11.9</b>	–	–	–	0.8	1.7	5.1
vt	2	–	–	–	42.9	–	<b>53.1</b>	–	–	–	2.9	–
bib	1	–	–	–	–	–	–	–	<b>50.0</b>	–	–	–
bb	1	–	–	–	–	–	–	–	<b>95.7</b>	–	–	–
fig	2	–	–	72.3	2.1	–	–	–	–	<b>2.1</b>	<b>2.1</b>	–
fm	1	–	–	–	40.0	–	–	–	–	–	–	–

## Macro Average

		article	bdy	sec+	p+	abs	vt	bib+	bb	fig	figc	fm
article	12	<b>24.7</b>	<b>15.0</b>	<b>29.3</b>	<b>10.5</b>	0.2	0.5	–	2.6	–	0.5	0.6
sec	20	<b>9.6</b>	<b>9.6</b>	<b>39.9</b>	<b>28.1</b>	0.5	0.1	0.2	0.1	0.3	0.2	0.5
p	5	<b>12.0</b>	<b>13.3</b>	<b>23.5</b>	<b>38.6</b>	1.1	–	–	–	–	0.7	0.9
'*'	2	<b>1.7</b>	<b>1.7</b>	<b>57.8</b>	<b>38.9</b>	–	–	–	–	–	–	–
abs	4	11.3	4.0	15.5	22.3	<b>11.8</b>	–	–	–	0.3	6.6	8.5
vt	2	–	–	–	22.0	–	<b>74.7</b>	–	–	–	1.5	–
bib	1	–	–	–	–	–	–	–	<b>50.0</b>	–	–	–
bb	1	–	–	–	–	–	–	–	<b>95.7</b>	–	–	–
fig	2	–	–	66.9	1.2	–	–	–	–	<b>10.0</b>	<b>10.0</b>	–
fm	1	–	–	–	40.0	–	–	–	–	–	–	–