

Range Results in XML Retrieval

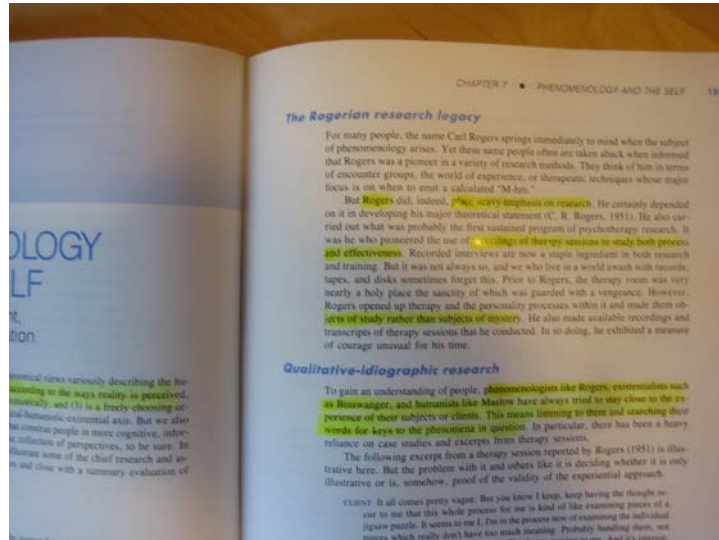
Charles L. A. Clarke
School of Computer Science
University of Waterloo
Waterloo, Canada

Ranges

- “pages 42-57”
- “chapters 2 and 3”
- “sections 4.2 to 4.7”
- “the first two sentences of the third paragraph”

all expressible in XPath

Highlighting



NOT arbitrary passage retrieval

Still want to present the user with
coherent document components:

sections, pages, subsections,
paragraphs, chapters, articles...

Range Results at INEX

INEX 2004 adhoc relevance judgments:

Of the 5229 elements judged as highly exhaustive and specific, at least 1700 (32%) are part of a larger range of elements with identical tag names.

Example: Topic 165

(3,3) elements of ex/2001/x3047.xml:

```
/article[1]/bdy[1]/sec[5]/ss1[1]/p[3]  
/article[1]/bdy[1]/sec[5]/ss1[1]/p[4]  
/article[1]/bdy[1]/sec[5]/ss1[1]/p[5]  
/article[1]/bdy[1]/sec[5]/ss1[1]/p[6]
```

Might be better expressed as:

```
/article[1]/bdy[1]/sec[5]/ss1[1]/p[3 to 6]
```

Range Results in XPath

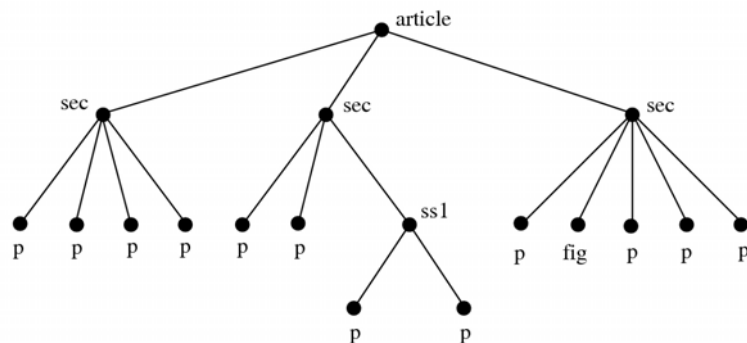
XPath may be too expressive:

```
//sec/p[2]
```

The second paragraph of every section has limited value as a retrieval result.

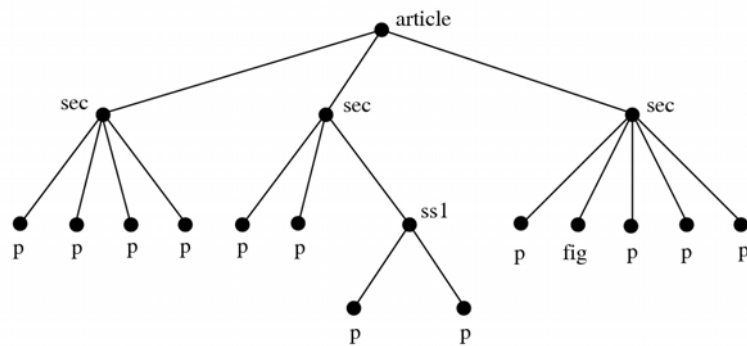
but expressing range results is difficult

First three paragraphs of 1st section



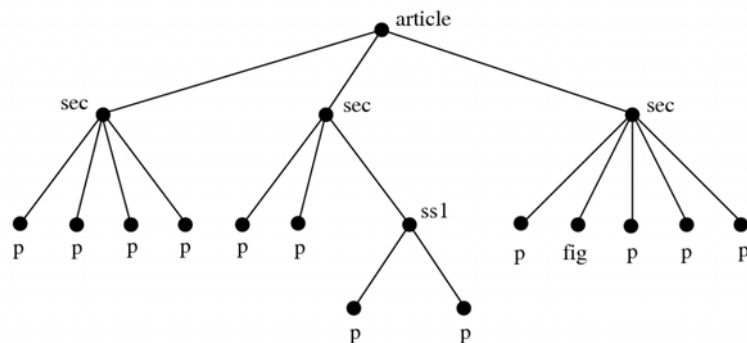
```
/article[1]/section[1]/p[position() <= 3]
```

First three paragraphs of 2nd section



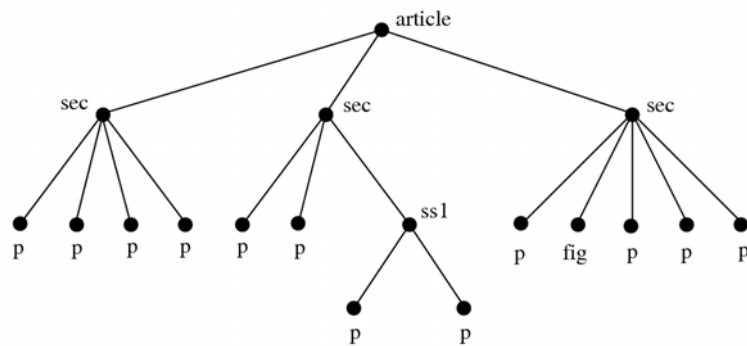
`/article[1]/section[2]/p[position() <= 3]`

First three paragraphs of 2nd section



`/article[1]/section[2]//p[position() <= 3]`

First three paragraphs of 3rd section



```
/article[1]/section[3]//p[position() <= 3]
```

Specifying ranges by endpoints

Given two location paths, X and Y , the XPath 2.0 expression

$X/\text{following}::* [. << Y]$

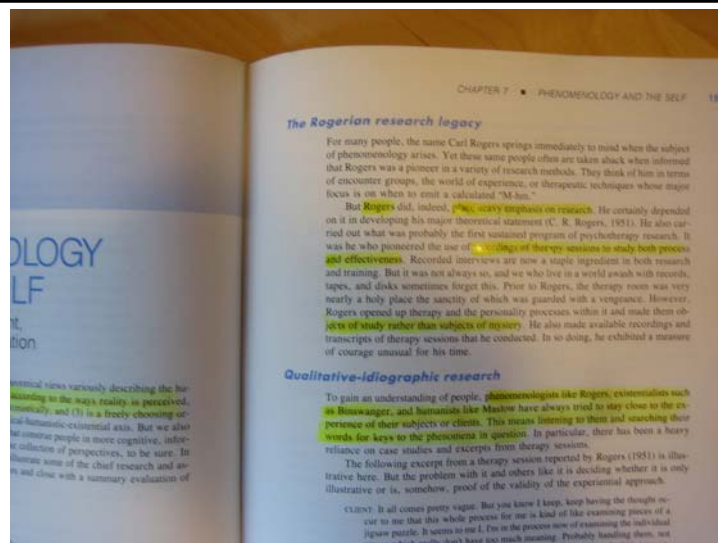
includes all the elements between them.

Results may be specified by (X,Y) pairs.

Evaluation

- $O(N^2)$ range components
- Reusing the test collection?
- Extending INEX methodology?

Evaluation



Evaluation

Given a range component X with relevant text highlighted:

$$\frac{\text{size of relevant text in } X}{\text{size of } X}$$

Specificity?

Evaluation

Given a range component X , contained in document D , with relevant text highlighted:

$$\frac{\text{size of relevant text in } X}{\text{size of relevant text in } D}$$

Exhaustivity? Coherence?

Questions?

- Would range results be useful?
- What other evaluation measures could be used?
- What retrieval techniques could be used?

Discussion

Range Results in XML Retrieval

Charles L. A. Clarke
School of Computer Science
University of Waterloo
Waterloo, Canada