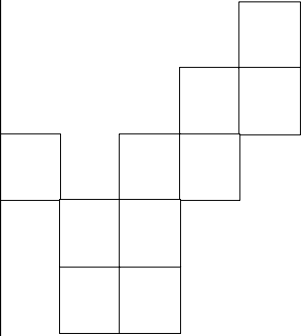
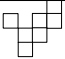


Russian Summer School  
in Information Retrieval



**RuSSIR 2010**  
**XML Retrieval**

Mounia Lalmas  
mounia@acm.org



**Outline**

- **Introduction to XML, basics and standards**
- **Document-oriented XML retrieval**
- **Evaluating XML retrieval effectiveness**
- **Going beyond XML retrieval**



## **Outline**

- **Introduction to XML, basics and standards**
- Document-oriented XML retrieval
- Evaluating XML retrieval effectiveness
- Going beyond XML retrieval



## **Introduction to XML, basics and standards**

- What is XML?
- Database vs. information retrieval
- Document Type Definition
- XML Schema
- Querying XML Data



## XML (eXtensible Markup Language)

- A **meta-language** (a language for describing other languages)  
XML is able to represent a mix of structured and text (unstructured) information
- Defined by the WWW Consortium (**W3C**)  
developed by a W3C working group, headed by James Clark.
- XML 1.0 became a W3C Recommendation on February 10, 1998
- At present XML is the *de facto* standard markup language.



## XML: eXtensible Mark-up Language

- Meta-language (user-defined tags) currently being adopted as the document format language by W3C
- Used to describe **content and structure** (and not layout)
- Grammar described in DTD (→ used for validation)

```
<lecture>
  <title> Structured Document Retrieval </title>
  <author> <fnm> Smith </fnm> <snm> John </snm> </author>
  <chapter>
    <title> Introduction into SDR </title>
    <paragraph> .... </paragraph>
    ...
  </chapter> ...
</lecture>
```

```
<!ELEMENT lecture (title, author+,chapter+)
>
<!ELEMENT author (fnm*,snm)>
<!ELEMENT fnm #PCDATA>
...
```



## XML: eXtensible Mark-up Language

- Use of XPath notation to refer to the XML structure

**chapter/title: title is a direct sub-component of chapter**  
**//title: any title**  
**chapter//title: title is a direct or indirect sub-component of chapter**  
**chapter/paragraph[2]: any direct second paragraph of any chapter**  
**chapter/\*: all direct sub-components of a chapter**

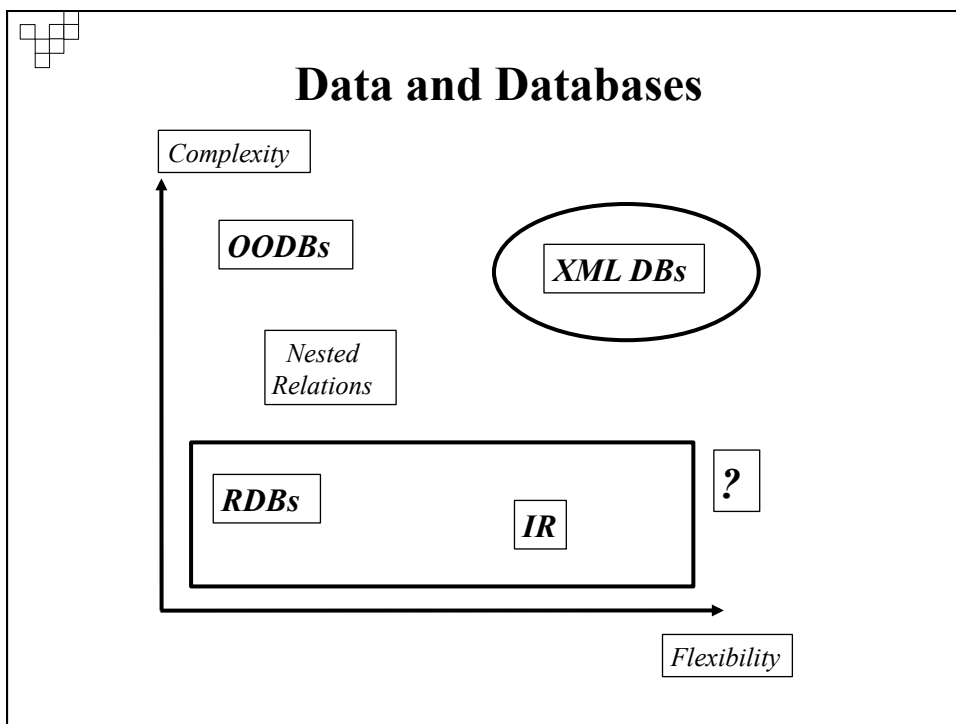
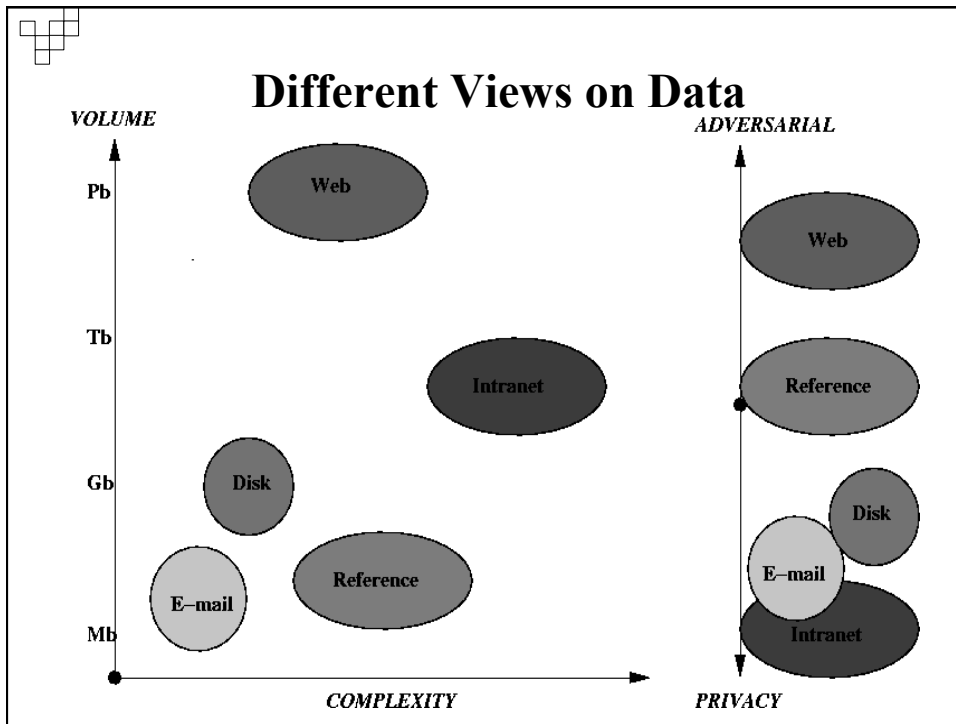
```
<lecture>
  <title> Structured Document Retrieval </title>
  <author> <fnm> Smith </fnm> <snm> John </snm> </author>
  <chapter>
    <title> Introduction into SDR </title>
    <paragraph> .... </paragraph>
    ...
  </chapter> ...
</lecture>
```



## XML

- XML applications: *data interchange, digital libraries, content management, complex documentation, etc.*
- XML repositories: *Library of Congress collection, SIGMOD DBLP, IEEE INEX collection, LexisNexis, ...*

(<http://www.w3.org/XML/>)





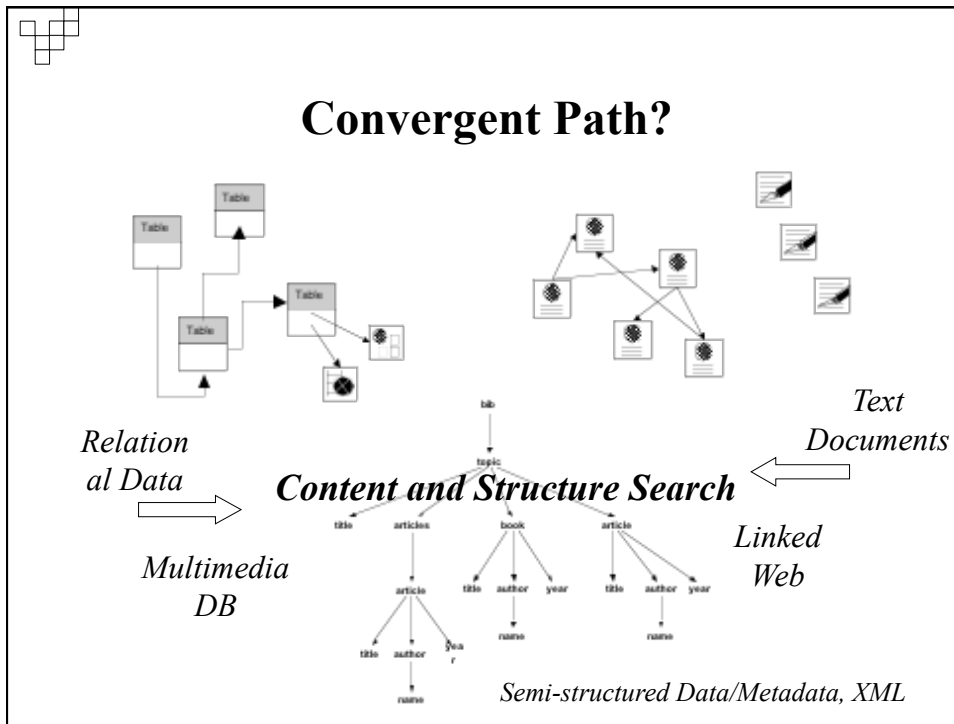
## DB vs. IR

- DBs allow structured querying
- Queries and results (tuples) are different objects
- Soundness & completeness expected
- All results are equally good
- User is expected to know the structure (Enterprise)
- IR only supports unstructured querying
- Queries and results are both documents
- Results are usually imprecise & incomplete
- Some results are more relevant than others
- User is expected to be dumb (Web)



## The Notion of Relevance

- Data retrieval: semantics tied to syntax
- Information retrieval: ambiguous semantics
- Relevance:
  - Depends on the user
  - Depends on the context (task, time, etc)
  - Corollary: The Perfect IR System does not exist



- ### Problems of the IR view
- Very simple query language
    - Is natural language the solution?
  - No query optimization
  - Does not handle the complete answer
  - No types



## Problems of the DB view

- The syndrome of the formal model
  - Model is possible because of structure
- The syndrome of “search then rank”
  - Large answers
  - Optimization is useless
  - Quality vs. Speed
  - E.g. XQuery
- What is a Database?
- Are RDBs really a special case of IR systems?
  - Full text over fields



## DB and IR view

- Data-centric view
  - XML as exchange format for structured data
  - Used for messaging between enterprise applications
  - Mainly a recasting of relational data
- Document-centric view
  - XML as format for representing the logical structure of documents
  - Rich in text
  - Demands good integration of text retrieval functionality
- Now increasingly both views (DB+IR)





## Possible Architectures

- IR on top of RDBs
- IR supported via functions in an RDB
- IR on top of a relational *storage engine*
  
- Middleware layer on top of RDB & IR systems
- RDB functionality on top of an IR system
- Integration via an XML database & query language



## Data-Centric XML Documents: Example

```
<CLASS name="DCS317" num_of_std="100">
  <LECTURER lecid="111">Thomas</LECTURER>
  <STUDENT marks="70" origin="Oversea">
    <NAME>Mounia</NAME>
  </STUDENT>
  <STUDENT marks="30" origin="EU">
    <NAME>Tony</NAME>
  </STUDENT>
</CLASS>
```



## Document-Centric XML Documents: Example

```
<CLASS name="DCS317" num_of_std="100">
  <LECTURER lecid="111">Mounia</LECTURER>
  <STUDENT studid="007" >
    <NAME>James Bond</NAME> is the best student in the
    class. He scored <INTERM>95</INTERM> points out of
    <MAX>100</MAX>. His presentation of <ARTICLE>Using
    Materialized Views in Data Warehouse</ARTICLE> was
    brilliant.
  </STUDENT>
  <STUDENT stuid="131">
    <NAME>Donald Duck</NAME> is not a very good
    student. He scored <INTERM>20</INTERM> points...
  </STUDENT>
</CLASS>
```



## Document-centric XML retrieval

- Documents marked up as XML
  - E.g., assembly manuals, journal issues ...
- Queries are user information needs
  - E.g., give me the Section (element) of the document that tells me how to change a brake light
- Different from well-structured XML queries where one tightly specifies what he/she is looking for.



## XML World

- DTD: Document Type Definition
- XSchema: Data Schema
- DOM: Document Object Model
- SOX: Schema for Object-oriented XML
- Others: XPointer, XSL, ....
  
- XLST: to transform XML
- XPath: to extract XML elements and content
- XQuery: to query XML



## XML

- Documents have **tags** giving extra information about sections of the document  
`<title> XML </title> <slide> Introduction ...</slide>`
- Derived from **SGML** (Standard Generalized Markup Language) but simpler to use
- Extensible, unlike **HTML**  
users can add new tags, and *separately* specify how the tag should be handled for display
- Goal was (is?) to replace HTML as the language for publishing documents on the Web



## XML

- The ability to specify new tags, and to create **nested tag structures** made XML a great way to exchange data, not just documents.

many of the use of XML has been in data exchange applications, and not just a replacement for HTML

- Tags make data **self-documenting**



## Example of an XML document (from database)

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<?xml:stylesheet type="text/xsl" href="staff_list.xsl"?>
<!DOCTYPE STAFFLIST SYSTEM "staff_list.dtd">
<STAFFLIST>
  <STAFF branchNo="B005">
    <STAFFNO>SL21</STAFFNO>
    <NAME>
      <FNAME>John</FNAME><LNAME>White</LNAME>
    </NAME>
    <POSITION>Manager</POSITION>
    <DOB>1-Oct-45</DOB>
    <SALARY>30000</SALARY>
  </STAFF>
  <STAFF branchNo="B003">
    <STAFFNO>SG37</STAFFNO>
    <NAME>
      <FNAME>Ann</FNAME><LNAME>Beech</LNAME>
    </NAME>
    <POSITION>Assistant</POSITION>
    <SALARY>12000</SALARY>
  </STAFF>
</STAFFLIST>
```



## XML - Elements

- **Tag:** label for a section of data
- **Element:** section of data beginning with <tagname> and ending with matching </tagname>
  
- Elements must be properly **nested**
  - Proper nesting  
    <account> ... <balance> .... </balance> </account>
  - Improper nesting  
    <account> ... <balance> .... </account> </balance>
  - Formally: every start tag must have a unique matching end tag that is in the context of the same parent element.
  
- Every document must have a single top-level element



## Example of Nested Elements

```
<bank>
  <customer>
    <customer-name> Monz </customer-name>
    <customer-street> Mile End </customer-street>
    <customer-city> London </customer-city>
    <account>
      <account-number> A-102 </account-number>
      <branch-name> QMUL </branch-name>
      <balance> 400 </balance>
    </account>
  </customer>
  ...
</bank>
```



## XML - Elements

### ■ Mixture of text with sub-elements:

```
<account>  
  This account is seldom used any more.  
  <account-number> A-102</account-number>  
  <branch-name> QMUL</branch-name>  
  <balance>400 </balance>  
</account>
```

- Useful for document markup but discouraged for data representation



## XML - Attributes

### ■ Elements can have attributes

```
<account acct-type = "checking" >  
  <account-number> A-102 </account-number>  
  <branch-name>QMUL </branch-name>  
  <balance> 400 </balance>  
</account>
```

- Attributes are specified by name=value pairs inside the starting tag of an element

- An element may have several attributes, but each attribute name can only occur once

```
<account acct-type = "checking" monthly-fee="5">
```



## XML - Attributes Vs. Elements

- In the context of **documents**, attributes are part of markup, while element contents are part of the basic document contents
- In the context of **data representation**, the difference is unclear and may be confusing
  - `<account account-number = "A-101"> .... </account>`
  - `<account>  
    <account-number>A-101</account-number> ...  
</account>`
- **Suggestion: use attributes for identifiers of elements, and use elements for contents**



## XML – Other Syntax

- Elements **without sub-elements or text content** can be abbreviated by ending the start tag with a `>` and deleting the end tag
  - `<account number="A-101" branch="QMUL" balance="200 />`
- **Comments:** enclosed in `<!--` and `-->` tags.
- **CDATA sections:** instructs XML processor to ignore markup characters and pass enclosed text directly to application.
  - `<![CDATA[<account> ... </account>]]>`



## XML – Ordering

- In XML, elements are ordered.
- In contrast, in XML attributes are unordered.



## Document Type Definition (DTD)

- **Type of an XML document can be specified using a DTD**
- **DTD constraints structure of XML data**
  - What elements can occur?
  - What attributes can/must an element have?
  - What sub-elements can/must occur inside each element, and how many times?
- **DTD does not constrain data types**
  - All values represented as strings in XML
- **DTD syntax**
  - <!ELEMENT element-name (subelements-specification) >
  - <!ATTLIST element-name (attributes) >





## Element Specification in DTD

- **Sub-elements can be specified as**

- names of elements
- #PCDATA (parsed character data), i.e., character strings
- EMPTY (no sub-elements) or ANY (anything can be a sub-element)

- **Example**

```
<!ELEMENT depositor (customer-name account-number)>
<!ELEMENT customer-name (#PCDATA)>
<!ELEMENT account-number (#PCDATA)>
```

- **Sub-element specification may have regular expressions**

```
<!ELEMENT bank ( ( account | customer | depositor)+)>
  “|” - alternatives
  “+” - 1 or more occurrences
  “*” - 0 or more occurrences
  “?” - 0 or 1 occurrence
```



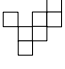
## Attribute Specification in DTD

- For each attribute

- Name
- Type of attribute
  - CDATA
  - ID (identifier) or IDREF (ID reference) or IDREFS (multiple IDREFs)
- Whether
  - mandatory (#REQUIRED)
  - has a default value (value),
  - or neither (#IMPLIED)

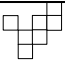
- Examples

```
<!ATTLIST account acct-type CDATA “checking”>
<!ATTLIST customer
  customer-id ID # REQUIRED
  accounts IDREFS # REQUIRED >
```



## DTD Example

mail.dtd	
<pre>&lt;!ELEMENT message   (urgent?, subject,   body)&gt; &lt;!ELEMENT subject   (#PCDATA)&gt; &lt;!ELEMENT body   (ref  #PCDATA)*&gt; &lt;!ELEMENT ref   (#PCDATA)&gt; &lt;!ELEMENT urgent   EMPTY&gt; &lt;!ATTLIST message   date DATE #IMPLIED   sender CDATA #REQUIRED   receiver CDATA #REQUIRED   mtype (TXT MM) ``TXT''&gt;</pre>	<p>Non-XML Language</p> <p>Elements</p> <p>Structure</p> <ul style="list-style-type: none"><li>Sequence</li><li>Nesting</li></ul> <p>Attributes</p>



## Namespaces

- XML data has to be exchanged between organizations
- Same tag name may have different meaning in different organizations, causing confusion on exchanged documents
- Specifying a unique string as an element name avoids confusion
- Better solution: use unique-name:element-name
- Avoid using long unique names all over document by using XML

**Namespaces**

```
<bank xmlns:FB='http://www.FirstBank.com?>
...
  <FB:branch>
    <FB:branchname>Downtown</FB:branchname>
    <FB:branchcity> Brooklyn </FB:branchcity>
  </FB:branch>
...
</bank>
```



## XML Schema

- Database schemas constrain what information can be stored, and the data types of stored values
- XML documents are not required to have an associated schema
- However, schemas are very important for XML data exchange
  - otherwise, a site cannot automatically interpret data received from another site
- Two mechanisms for specifying schema language
  - Document Type Definition (DTD)
    - Widely used
  - XML Schema
    - Newer, increasing use



## XML Schema

- XML Schema is a more sophisticated schema language which addresses the drawbacks of DTDs.
  - Typing of values
    - E.g. integer, string, etc
    - Also, constraints on min/max values
  - User defined types
  - Is itself specified in XML syntax, unlike DTDs
  - Is integrated with namespaces
  - Many more features
    - List types, uniqueness and foreign key constraints, inheritance ..
- BUT: significantly more complicated than DTDs, not yet as widely used.

## XML Schema -Example (from database)

```

<xsd:schema xmlns:xsd=http://www.w3.org/2001/XMLSchema>
<xsd:element name="bank" type="BankType"/>
<xsd:element name="account">
  <xsd:complexType>
    <xsd:sequence>
      <xsd:element name="account-number" type="xsd:string"/>
      <xsd:element name="branch-name" type="xsd:string"/>
      <xsd:element name="balance" type="xsd.decimal"/>
    </xsd:sequence>
  </xsd:complexType>
</xsd:element>
.... definitions of customer and depositor ....
<xsd:complexType name="BankType">
  <xsd:sequence>
    <xsd:element ref="account" minOccurs="0" maxOccurs="unbounded"/>
    <xsd:element ref="customer" minOccurs="0" maxOccurs="unbounded"/>
    <xsd:element ref="depositor" minOccurs="0" maxOccurs="unbounded"/>
  </xsd:sequence>
</xsd:complexType>
</xsd:schema>

```

## XML Schema Example

<pre> &lt;?xml version="1.0"?&gt; &lt;purchaseOrder orderDate="2006-06-20"&gt;   &lt;shipTo country="US"&gt;     &lt;name&gt;Alice Smith&lt;/name&gt;     &lt;street&gt;123 Maple Street&lt;/street&gt;     &lt;city&gt;Mill Valley&lt;/city&gt;     &lt;state&gt;CA&lt;/state&gt;     &lt;zip&gt;90952&lt;/zip&gt;   &lt;/shipTo&gt;   &lt;billTo country="US"&gt;     &lt;name&gt;Robert Smith&lt;/name&gt;     &lt;street&gt;8 Oak Avenue&lt;/street&gt;     &lt;city&gt;Old Town&lt;/city&gt;     &lt;state&gt;PA&lt;/state&gt;     &lt;zip&gt;95819&lt;/zip&gt;   &lt;/bill&gt;   &lt;comment&gt;Hurry, my lawn is going wild!&lt;/comment&gt; </pre>	<pre> &lt;items&gt;   &lt;item partNum="872-AA"&gt;     &lt;productName&gt;Lawnmower     &lt;/productName&gt;     &lt;quantity&gt;1&lt;/quantity&gt;     &lt;USPrice&gt;148.95&lt;/USPrice&gt;     &lt;comment&gt;Confirm this is electric&lt;/comment&gt;   &lt;/item&gt;   &lt;item partNum="926-AA"&gt;     &lt;productName&gt; Baby Monitor     &lt;/productName&gt;     &lt;quantity&gt;1&lt;/quantity&gt;     &lt;USPrice&gt;39.98&lt;/USPrice&gt;     &lt;shipDate&gt;2006-08-21&lt;/shipDate&gt;   &lt;/item&gt; &lt;/items&gt; &lt;/purchaseOrder&gt; </pre>
--	---



## Querying and Transforming XML Data

- Translation of information from one XML schema to another
- Querying on XML data
- Standard XML querying/translation languages
  - XSLT
    - Simple language designed for translation from XML to XML and XML to HTML
  - XPath
    - Simple language consisting of path expressions
  - XQuery
    - An XML query language with a rich set of features
- Wide variety of other languages have been proposed, and some served as basis for the XQuery standard (XML-QL, Quilt, XQL, ...)



## Tree Model of XML Data

- Query and transformation languages based on **tree model** of XML data
- An XML document is modeled as a tree, with nodes corresponding to elements and attributes
  - Element nodes have children nodes, which can be attributes or sub-elements
  - Text in an element is modeled as a text node child of the element
  - Children of a node are ordered according to their order in the XML document
  - Element and attribute nodes (except root node) have a single parent, which is an element node
  - Root node has single child = root element of the document
- Terminology: node, children, parent, sibling, ancestor, descendant.



## XPath

- XPath used to select document parts using path expressions
- Path expression = sequence of steps separated by “/”
- Result of path expression: set of values that along with their containing elements/attributes match the specified path
  
- Examples
  - /bank/customer/customer-name
    - <customer-name>Joe</customer-name>
    - <customer-name>Mary</customer-name>
  - bank/customer/customer-name/text()
    - returns the same names, but without the enclosing tags



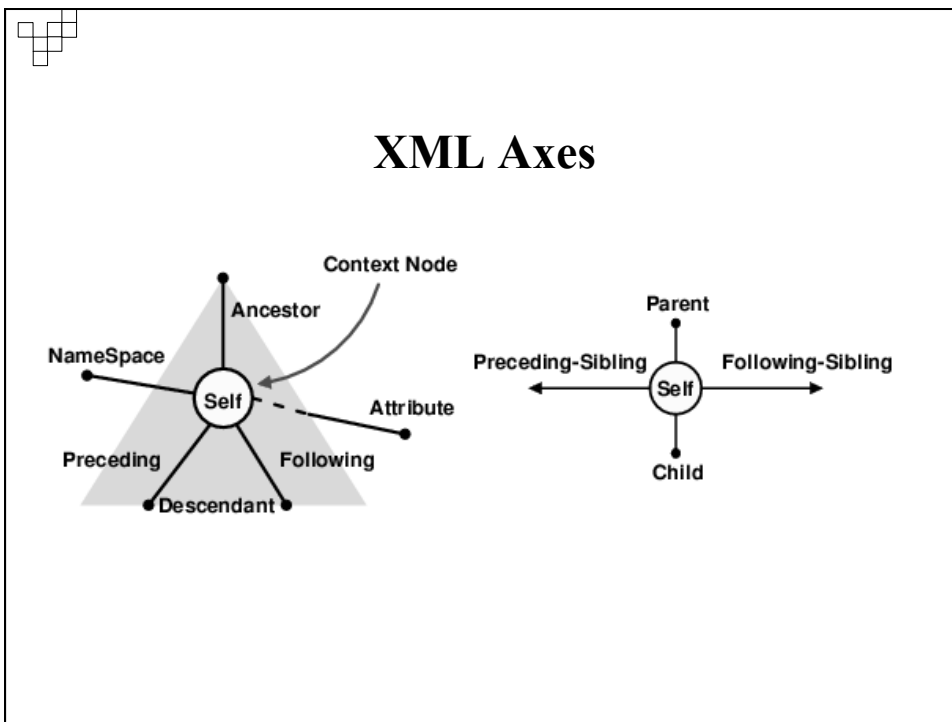
## XPath - Examples

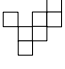
- /bank/account[balance > 400]
  - returns account elements with a balance value greater than 400
- /bank/account[balance]
  - returns account elements containing a balance sub-element
- /bank/account[balance > 400]/@account-number
  - returns the account numbers of those accounts with balance > 400
- /bank/account[customer/count() > 2]
  - returns accounts with > 2 customers

## XPath

**Table 30.2** Some examples of location paths.

Location path	Meaning
.	Selects the context node
..	Selects the parent of the context node
/	Selects the root node, or a separator between steps in a path
//	Selects descendants of the current node
/child::STAFF (or just /STAFF)	Selects all the STAFF elements that are children of the root
child::STAFF (or just STAFF)	Selects the STAFF element children of the context node
attribute::branchNo (or just @branchNo)	Selects the branchNo attribute of the context node
attribute::* (or just @*)	Selects all the attributes of the context node
child::STAFF[3]	Selects the third STAFF element that is a child of the context node
/child::STAFF[@branchNo = "B005"]	Selects all the STAFF elements that have an attribute with a branchNo value of B005
/child::STAFF[@branchNo = "B005"> [position()=1]	Selects first STAFF element that has an attribute with a branchNo value of B005





## Path Expressions

```

<bib>
  <book year="1994">
    <title>TCP/IP
    <author>
      <last>Steve
      <first>W.</
    </author>
    <publisher>A
    <price> 65.95
  </book>

```

**{-- XQuery uses the abbreviated syntax of XPath for path expressions --}**

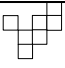
**document("bib.xml")**

**/bib/book/author**

**/bib/book/\*\***

**//author[last="Stevens" and first="W."]**

**document("bib.xml")//author**



## XQuery

- General purpose query language for XML data
- Currently being standardized by World Wide Web Consortium (**W3C**)
- Derived from the Quilt query language, itself based on features from XPath, XML-QL, SQL, OQL, Lorel, XQL, and YATL.



**SQL vs. XQuery**

**"Find item numbers of books"**

**SQL:**

```
SELECT itemno
FROM items AS i
WHERE description LIKE 'Book'
ORDER BY itemno;
```

**XQuery:**

```
FOR $i IN //item_tuple
WHERE contains($i/description, "Books")
RETURN $i/itemno ORDERBY(.)
```

**XQuery – FLWOR Expressions**

■ **FLWOR** (“flower”) expression is constructed from

- **FOR,**
- **LET,**
- **WHERE,**
- **ORDER BY,**
- **RETURN** clauses.

```
graph TD
    A[FOR/LET clauses] --> B["WHERE clause (optional)"]
    B --> C["ORDER BY clause (optional)"]
    C --> D[RETURN clause]
    A --- A1[list of tuples of bound variables]
    B --- B1["restr cted list of tuples of bound variables"]
    C --- C1[sort the tuples in the tuple stream]
    D --- D1["nstance of XML query data model"]
```

The flowchart illustrates the construction of a FLWOR expression. It starts with 'FOR/LET clauses', which produces a 'list of tuples of bound variables'. An optional 'WHERE clause' then filters this into a 'restr cted list of tuples of bound variables'. An optional 'ORDER BY clause' then 'sort the tuples in the tuple stream'. Finally, the 'RETURN clause' produces an 'nstance of XML query data model'.



### **Example - FLWOR Expressions**

**List staff at branch B005 with salary > £15,000.**

```
FOR $$ IN doc("staff_list.xml")//STAFF
WHERE $$/SALARY > 15000 AND
      $$/@branchNo = "B005"
RETURN $$/STAFFNO
```



### **Example - FLWOR Expressions**

**List all staff in descending order of staff number.**

```
FOR $$ IN doc("staff_list.xml")//STAFF
ORDER BY $$/STAFFNO DESCENDING
RETURN $$/STAFFNO
```



## Example - FLWOR Expressions

**List each branch office and average salary at branch.**

```
FOR $B IN distinct-values(doc("staff_list.xml")//@branchNo)
LET $avgSalary :=
  avg(doc("staff_list.xml")//STAFF[@branchNo = $B]/SALARY)
RETURN
  <BRANCH>
    <BRANCHNO>{ $B/text() }</BRANCHNO>,
    <AVGSALARY>$avgSalary</AVGSALARY>
  </BRANCH>
```



## Example - FLWOR Expressions

**List branches that have more than 20 staff.**

```
<LARGEBRANCHES>
  FOR $B IN
    distinct-values(doc("staff_list.xml")//@branchNo)
  LET $$:= doc("staff_list.xml")//STAFF/[@branchNo = $B]
  WHERE count($$) > 20
  RETURN
    <BRANCHNO>{ $B/text() }</BRANCHNO>
</LARGEBRANCHES>
```



## Example – Joining Two Documents

**List staff along with details of their next of kin.**

```
FOR $$ IN doc("staff_list.xml")//STAFF,  
  $NOK IN doc("nok.xml")//NOK  
WHERE $$/STAFFNO = $NOK/STAFFNO  
RETURN  
  <STAFFNO>{ $$, $NOK/NAME }</STAFFNO>
```



## Example – Joining Two Documents

**List all staff along with details of their next of kin.**

```
FOR $$ IN doc("staff_list.xml")//STAFF  
RETURN  
  <STAFFNOK>  
  { $$ }  
  FOR $NOK IN doc("nok.xml")//NOK  
  WHERE $$/STAFFNO = $NOK/STAFFNO  
  RETURN $NOK/NAME  
  </STAFFNOK>
```



## Why XQuery?

- Expressive power
- Easy to learn (?)
- Easy to implement (?)
- Optimizable in many environments
- Related to concepts that people already know
- Several current implementations
- The accepted W3C XML Query Language



## Recap

- Components of the XML World
- Virtues and setbacks of XML Query Language
  - Powerful query language
  - But, too complex for many applications
  - Many implementations
  - Future: XQuery core?
- Any formal background?
  - Structured text models



## XQuery Implementations

- Software AG's Tamino XML Query
- Microsoft, Oracle,
- Lucent Galax
- GMD-IPSI
- X-Hive
- XML Global
- SourceForge XQuench, Saxon, eXist, XQuery Lite
- Fatdog
- Qexo (GNU Kawa) - compiles to Java byte code
- Openlink, CL-XML (Common Lisp), Kweelt,...
- Soda3, DB4XML and about 15 more



## Storing XML documents in databases

- Data centric and document centric XML documents
- Different ways to store XML documents
  - Flat files
  - BLOBs
  - Relational databases
  - Object-Relational databases
  - Native XML databases

<http://www.rpbouret.com/xml/XMLAndDatabases.htm>



## Outline

- Introduction to XML, basics and standards
- **Document-oriented XML retrieval**
- Evaluating XML retrieval effectiveness
- Going beyond XML retrieval



## Document-oriented XML retrieval

- Document vs. data- centric XML retrieval (recall)
- Focused retrieval
- Structured documents
- Structured document (text) retrieval
- XML query languages
- XML element retrieval
- (A bit about) user aspects



## Data-Centric and Document-Centric XML

- Data with partial structure is called **semi-structured**
- XML documents are considered to be **semi-structured**
  
- XML documents classified as:
  - **Data centric**
  - **Document centric**
  
- Nowadays border between data and document centric XML documents is not always clear



## Document-centric XML documents

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<CLASS name="DCS317" num_of_std="100">
  <LECTURER lecid="111">Mounia</LECTURER>
  <STUDENT studid="007" >
    <NAME>James Bond</NAME> is the best student in the
    class. He scored <INTERM>95</INTERM> points out of
    <MAX>100</MAX>. His presentation of <ARTICLE>Using
    Materialized Views in Data Warehouse</ARTICLE> was
    brilliant.
  </STUDENT>
  <STUDENT stuid="131">
    <NAME>Donald Duck</NAME> is not a very good
    student. He scored <INTERM>20</INTERM> points...
  </STUDENT>
</CLASS>
```



## Focused retrieval: Scientific Collection

- **Query**  
model checking  
aviation systems
  
- **Answer**  
one section in a  
workshop report

## Focused Retrieval: Encyclopedia

- **Information need**  
volcanic eruption  
prediction
  
- **Answer**  
relatively small  
portion of the  
volcano topic

## Focused retrieval: Technical Manual

- **Query**  
segmentation fault  
windows services  
for unix
- **Answer**  
only a single  
paragraph in a long  
manual

**SFU-Specific Issues**  
There are some code practices that can cause more problems on SFU's subsystem as compared to traditional UNIX. Specifically, dereferencing an uninitialized pointer can often be unnoticed in a traditional UNIX system but will result in a segmentation fault in SFU. Even when unnoticed this practice leads to subtle problems, and for this reason such dereferencing is considered bad coding and is not advisable for your application. While this may mean some work during the initial part, the resulting application will be more robust.

The security system used in Windows can be stricter than those on traditional UNIX systems. The concept of a user identifier (UID)/group identifier (GID) pair, defined in the `/etc/passwd` and `/etc/group` files, is replaced by a single security identifier (SID) that includes a unique domain-name/login-name pair and domain-name/group-name pair that are both part of a single namespace. As there are no `/etc/passwd` or `/etc/group` files, developers need to handle `struct passwd` differently. For instance, the `getpwnam()` and `getgrnam()` calls accept `domainname+username` pair instead of `username` or `group name` only, and `pw_gecos` is populated with text taken from the Description field in the Windows user account.

**X Windows**  
The Interix SDK includes X11 libraries, header files, and various tools for building X Windows applications. However, SFU does not include an X Windows server, which means that X Windows applications that need to be displayed on the local workstation will need an X Windows server installed. There are a number of excellent X Windows server products available, including a version specifically written for Interix—Interop X Server 8.0 for SFU/Interix from Interop Systems.

Most code written for X Windows assumes a directory structure of `/usr/X11`, but Interix uses a version specific directory structure, `/usr/X11n`, where `n` is replaced with the release level of X11. This difference is best handled by creating a symbolic link to point to the new directory and doesn't require any code changes in applications. Version 3.0 of SFU only included X11R5, but X11R6 is shipping natively in the 3.5 version of Services for UNIX (to be released January 2004) and is already available for download from the Interop Systems Tool Warehouse at <http://www.interopsystems.com/tools/warehouse.htm> for those applications that require a later version of X11 than that originally shipped with SFU 3.0.

**Curses**  
For character mode applications that use curses, the Interix SDK includes the curses implementation of curses written by Eric S. Raymond and Zeyd M. Ben-Halim. This highly compatible and robust implementation of curses is detailed, along with full documentation about writing curses applications and the specifics of the curses implementation, in the SFU Help

## Focused retrieval: Right level of granularity

**Query:** wordnet information retrieval

**OntoSeek: Content-Based Access to the Web**  
Nicola Guarino, Claudio Masolo, Guido Vetere

- article[1]
- fm[1]
- abs[1]
- p[1]
- body[1]
- sec[3] (THE ROLE OF LINGUISTIC ONTOLOGIES)
- ss[1] (Some advantages)
- p[4]
- b[7]
- sec[4] (ONTOSEEK)

... and precision of content-based retrieval. Our OntoSeek system adopts ... large ontology based on **WordNet** for content matching.


The **retrieval** quality improves considerably if ... linguistic ontology such as **WordNet**. For example, let's add **WordNet** to a simple matching ...


... linguistic ontologies such as **WordNet** and structured representation formalisms can help an **information-retrieval** system to

... of a project on **retrieval** and reuse of object-oriented ... system designed for content-based **information retrieval** from online yellow pages ... mostly resulting from merging **WordNet's** thesaurus into the Penman ... broad ontology endowed with **WordNet's** powerful lexical interface, which ...

**Structured Document Retrieval (SDR)**

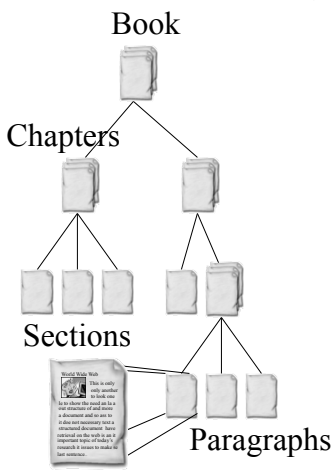
- Traditional IR is about finding relevant documents to a user's information need, e.g. entire book.
 


- SDR allows users to retrieve document components that are more focussed to their information needs, e.g. a chapter of a book instead of an entire book.
 


- The structure of documents is exploited to identify which document components to retrieve.
 

- Structure improves precision
  - Exploit visual memory

**Structured Documents**



- ◆ In general, any document can be considered structured according to one or more structure-type
  - **Linear order of words, sentences, paragraphs ...**
  - **Hierarchy or logical structure of a book's chapters, sections ...**
  - Links (hyperlink), cross-references, citations ...
  - Temporal and spatial relationships in multimedia documents

## Structured Documents

- The structure can be implicit or explicit
- Explicit structure is formalised through document representation standards (Mark-up Languages)
  - Layout
    - LaTeX (publishing), HTML (Web publishing)
  - Structure
    - **SGML, XML (Web publishing, engineering), MPEG-7 (broadcasting)**
  - Content/Semantic
    - RDF (ontology)

## Microformats

- Community data formats
  - Personal Data: hCard (vCard)
  - Calendar and Events: hCal (iCal)
  - Social Networking: XFN
  - Reviews: hReview
  - Licenses: rel-license
  - Folksonomies: rel-tag
- Embedded in XHTML pages and RSS feeds
  - Also RSS Extensions (iTunes, Yahoo! Media, Geo, Google Base, 20+ more in use)



## Example: hCal

```
<strong class="summary">Fashion Expo</strong> in  
<span class="location">Paris, France</span>:  
<abbr class="dtstart" title="2006-10-20">Oct 20</abbr>  
to <abbr class="dtend" title="2006-10-23">22</abbr>
```

- Large and growing list of websites
  - [Eventful.com](#)
  - [LinkedIn](#)
  - [Yedda](#)
  - [upcoming.yahoo.com](#)
  - [Yahoo! Local](#), [Yahoo! Tech Reviews](#)
- Benefit from shared tools, practices (hCalendar creator, iCal Extraction)

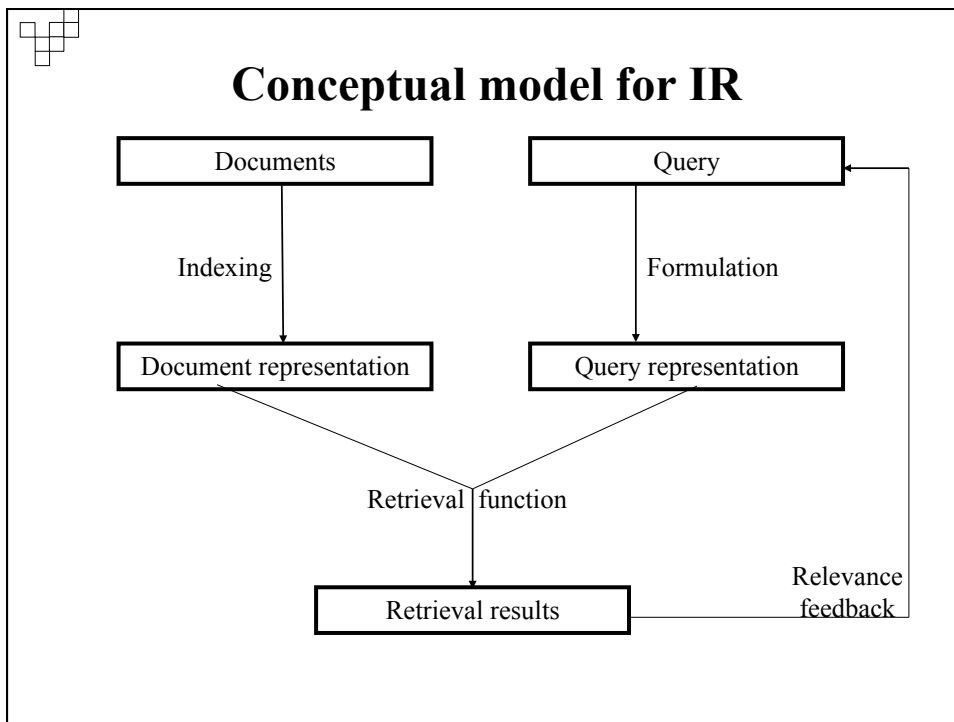


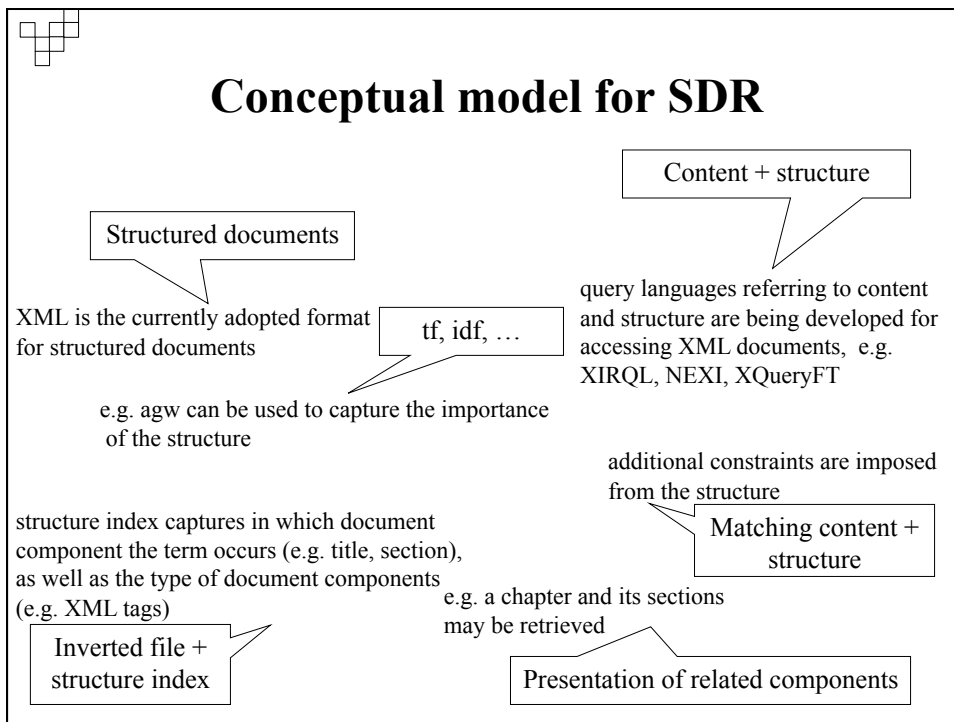
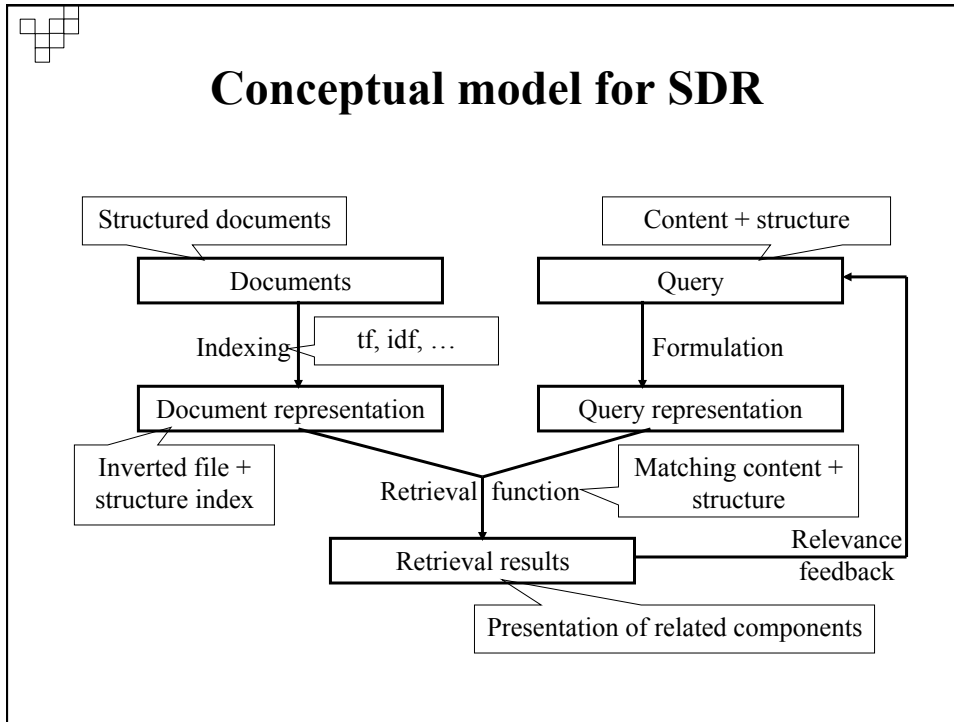
## Queries in SDR

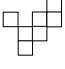
- Three types of queries:
  - Content-only (CO) queries
    - Standard IR queries but here we are retrieving document components
    - “London tube strikes”
  - Structure-only queries
    - Usually not that useful from an IR perspective
    - “Paragraph containing a diagram next to a table”

**Queries in SDR**

- Three types of queries:
  - Content-and-structure (CAS) queries
    - Put on constraints on which types of components are to be retrieved
      - E.g. “Sections of an article in the Times about congestion charges”
      - E.g. Articles that contain sections about congestion charges in London, and that contain a picture of Ken Livingstone, and return titles of these articles”
  - Inner constraints (support elements), target elements

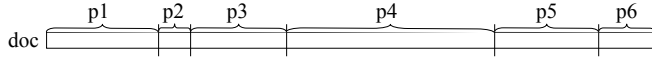






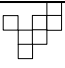
## Passage retrieval

- Passage: continuous part of a document,  
Document: set of passages



- A passage can be defined in several ways:
  - Fixed-length e.g. (300-word windows, overlapping)
  - Discourse (e.g. sentence, paragraph) ← e.g. according to logical structure but fixed (e.g. passage = sentence, or passage = paragraph)
  - Semantic (TextTiling based on sub-topics)
- Apply IR techniques to passages
  - Retrieve passage or document based on highest ranking passage or sum of ranking scores for all passages
  - Deal principally with content-only queries

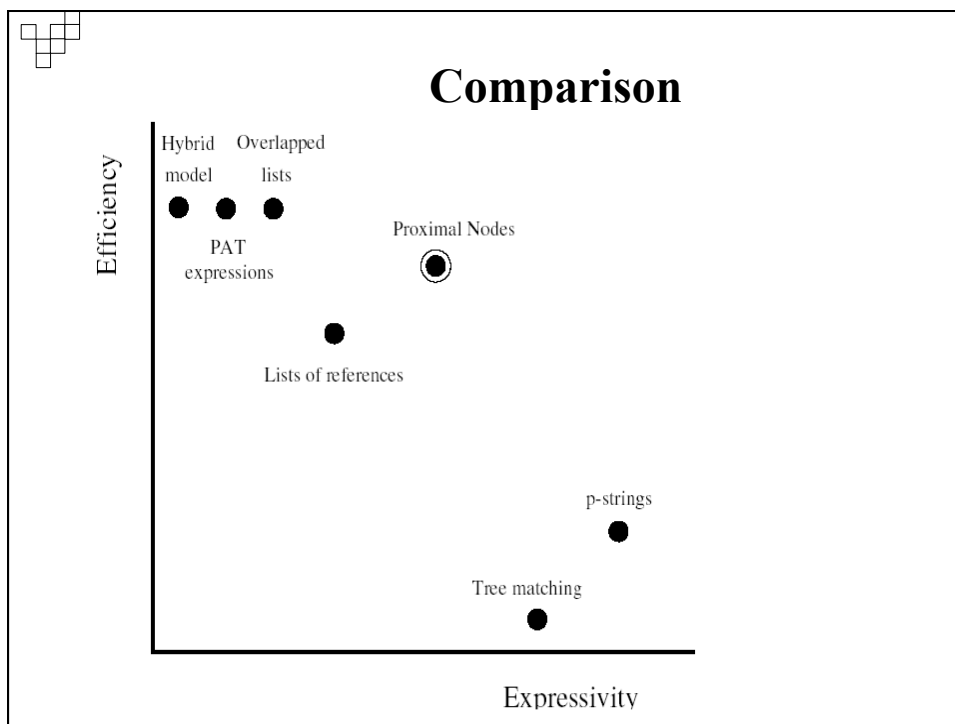
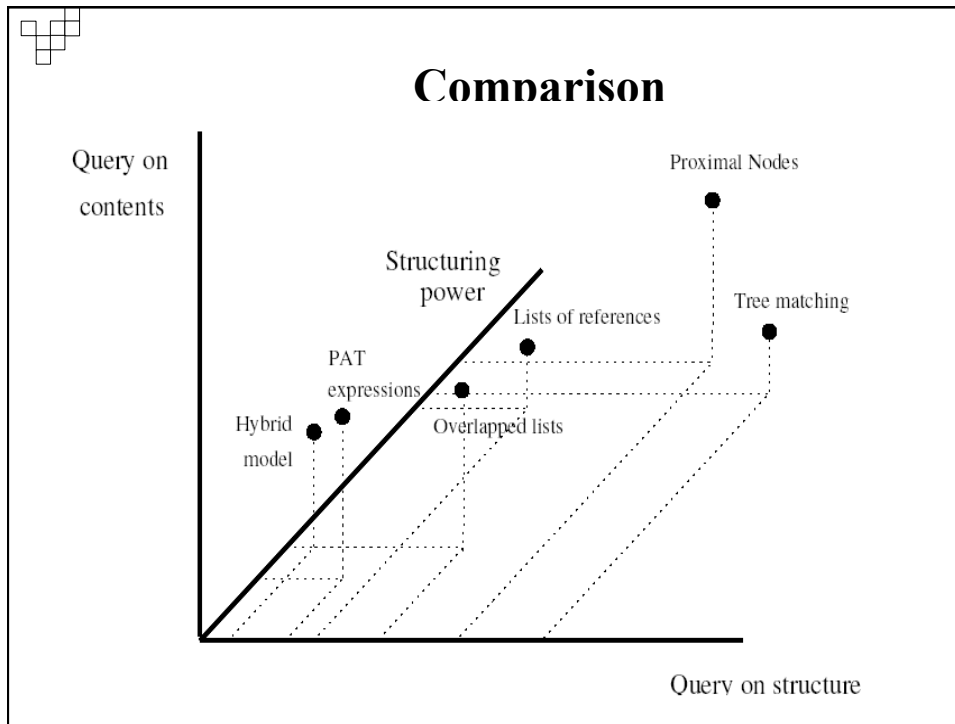
(Callan, SIGIR 1994; Wilkinson, SIGIR 1994; Salton et al, SIGIR 1993; Hearst & Plaunt, SIGIR 1993; ...)



## Structured document (text) retrieval

- Trade-off: expressiveness vs. efficiency
- Models (1989-1995)
  - Hybrid model (flat fields)
  - PAT expressions
  - Overlapped lists
  - Reference lists
  - Proximal nodes
  - Region algebra
    - Proposed as Algebra for XML-IR-DB Sandwich
  - p-strings
  - Tree matching







## Example: Proximal Nodes

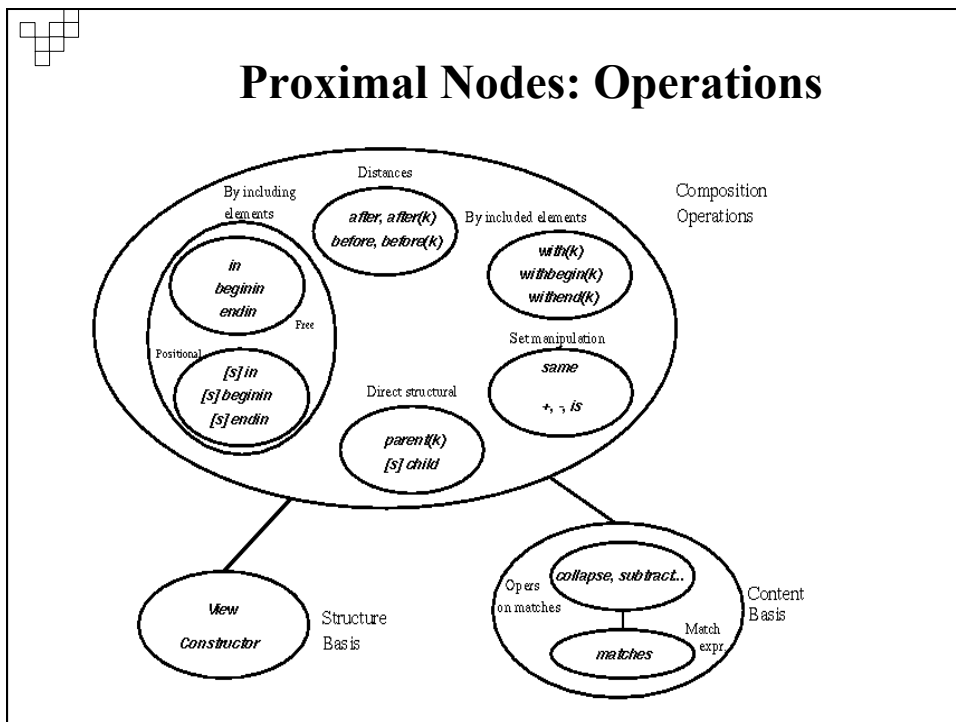
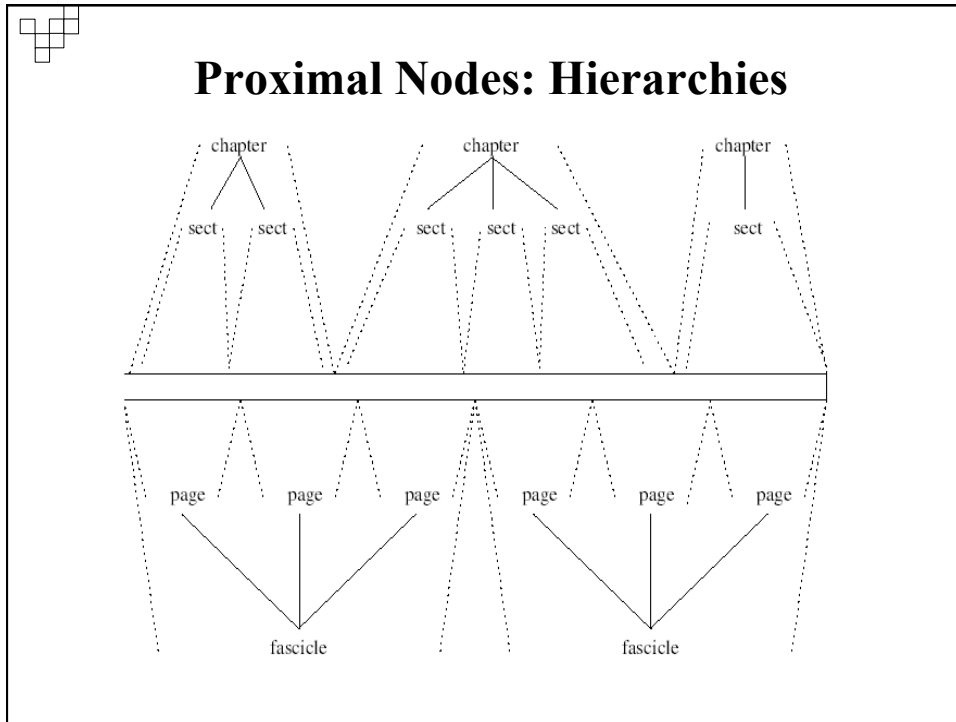
- Hierarchical structure
- Set-oriented language
- Avoid traversing the whole database
- Bottom-up strategy
- Solve leaves with indexes
- Operators work with near-by nodes
- Operators cannot use the text contents
- Most XPath and XQuery expressions can be solved using this model

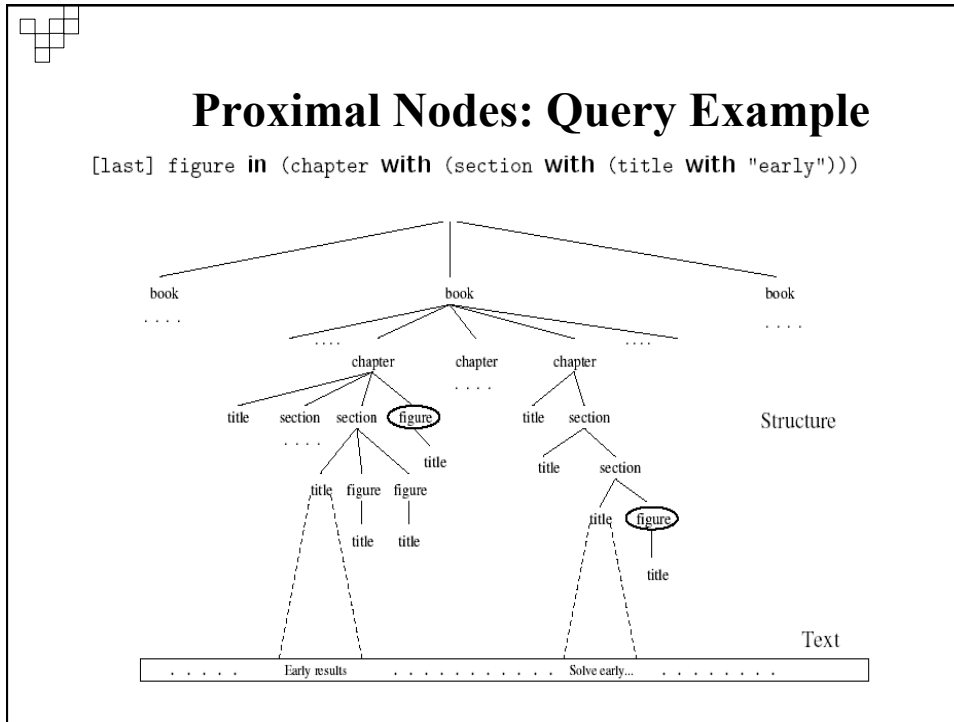
(Navarro & Baeza-Yates, 1995)



## Proximal Nodes: Data Model

- Text = sequence of symbols (filtered)
- Structure = set of independent and disjoint hierarchies or “views”
- Node = Constructor + Segment
- Segment of node  $\subseteq$  segment of children
- Text view, to modelize pattern-matching queries
- Query result = subset of some view





### Region algebra

- Manipulates text intervals - “between which positions in the document?”; and uses containment relationships - “in which components?”
  - Various methods but with similar aims: Simple Concordance List, Generalised Concordance List, Proximal Nodes ...

```

100.1 100.2 101      102      103      103.1
<lecture><title> Structured Document Retrieval </title>
103.2 103.3 104      105 106      107      108      108.1
<chapter><title> Introduction into Structured Document Retrieval </title>
... SDR ...
167.2 167.3
</chapter></lecture>
            
```

- Ranking based on word distances
- Suited for CO and CAS queries

```

Query: “document” and “retrieval”
Intervals: {(102, 103)(107, 108)}

Query: [chapter] containing SDR
Intervals: {(103.2, 167.2)}
            
```

(SIGIR 1992, but see also XML retrieval Mihajlovic etal CIKM 2005)



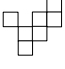
## Query languages for XML

- Four “levels” of expressiveness
  - Keyword search (CO Queries)
    - “xml”
  - Tag + Keyword search
    - book: xml
  - Path Expression + Keyword search (CAS Queries)
    - /book[./title about “xml db”]
  - XQuery + Complex full-text search
    - for \$b in /book  
let score \$s := \$b ftcontains “xml” && “db”  
distance 5



## Query languages for XML

- Keyword search (CO Queries)
  - “xml”
- Tag + Keyword search
  - book: xml
- Path Expression + Keyword search (CAS Queries)
  - /book[./title about “xml db”]
- XQuery + Complex full-text search
  - for \$b in /book  
let score \$s := \$b ftcontains “xml” && “db”  
distance 5

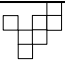


## XRank

```

<workshop date="28 July 2000">
  <title> XML and Information Retrieval: A SIGIR 2000 Workshop </title>
  <editors> David Carmel, Yoelle Maarek, Aya Soffer </editors>
  <proceedings>
    <paper id="1">
      <title> XQL and Proximal Nodes </title>
      <author> Ricardo Baeza-Yates </author>
      <author> Gonzalo Navarro </author>
      <abstract> We consider the recently proposed language ... </abstract>
      <section name="Introduction">
        Searching on structured text is becoming more important with XML ...
        <subsection name="Related Work">
          The XQL language ...
        </subsection>
      </section>
      ...
      <cite xmlns:xlink="http://www.acm.org/www8/paper/xmlql"> ... </cite>
    </paper>
    ...
  
```

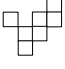
(Guo et al, SIGMOD 2003)



## XRank

```

<workshop date="28 July 2000">
  <title> XML and Information Retrieval: A SIGIR 2000 Workshop </title>
  <editors> David Carmel, Yoelle Maarek, Aya Soffer </editors>
  <proceedings>
    <paper id="1">
      <title> XQL and Proximal Nodes </title>
      <author> Ricardo Baeza-Yates </author>
      <author> Gonzalo Navarro </author>
      <abstract> We consider the recently proposed language ... </abstract>
      <section name="Introduction">
        Searching on structured text is becoming more important with XML ...
        <subsection name="Related Work">
          The XQL language ...
        </subsection>
      </section>
      ...
      <cite xmlns:xlink="http://www.acm.org/www8/paper/xmlql"> ... </cite>
    </paper>
    ...
  
```

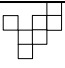


## XIRQL

```

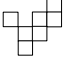
<workshop date="28 July 2000">
  <title> XML and Information Retrieval: A SIGIR 2000 Workshop </title>
  <editors> David Carmel, Yoelle Maarek, Aya Soffer </editors>
  <proceedings>
    <paper id="1">
      <title> XQL and Proximal Nodes </title>
      <author> Ricardo Baeza-Yates </author>
      <author> Gonzalo Navarro </author>
      <abstract> We consider the recently proposed language ... </abstract>
      <section name="Introduction">
        Searching on structured text is becoming more important with XML ...
        <em>The XQL language</em>
      </section>
      ...
      <cite xmlns:xlink="http://www.acm.org/www8/paper/xmlql"> ... </cite>
    </paper>
    ...
  
```

(Fuhr & Großjohann, SIGIR 2001)



## Query languages for XML

- Keyword search (CO Queries)
  - "xml"
- Tag + Keyword search
  - book: xml
- Path Expression + Keyword search (CAS Queries)
  - /book[./title about "xml db"]
- XQuery + Complex full-text search
  - for \$b in /book
    - let score \$s := \$b fcontains "xml" && "db"
    - distance 5



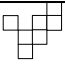
## XSearch

```

<workshop date="28 July 2000">
  <title> XML and Information Retrieval: A SIGIR 2000 Workshop </title>
  <editors> David Carmel, Yoelle Maarek, Aya Soffer </editors>
  <proceedings>
    <paper id="1">
      <title> XQL and Proximal Nodes </title>
      <author> Ricardo Baeza-Yates </author>
      <author> Gonzalo Navarro </author>
      <abstract> We consider the recently proposed language ... </abstract>
      <section name="Introduction">
        Searching on structured text is becoming more important with XML ...
      </section>
    </paper>
    <paper id="2">
      <title> XML Indexing </title>
      (Cohen etal, VLDB 2003)
    </paper>
  </proceedings>
</workshop>

```

**Not a meaningful result**



## Query languages for XML

- Keyword search (CO Queries)
  - "xml"
- Tag + Keyword search
  - book: xml
- Path Expression + Keyword search (CAS Queries)
  - /book[./title about "xml db"]
- XQuery + Complex full-text search
  - for \$b in /book
    - let score \$s := \$b ftcontains "xml" && "db"
    - distance 5





## XPath

- `fn:contains($e, string)` returns true iff `$e` contains string

```
//section[fn:contains(/title, "XML Indexing")]
```

(W3C 2005)



## XIRQL

- Weighted extension to XQL (precursor to XPath)

```
//section[0.6 · //* $cw$ "XQL" +  
0.4 · //section $cw$ "syntax"]
```

(Fuhr & Großjohann, SIGIR 2001)



## XXL

- Introduces similarity operator ~

Select Z

From <http://www.myzoos.edu/zoos.html>

Where zoos.#.zoo As Z and

Z.animals.(animal)?.specimen as A and

A.species ~ “lion” and

A.birthplace.#.country as B and

A.region ~ B.content

(Theobald & Weikum, EDBT 2002)



## NEXI

- Narrowed Extended XPath I
- INEX Content-and-Structure (CAS) Queries
- Specifically targeted for content-oriented XML search (i.e. “aboutness”)

`//article[about(./title, apple) and  
about(./sec, computer)]`

(Trotman & Sigurbjornsson, INEX 2004)



## Query languages for XML

- Keyword search (CO Queries)
  - “xml”
- Tag + Keyword search
  - book: xml
- Path Expression + Keyword search (CAS Queries)
  - /book[./title about “xml db”]
- XQuery + Complex full-text search
  - for \$b in /book  
let score \$s := \$b fcontains “xml” && “db”  
distance 5



## Schema-Free XQuery

- Meaningful least common ancestor (mlcas)

```
for $a in doc(“bib.xml”)//author
  $b in doc(“bib.xml”)//title
  $c in doc(“bib.xml”)//year
where $a/text() = “Mary” and
  exists mlcas($a,$b,$c)
return <result> {$b,$c} </result>
```

(Li et al, VLDB 2003)



## XQuery Full-Text

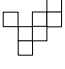
- Two new XQuery constructs
  - 1) **FTContainsExpr**
    - Expresses “Boolean” full-text search predicates
    - Seamlessly composes with other XQuery expressions
  - 3) **FTScoreClause**
    - Extension to FLWOR expression
    - Can score FTContainsExpr *and* other expressions

(W3C 2005)



## FTContainsExpr

```
//book ftcontains “Usability” && “testing” distance 5  
  
//book[./content ftcontains “Usability” with stems]/title  
  
//book ftcontains /article[author=“Dawkins”]/title
```



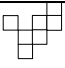
## FTScore Clause

In any order {  
FOR \$v [SCORE \$s]? IN [FUZZY] Expr  
LET ...  
WHERE ...  
ORDER BY ...  
RETURN

Example

FOR \$b SCORE \$s in  
/pub/book[. ftcontains "Usability" && "testing"  
and ./price < 10.00]

ORDER BY \$s  
RETURN \$b



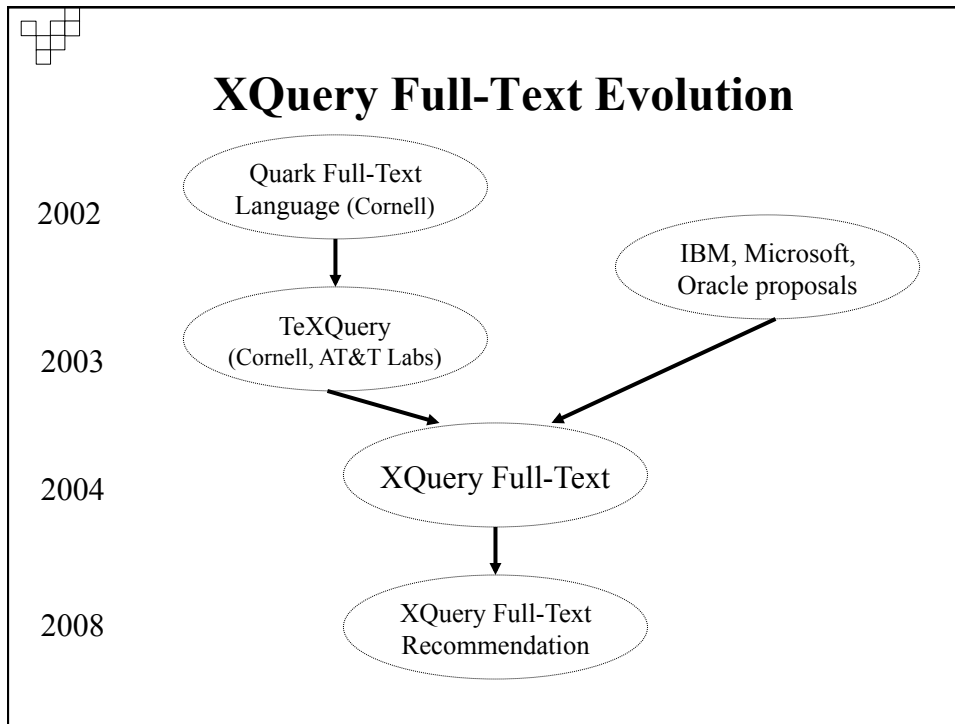
## FTScore Clause

In any order {  
FOR \$v [SCORE \$s]? IN [FUZZY] Expr  
LET ...  
WHERE ...  
ORDER BY ...  
RETURN

Example

FOR \$b SCORE \$s in FUZZY  
/pub/book[. ftcontains "Usability" && "testing"]

ORDER BY \$s  
RETURN \$b



## Query languages for XML - Recap

- Virtues and setbacks of XML query languages
  - Expressive query languages
  - But, too complex for many applications
  - Different interpretations

■ ■ ■

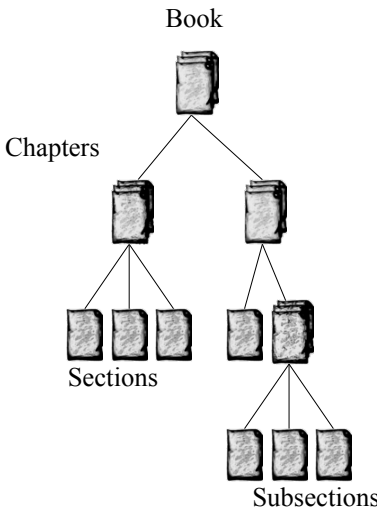
## Element retrieval

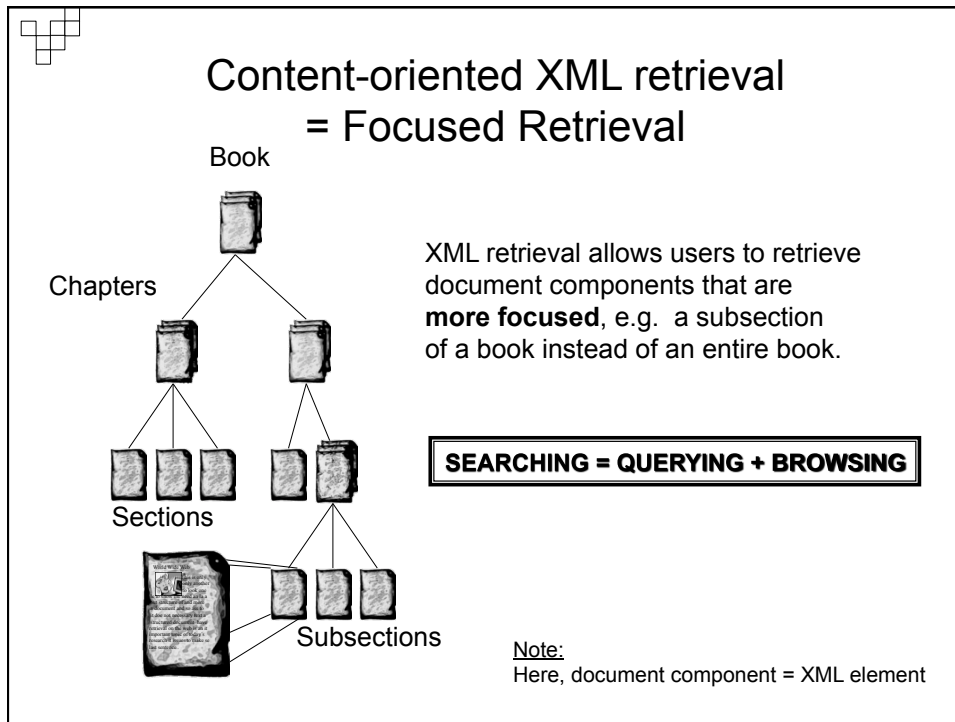
- XML retrieval vs. document retrieval
- XML retrieval = **Focused** retrieval
- Challenges
  1. Term statistics
  2. Relationship statistics
  3. Structure statistics
  4. Overlapping elements
  5. Interpretations of structural constraints
- Ranking
  1. Retrieval units
  2. Combination of evidence
  3. Post-processing – Presentation of XML search results

■ ■ ■

## XML retrieval vs. document retrieval

- No predefined unit of retrieval
- Dependency of retrieval units
- Aims of XML retrieval:
  - Not only to find relevant elements
  - But those at the appropriate level of granularity





**Focused Retrieval for XML: Principle**

- A XML retrieval system should always retrieve the most specific part of a document answering a query.
- Example query: **football**
- Document

```
<chapter> 0.3 football
  <section> 0.5 history </section>
  <section> 0.8 football 0.7 regulation </section>
</chapter>
```
- Return <section>, not <chapter>



**Content-oriented XML retrieval**  
**= Focused Retrieval**

Return document components of **varying granularity** (e.g. a book, a chapter, a section, a paragraph, a table, a figure, etc), relevant to the user's information need both with regards to **content and structure**.

**SEARCHING = QUERYING + BROWSING**

**Challenge 1: Term statistics**

```
graph TD; Article[Article] --> Title[Title]; Article --> Section1[Section 1]; Article --> Section2[Section 2];
```

?XML, ?retrieval ?authoring

0.9 XML  
0.4 retrieval

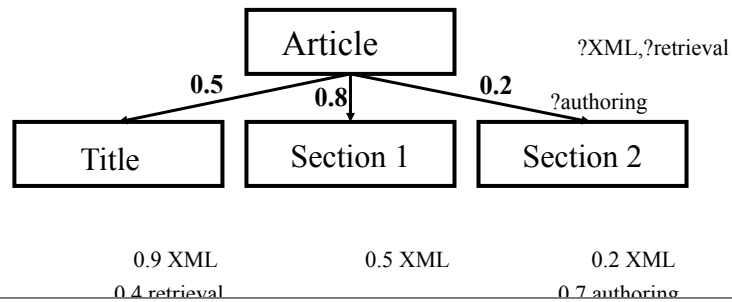
0.5 XML

0.2 XML  
0.7 authoring

**No fixed retrieval unit + nested document components:**

- how to obtain element and collection statistics (e.g. tf, idf)?
- which aggregation formalism to use?
- inner or outer aggregation?

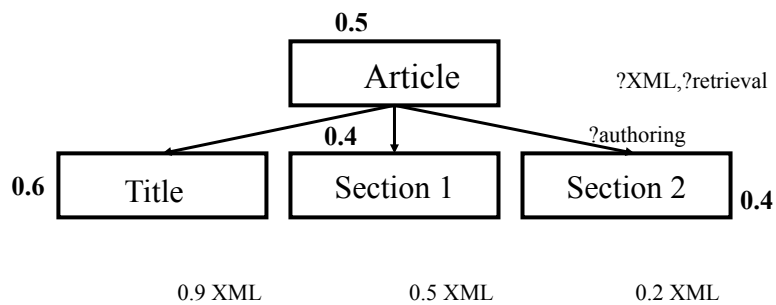
## Challenge 2: Relationship statistics



### Relationship between elements:

- which sub-element(s) contribute best to content of its parent element and vice versa?
- how to estimate (or learn) relationship statistics (e.g. size, number of children, depth, distance)?
- how to aggregate term and/or relationship statistics?

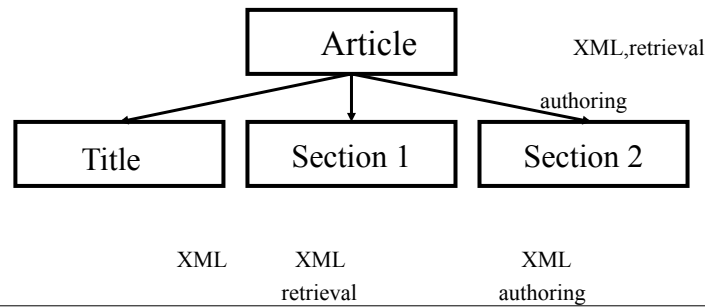
## Challenge 3: Structure statistics



### Different types of elements:

- which element is a good retrieval unit?
- is element size an issue?
- how to estimate (or learn) structure statistics (frequency, user studies, size, depth)?
- how to aggregate term, relationship and/or structure statistics?

## Challenge 4: Overlapping elements

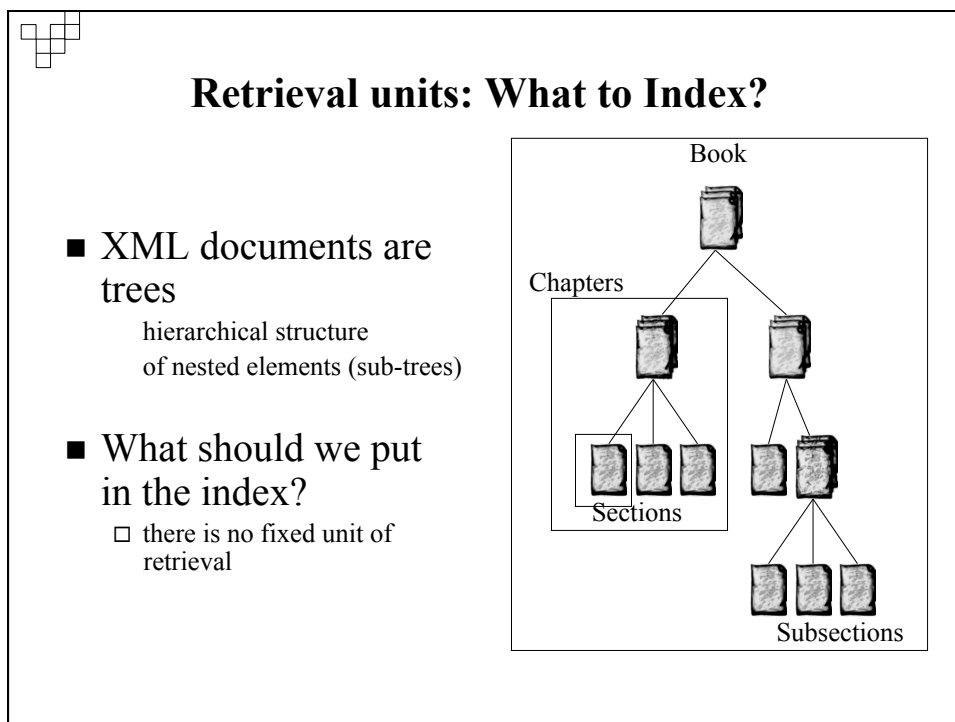
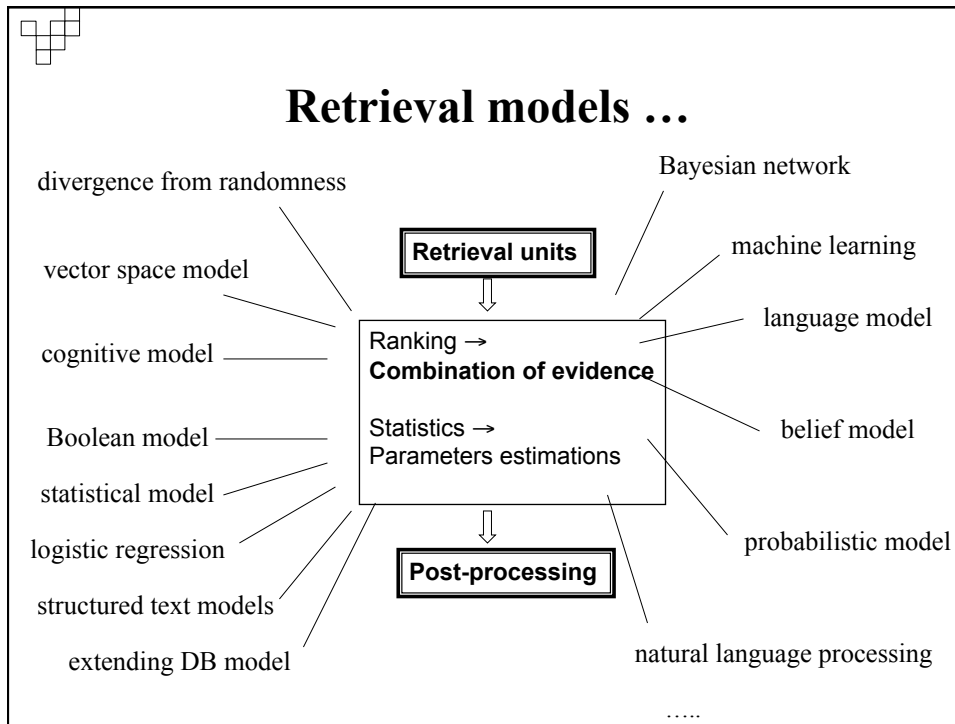


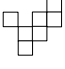
### **Nested (overlapping) elements:**

- section 1 and article are both relevant to “XML retrieval”
- which one to return so that to reduce overlap?
- should the decision be based on user studies, size, types, etc?

## Challenge 5: Expressing and interpreting structural constraints

- Ideally:
  - There is one DTD/schema
  - User understands DTD/schema
- In practice: rare
  - Many DTs/schemas
  - DTDs/Schema not known in advance
  - DTDs/Schema change
  - Users do not understand DTDs/schema
- Need to identify “similar/synonym” elements/tags
- Importance (weight) of tags
- Strict or vague interpretation of the structure
- Relevance feedback/blind feedback?





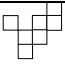
## Retrieval units: XML sub-trees

**Assume a document like**

```
<article>
  <title>XXX</title>
  <abstract>YYY</abstract>
  <body>
    <sec>ZZZ</sec>
    <sec>ZZZ</sec>
  </body>
</article>
```

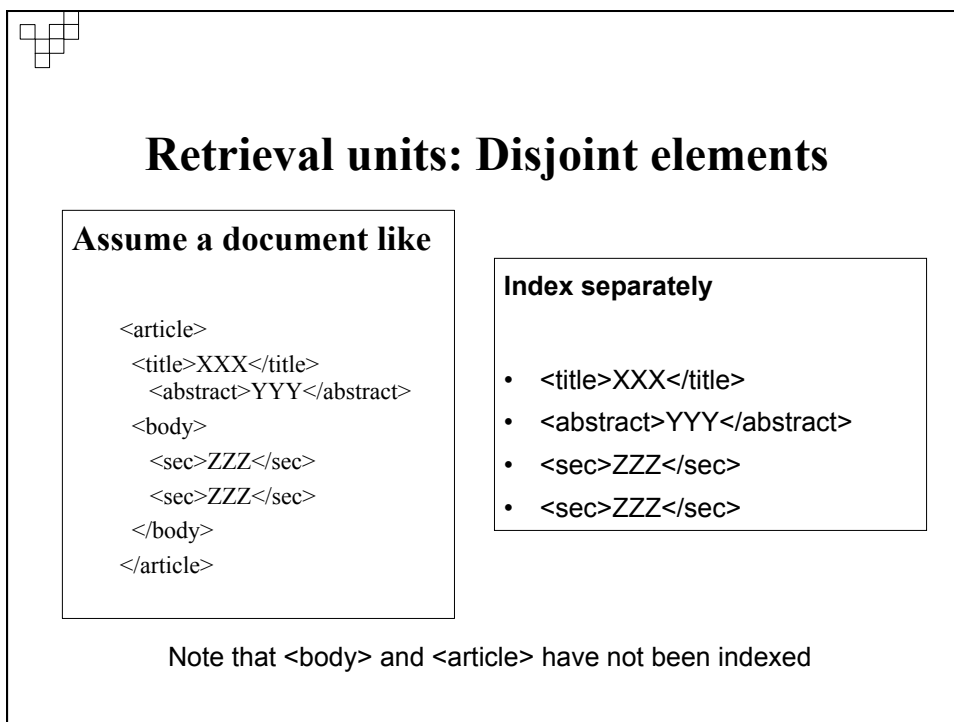
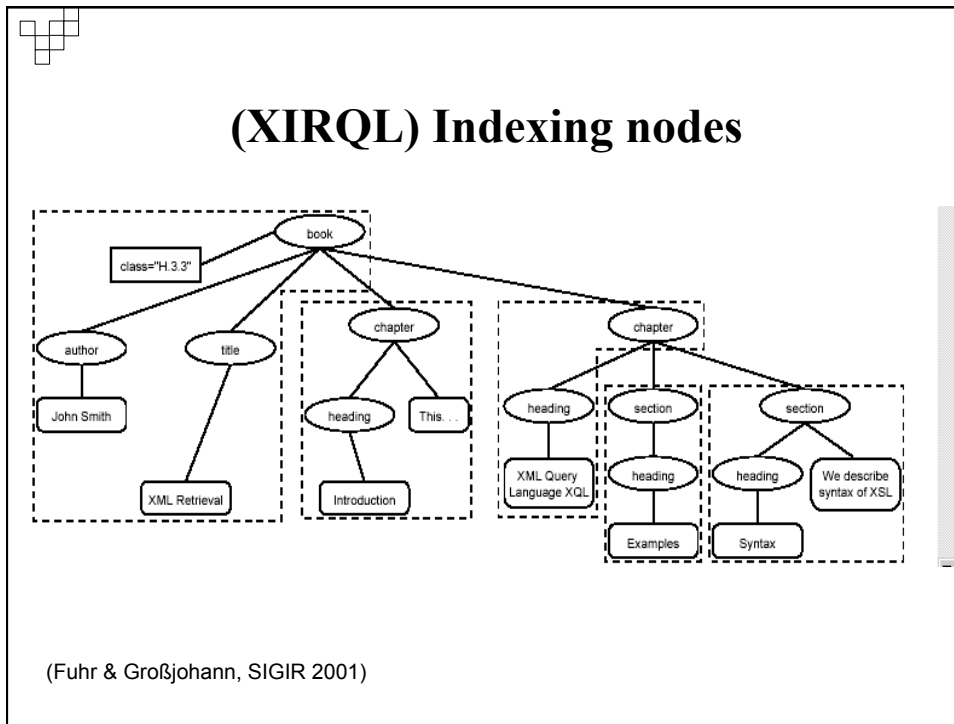
**Index separately**

- <article>XXX YYY ZZZ ZZZ </article>
- <title>XXX</title>
- <abstract>YYY</abstract>
- <body>ZZZ ZZZ</body>
- <sec>ZZZ</sec>
- <sec>ZZZ</sec>



## Retrieval units: XML sub-trees

- Indexing sub-trees is closest to traditional IR
  - each XML elements is bag of words of itself and its descendants
  - and can be scored as ordinary plain text document
  
- Advantage: well-understood problem
  
- Negative:
  - redundancy in index
  - terms statistics
  
  - Led to the notion of indexing nodes
  - Problem: how to select them?
    - manually, frequency, relevance data





## Retrieval units 2: Disjoint elements

- Main advantage and main problem
  - (most) article text is not indexed under /article
  - avoids redundancy in the index
  
- But how to score higher level (non-leaf) elements?
  - Propagation/Augmentation approach
  - Element specific language models



(Geva, INEX 2004, INEX 2005)

## Propagation - GPX model

Leaf elements score

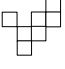
$$L = N^{n-1} \sum_{i=1}^n \frac{t_i}{f_i}$$

$n$  : the number of unique query terms  
 $N$ : a small integer ( $N=5$ , but any  $10 > N > 2$  works)  
 $t_i$  : the frequency of the term in the leaf element  
 $f_i$  : the frequency of the term in the collection

Branch elements score

$$RSV = D(n) \sum_{i=1}^n L_i$$

$n$  : the number of children elements  
 $D(n) = 0.49$  if  $n = 1$   
 $0.99$  Otherwise  
 **$D(n)$  = relationship statistics**  
 $L_i$  : child element score  
 scores are recursively propagated up the tree



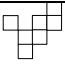
## Element specific language model (simplified)

Assume a document

```
<bdy>
  <sec>cat...</sec>
  <sec>dog...</sec>
</bdy>
```

Query: cat dog  
(Ogilvie & Callan, INEX 2004)

- Assume
  - $P(\text{dog}|\text{bdy}/\text{sec}[1])=0.7$
  - $P(\text{cat}|\text{bdy}/\text{sec}[1])=0.3$
  - $P(\text{dog}|\text{bdy}/\text{sec}[2])=0.3$
  - $P(\text{cat}|\text{bdy}/\text{sec}[2])=0.7$
- Mixture  $P(w|e) = \sum \lambda_i P(w|e_i)$ 
  - With uniform weights ( $\lambda=0.5$ )
  - $\lambda = \text{relationship statistics}$
  - $P(\text{cat}|\text{bdy})=0.5$
  - $P(\text{dog}|\text{bdy})=0.5$
  - So /bdy will be returned



## Retrieval units: Distributed

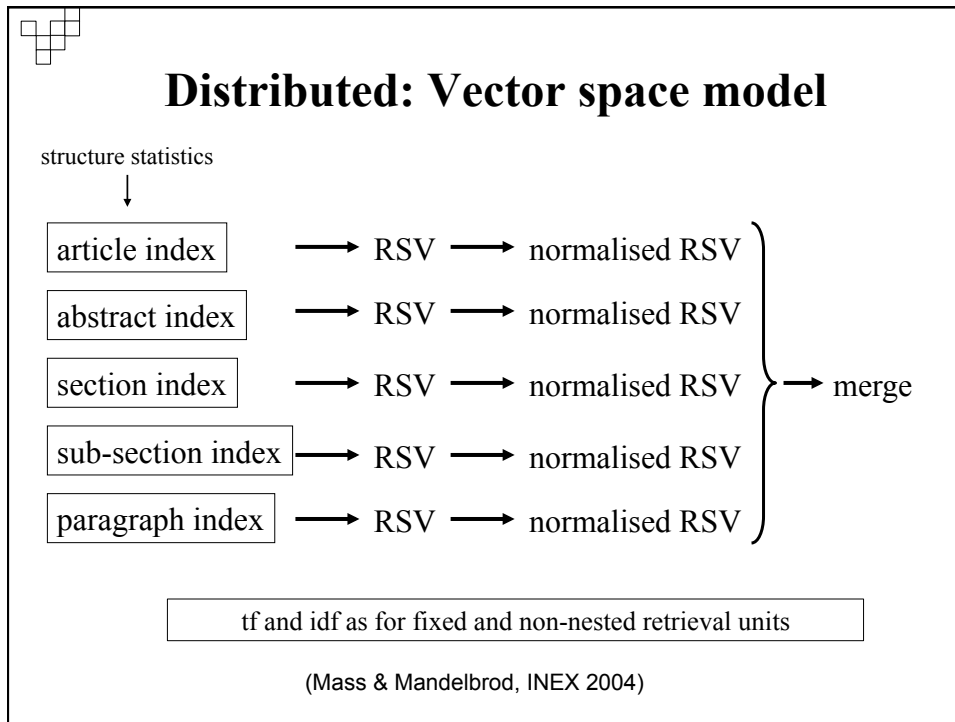
- Index separately particular types of elements
- E.g., create separate indexes for
 

articles  
 abstracts  
 sections  
 subsections  
 subsubsections  
 paragraphs ...

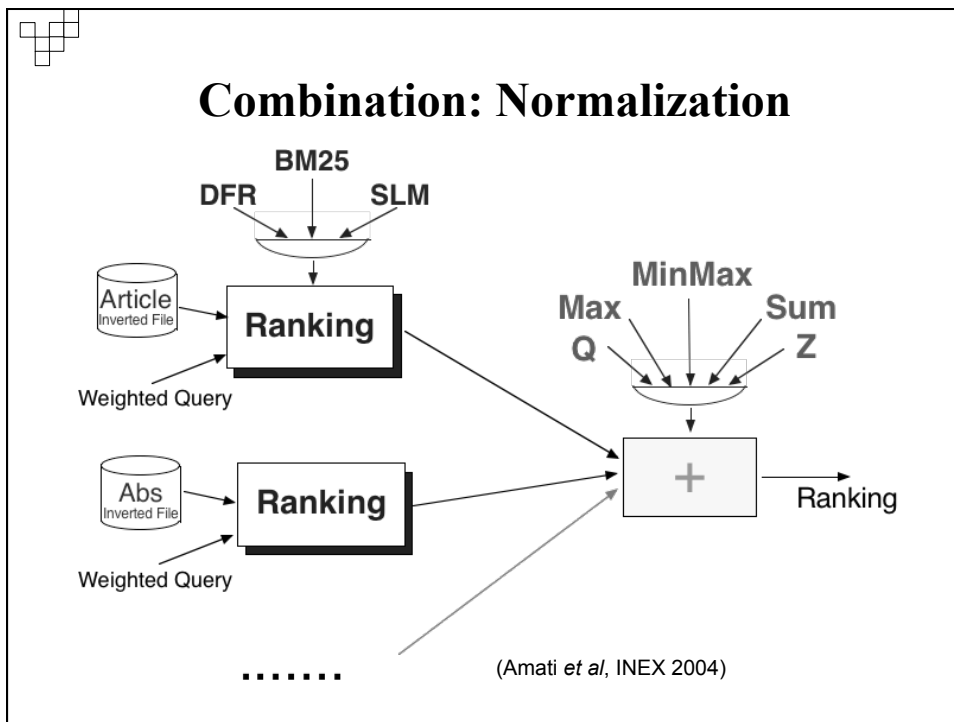
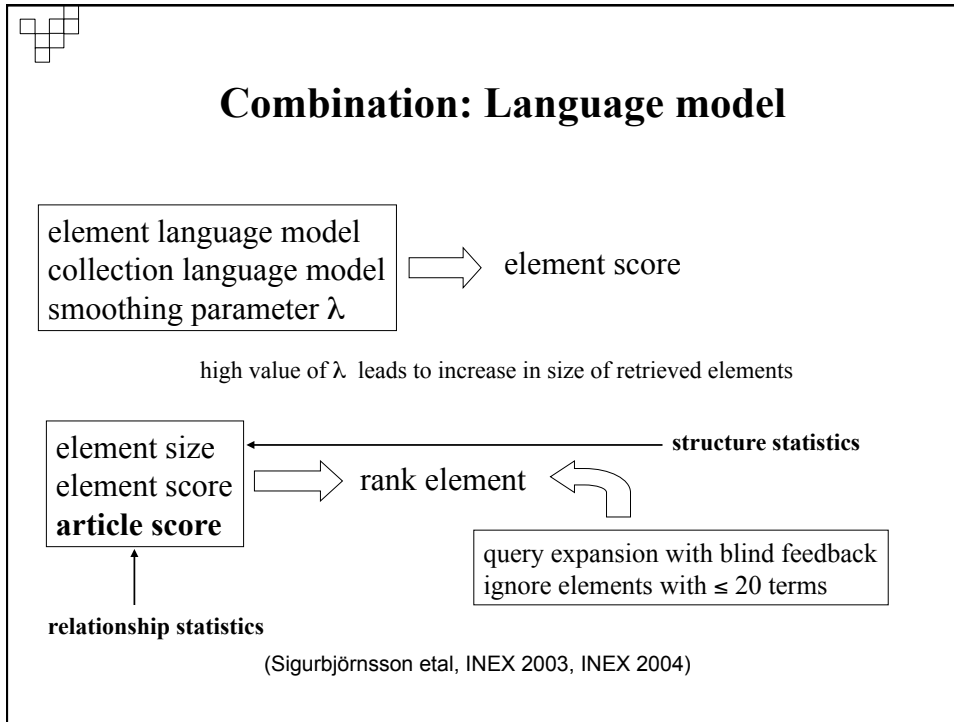
}

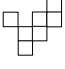
structure statistics
- Each index provides statistics tailored to particular types of elements
  - language statistics may deviate significantly
  - queries issued to all indexes
  - results of each index are combined (after score normalization)





- 
- ### Retrieval units: Distributed
- Only part of the structure is used
    - Element size
    - Relevance assessment
    - Others
  - Main advantages compared to disjoint element strategy:
    - avoids score propagation which is expensive at run-time
    - index redundancy is basically pre-computing propagation
    - XML specific propagation requires nontrivial parameters to train
  - Indexing methods and retrieval models are “standard” IR
    - although issue of merging - normalization





## Combination: Machine learning

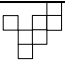
- Use of standard machine learning to train a function that combines
  - Parameter for a given element type
  - Parameter \* score(element)
  - Parameter \* score(parent(element))
  - Parameter \* score (document)

structure statistics

relationship statistics

- Training done on relevance data (previous years)
- Scoring done using OKAPI

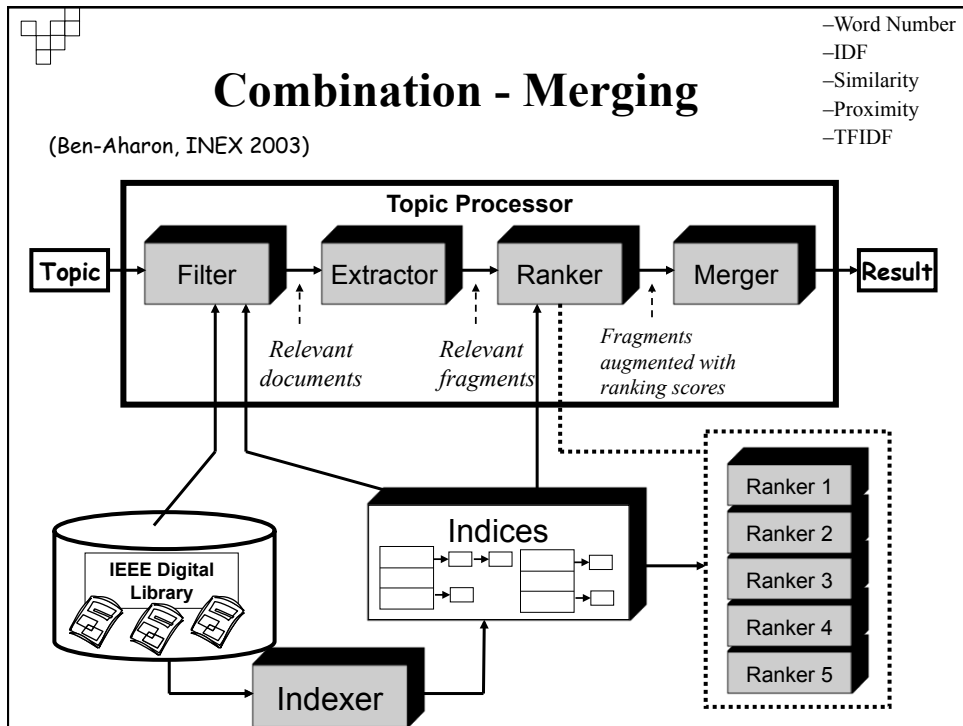
(Vittaut & Gallinari, ECIR 2006)



## Combination: Contextualization

- Basic ranking by adding weight value of all query terms in element.
- Re-weighting is based on the idea of using the ancestors of an element as a context.
  - Root: combination of the weight of an element its 1.5 \* root.
  - Parent: average of the weights of the element and its parent.
  - Tower: average of the weights of an element and all its ancestors.
  - Root + Tower: as above but with 2 \* root.
- Here root is the document

(Arvola et al, CIKM 2005, INEX 2005)



## Presenting XML retrieval results

dbdk\_training in Baseline System

query was: text classification naive bayes  
 Results 1 - 10 of 100.  
 Result pages: 1 2 3 4 5 6 7 8 9 10 next

### Search Result

- 1: (0.247) **Scalable Feature Mining for Sequential Data**  
 Neal Lesh *Mitsubishi Electric Research Lab* Mohammed J. Zaki *Rensselaer Polytechnic Institute* Mitsunori Oghara *University of Rochester*  
 Result path: /article[1]/body[4]/sec[5]
- 2: (0.204) **Probability and Agents**  
 Marco G. Valtorta *University of South Carolina* mgv@cse.sc.edu Michael N. Huhns *University of South Carolina* huhns@sc.edu  
 Result path: /article[1]/body[4]/sec[3]
- 3: (0.176) **Combining Image Compression and Classification Using Vector Quantization**  
 Karen L. Oehler *Member IEEE* Robert M. Gray *Fellow IEEE*  
 Result path: /article[1]/body[4]/sec[4]/ss1[2]/ss2[4]
- 4: (0.175) **Text-Learning and Related Intelligent Agents: A Survey**  
 Dunja Mladenic *J. Stefan Institute*  
 Result path: /article[1]/body[5]/app[4]/sec[5]
- 5: (0.175) **Detecting Faces in Images: A Survey**  
 Ming-Hsuan Yang *Member IEEE* David J. Kriegman *Senior Member IEEE* Narendra Ahuja *Fellow IEEE*  
 Result path: /article[1]/body[4]/sec[2]/ss1[9]/ss2[10]

(Goevert etal, INEX 2002)



## Presenting XML retrieval results

- XML element retrieval is a core task
  - how to estimate the relevance of individual elements
- However, it may not be the end task
  - Simply returning a ranked list of elements results seems insufficient
    - **remove or reduce overlapping elements**
    - elements from the same article may be grouped
    - return one element per article (best entry point from where start to read relevant content)
- Note
  - Presentation and interface



## New retrieval tasks (at INEX)

- INEX 2005 ... addressed new retrieval tasks
  - Thorough is 'pure' XML element retrieval as before
  - **Focused does not allow for overlapping elements to be returned**
  - **Fetch and Browse requires results to be clustered per article**
    - Various variants
  - **Passage retrieval**
- New tasks require post-processing of 'pure' XML element runs
  - geared toward displaying them in a particular interface



## Post-processing: Controlling Overlap

What most approaches are doing:

- Given a ranked list of elements:

1. select element with the highest score within a path
2. discard all ancestors and descendants
3. go to step 1 until all elements have been dealt with

- (Also referred to as brute-force filtering)



## “Post”-Processing: Removing overlap

- Sometimes with some “prior” processing to affect ranking:

- Use of a utility function that captures the amount of useful information in an element

Element score \* Element size \* Amount of relevant information

- Used as a prior probability
- Then apply “brute-force” overlap removal

(Mihajlovic etal, INEX 2005; Ramirez etal, FQAS 2006))



### Post-processing: Controlling Overlap

- Start with a component ranking, elements are re-ranked to control overlap.
- Retrieval status values of those components containing or contained within higher ranking components are iteratively adjusted
- (depends on amount of overlap “allowed”)

1. Select the highest ranking component.
2. Adjust the retrieval status value of the other components.
3. Repeat steps 1 and 2 until the top *m* components have been selected.

(Clarke, SIGIR 2005)



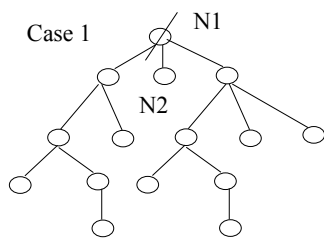
### Post-Processing: Removing overlap

(Mass & Mandelbrod, INEX 2005)

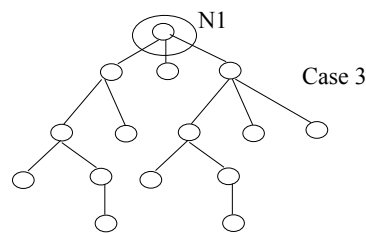
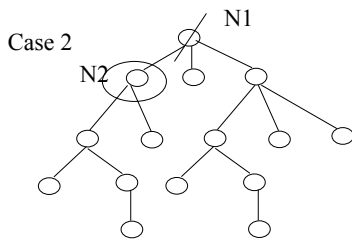
**Smart filtering**

Given a list of rank elements

- group elements per article
- build a result tree
- “score grouping”: ○ ○ ○



- for each element N1
1. score N2 > score N1
  2. concentration of good elements
  3. even distribution of good elements





## CAS query processing: sub-queries

- Sub-queries decomposition
    - *//article [search engines] // sec [Internet growth] AND sec [Yahoo]*
- ↓
- article [search engines]
  - sec [Internet growth]
  - sec [Yahoo]
- 
- Run each sub-queries and then combine
  - Reward structure matching (strict vs vague)

(Sauvagnat etal, INEX 2005)



## Example of combination: Probabilistic algebra

*// article [about(.,bayesian networks)] // sec [about(., learning structure)]*

$$\begin{aligned} &R(\text{learning structure}) \cap \text{label}^{-1}(\text{sec}) \\ &\cap \text{descendants}(R(\text{bayesian networks}) \cap \text{label}^{-1}(\text{article})) \end{aligned}$$

- “Vague” sets
  - $R(\dots)$  defines a vague set of elements
  - $\text{label}^{-1}(\dots)$  can be defined for strict or vague interpretation
- Intersections and Unions are computed as probabilistic “and” and fuzzy-or.

(Vittaut etal, INEX 2004)





### XML Query Relaxation (FlexPath) where DB and IR meet

■ **Tree pattern relaxations:**

- Leaf node deletion
- Edge generalization
- Subtree promotion

**Query**

```

graph TD
    book --> info
    book --> edition
    info --> author
    edition --> paperback
    author --- Dickens
    
```

---

**Data**

```

graph TD
    book --> info
    book --> edition
    info --> author
    edition --> paperback
    author --- Dickens
    
```

```

graph TD
    book --> info
    info --> author
    author --- Dickens
    
```

```

graph TD
    book --> info
    book --> author
    info --> paperback
    author --- Dickens
    
```

(Amer-Yahia, SIGMOD 2004) (Schlieder, EDBT 2002)  
(Delobel & Rousset, 2002) (Amer-Yahia et al, VLDB 2005)

### A Family of XML Scoring Methods

■ **Twig scoring**

- High quality
- Expensive computation

■ **Path scoring**

■ **Binary scoring**

- Low quality
- Fast computation

**Query**

```

graph TD
    book --> info
    book --> edition
    info --> author
    edition --> paperback
    author --- Dickens
    
```

```

graph TD
    book --> info
    book --> edition
    info --> author
    edition --> paperback
    author --- Dickens
    
```

book + book

```

graph TD
    book1 --> info1
    book2 --> edition2
    info1 --> author1
    edition2 --> paperback2
    author1 --- Dickens1
    
```

book + book + book

```

graph TD
    book1 --> info1
    book2 --> info2
    book3 --> edition3
    info1 --> author1
    info2 --> info2
    edition3 --> paperback3
    author1 --- Dickens1
    
```

(Amer-Yahia, VLDB 2005)



## XML Element retrieval - Recap

- Choice of retrieval units can affect the “type” of retrieval models
- XML retrieval can be viewed as a **combination of evidence** problem
- No “clear winner” in terms of retrieval models
  - We still miss the benchmark/baseline approach
  - Lots of heuristics
- BUT WHAT SEEM TO WORK WELL:
  - Element
  - Document
  - Size
- *Thorough investigation for all ranking models, all indexing approaches, and all evidence needed*



## User aspects

- User study - INEX interactive track
- Incorporating user behaviour



## Evaluation of XML retrieval: INEX

- Evaluating the effectiveness of **content-oriented** XML retrieval approaches
- Similar methodology as for TREC, but adapted to XML retrieval  
(to be described later)



## Interactive Track in 2004

- **Investigate behaviour of searchers when interacting with XML components**
- Content-only Topics
  - topic type an additional source of context
    - Background topics / Comparison topics
  - 2 topic types, 2 topics per type
  - 2004 INEX topics have added task information
- Searchers
  - “distributed” design, with searchers spread across participating sites

## Topic Example

```

<title>+new +Fortran +90 +compiler</title>
<description> How does a Fortran 90 compiler differ from a compiler
for the Fortran before it. </description>
<narrative> I've been asked to make my Fortran compiler compatible
with Fortran 90 so I'm interested in the features Fortran 90 added to
the Fortran standard before it. I'd like to know about compilers (they
would have been new when they were introduced), especially
compilers whose source code might be available. Discussion of
people's experience with these features when they were new to them
is also relevant. An element will be judged as relevant if it discusses
features that Fortran 90 added to Fortran. </narrative>
<keywords>new Fortran 90 compiler</keywords>

```

## Baseline system

dbdk\_training in Baseline  
System

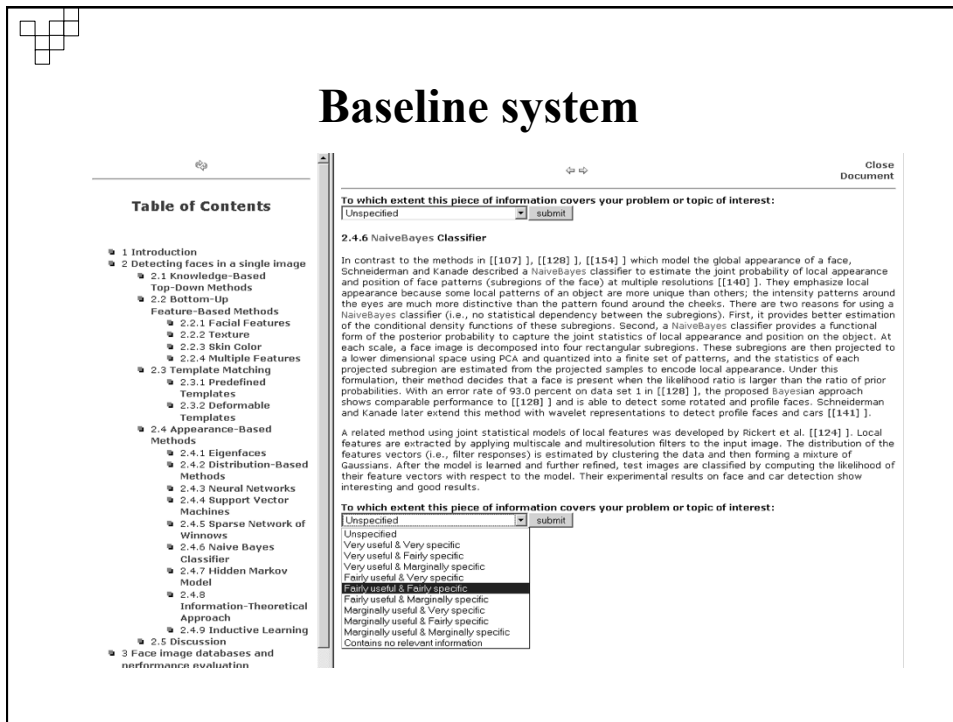
 Search


query was: text classification naive bayes  
Results 1 - 10 of 100.  
Result pages: 1 2 3 4 5 6 7 8 9 10 next



### Search Result

- 1: (0.247) **Scalable Feature Mining for Sequential Data**  
*Neal Lesh Mitsubishi Electric Research Lab Mohammed J. Zaki Rensselaer Polytechnic Institute Mitsunori Ogihara University of Rochester*  
Result path: /article[1]/body[4]/sec[5]
- 2: (0.204) **Probability and Agents**  
*Marco G. Valtorta University of South Carolina mgv@cse.sc.edu Michael N. Huhns University of South Carolina huhns@sc.edu*  
Result path: /article[1]/body[4]/sec[3]
- 3: (0.176) **Combining Image Compression and Classification Using Vector Quantization**  
*Karen L. Oehler Member IEEE Robert M. Gray Fellow IEEE*  
Result path: /article[1]/body[4]/sec[4]/ss1[2]/ss2[4]
- 4: (0.175) **Text-Learning and Related Intelligent Agents: A Survey**  
*Dunja Mladenic J. Stefan Institute*  
Result path: /article[1]/htm[5]/app[4]/sec[5]
- 5: (0.175) **Detecting Faces in Images: A Survey**  
*Ming-Hsuan Yang Member IEEE David J. Kriegman Senior Member IEEE Narendra Ahuja Fellow IEEE*  
Result path: /article[1]/body[4]/sec[2]/ss1[9]/ss2[10]



**Baseline system**

Close Document

**Table of Contents**

- 1 Introduction
- 2 Detecting faces in a single image
  - 2.1 Knowledge-Based Top-Down Methods
  - 2.2 Bottom-Up Feature-Based Methods
    - 2.2.1 Facial Features
    - 2.2.2 Texture
    - 2.2.3 Skin Color
    - 2.2.4 Multiple Features
  - 2.3 Template Matching
    - 2.3.1 Predefined Templates
    - 2.3.2 Deformable Templates
  - 2.4 Appearance-Based Methods
    - 2.4.1 Eigenfaces
    - 2.4.2 Distribution-Based Methods
    - 2.4.3 Neural Networks
    - 2.4.4 Support Vector Machines
    - 2.4.5 Sparse Network of Winnows
    - 2.4.6 Naive Bayes Classifier
    - 2.4.7 Hidden Markov Model
    - 2.4.8 Information-Theoretical Approach
      - 2.4.9 Inductive Learning
    - 2.5 Discussion
  - 3 Face image databases and performance evaluation

To which extent this piece of information covers your problem or topic of interest:

Unspecified submit

**2.4.6 NaiveBayes Classifier**

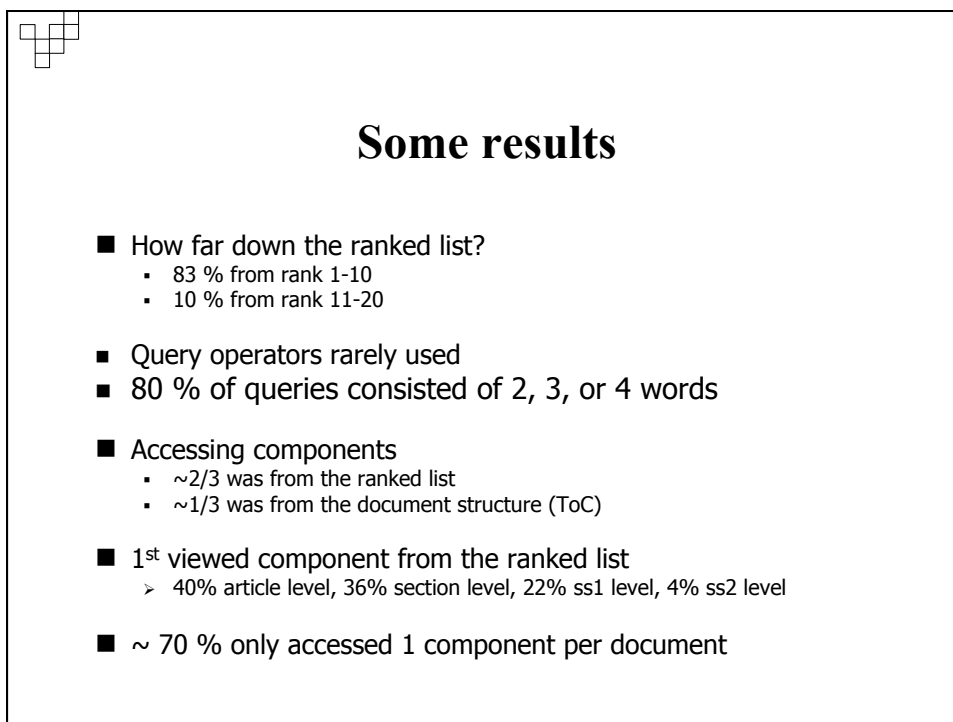
In contrast to the methods in [107], [128], [154] which model the global appearance of a face, Schneideman and Kanade described a NaiveBayes classifier to estimate the joint probability of local appearance and position of face patterns (subregions of the face) at multiple resolutions [140]. They emphasize local appearance because some local patterns of an object are more unique than others; the intensity patterns around the eyes are much more distinctive than the pattern found around the cheeks. There are two reasons for using a NaiveBayes classifier (i.e., no statistical dependency between the subregions). First, it provides better estimation of the conditional density functions of these subregions. Second, a NaiveBayes classifier provides a functional form of the posterior probability to capture the joint statistics of local appearance and position on the object. At each scale, a face image is decomposed into four rectangular subregions. These subregions are then projected to a lower dimensional space using PCA and quantized into a finite set of patterns, and the statistics of each projected subregion are estimated from the projected samples to encode local appearance. Under this formulation, their method decides that a face is present when the likelihood ratio is larger than the ratio of prior probabilities. With an error rate of 93.0 percent on data set 1 in [128], the proposed Bayesian approach shows comparable performance to [128] and is able to detect some rotated and profile faces. Schneideman and Kanade later extend this method with wavelet representations to detect profile faces and cars [141].

A related method using joint statistical models of local features was developed by Rickert et al. [124]. Local features vectors (i.e., filter responses) is estimated by applying multiscale and multiresolution filters to the input image. The distribution of Gaussians. After the model is learned and further refined, test images are classified by computing the likelihood of their feature vectors with respect to the model. Their experimental results on face and car detection show interesting and good results.

To which extent this piece of information covers your problem or topic of interest:

Unspecified submit

- Unspecified
- Very useful & Very specific
- Very useful & Fairly specific
- Very useful & Marginally specific
- Fairly useful & Very specific
- Fairly useful & Fairly specific**
- Fairly useful & Marginally specific
- Marginally useful & Very specific
- Marginally useful & Fairly specific
- Marginally useful & Marginally specific
- Contains no relevant information



**Some results**

- How far down the ranked list?
  - 83 % from rank 1-10
  - 10 % from rank 11-20
- Query operators rarely used
- 80 % of queries consisted of 2, 3, or 4 words
- Accessing components
  - ~2/3 was from the ranked list
  - ~1/3 was from the document structure (ToC)
- 1<sup>st</sup> viewed component from the ranked list
  - 40% article level, 36% section level, 22% ss1 level, 4% ss2 level
- ~ 70 % only accessed 1 component per document



## **Document-centric XML retrieval: Conclusions**

- SDR → now mostly about XML retrieval
- Efficiency:
  - Not just documents, but all its elements
- Models
  - Units
  - Statistics
  - Combination
- User tasks
  
- Link to web retrieval / novelty retrieval
- Interface and visualisation
- Clustering, categorisation, summarisation



## **Outline**

- Introduction to XML, basics and standards
  
- Document-oriented XML retrieval
  
- **Evaluating XML retrieval effectiveness**
  
- Going beyond XML retrieval



## Evaluating XML retrieval effectiveness

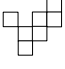
- Structured document retrieval and evaluation
- XML retrieval evaluation
  - Collections
  - Topics
  - Retrieval tasks
  - Relevance and assessment procedures
  - Metrics
- INEX tracks



## SDR and Evaluation

- Passage retrieval
  - Test collection built for that purpose, where passages in relevant documents were assessed (Wilkinson SIGIR 1994)
- Structured document retrieval
  - Web retrieval collection (museum) (Lalmas & Moutogianni, RIAO 2000)
  - Fictitious collection (Roelleke et al, ECIR 2002; Ruthven & Lalmas JDoc 1998)
  - Shakespeare collection (Kazai *et al*, ECIR 2003)
- **INEX initiative** (Kazai *et al*, JASIST 2004; INEX proceedings; SIGIR forum reports, ...)
  - “Real” large test collection following TREC methodology
  - Evaluation campaign
  - XML

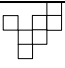




### Evaluation of XML retrieval: INEX

- Evaluating the effectiveness of **content-oriented** XML retrieval approaches
- **Collaborative** effort ⇒ participants contribute to the development of the collection
  - queries
  - relevance assessments
  - methodology
- Similar methodology as for TREC, but adapted to XML retrieval

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



### Document collections

Year	number documents	number elements	size	average number elements	average element depth
2002-2004	12,107	8M	494MB	1,532	6.9
2005	16,819	11M	764MB	“	“
2006-2007	659,388	52M	60 (4.6) GB	161.35	6.72

**IEEE**

---

**Wikipedia**  
(Denoyer & Gallinari, SIGIR Forum, June 2006)

INEX 2009: a new larger collection, based on Wikipedia, with richer set of tags



## Topics

In IR (TREC - <http://trec.nist.gov/>) evaluation, topics are made of:

- Title field:**
  - short explanation of the information need.
- Description field:**
  - one or two sentence natural language definition of the information need.
- Narrative field:**
  - detailed explanation of information need
  - description of what makes something relevant
  - work task it might help to solve
  
- Keywords obtained during collection exploration for the topic creation
- On and off- topic keywords (Amitay *et al*, SIGIR 2004)



## Two types of topics

- **Content-only (CO) topics**
  - ignore document structure
  - simulates users, who do not have any knowledge of the document structure or who choose not to use such knowledge
  
- **Content-and-structure (CAS) topics**
  - contain conditions referring both to content and structure of the sought elements
  - simulate users who do have some knowledge of the structure of the searched collection



## CO topics

```
<title>
  "Information Exchange", +"XML", "Information Integration"
</title>
<description>
  How to use XML to solve the information exchange (information integration) problem,
  especially in heterogeneous data sources?
</description>
<narrative>
  Relevant documents/components must talk about techniques of
  using XML to solve information exchange (information integration)
  among heterogeneous data sources where the structures of participating
  data sources are different although they might use the same ontologies
  about the same content.
</narrative>
```



## CAS topics

```
<title>
  //article[(./fm//yr = '2000' OR ./fm//yr = '1999') AND about(., "'intelligent
  transportation system'")]//sec[about(.,'automation +vehicle')]
</title>
<description>
  Automated vehicle applications in articles from 1999 or 2000 about
  intelligent transportation systems.
</description>
<narrative>
  To be relevant, the target component must be from an article on intelligent
  transportation systems published in 1999 or 2000 and must include a
  section which discusses automated vehicle applications, proposed or
  implemented, in an intelligent transportation system.
</narrative>
```



## NEXI

- Narrowed Extended XPath I
- INEX Content-and-Structure (CAS) Queries
- Specifically targeted for content-oriented XML search (i.e. “**aboutness**”)

**//article[about(./title, apple) and  
about(./sec, computer)]**

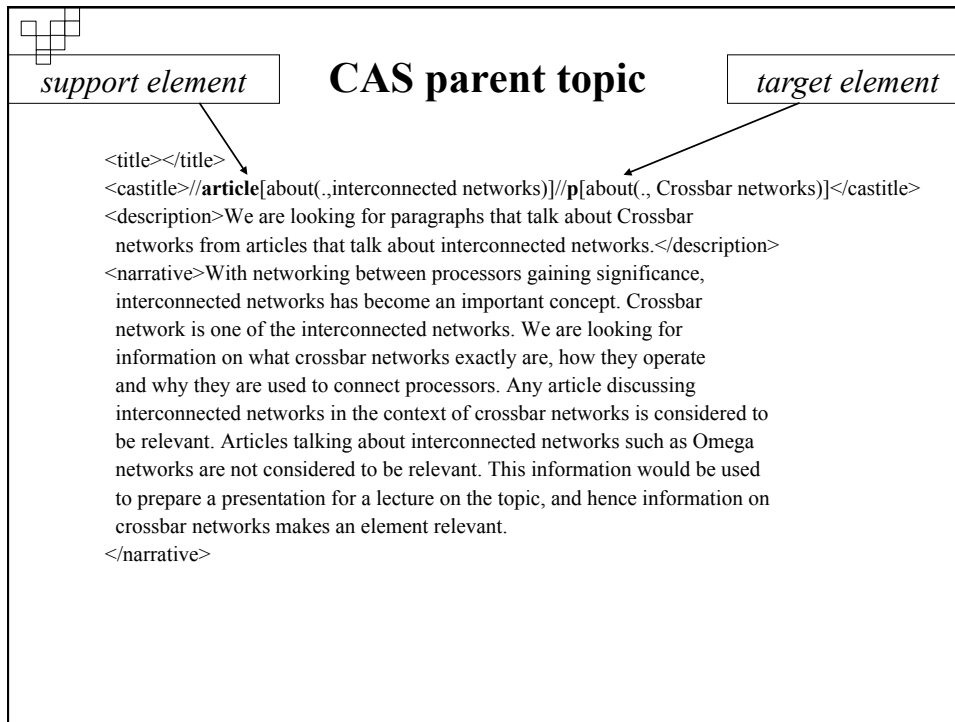
(Trotman & Sigurbjörnsson, INEX 2004)  
(Sigurbjörnsson & Trotman, INEX 2003)



## How to interpret structural constraints?

- Strict vs. vague interpretation of the structure led to:
  - CO+S topics
  - CAS topics

defined in INEX 2005

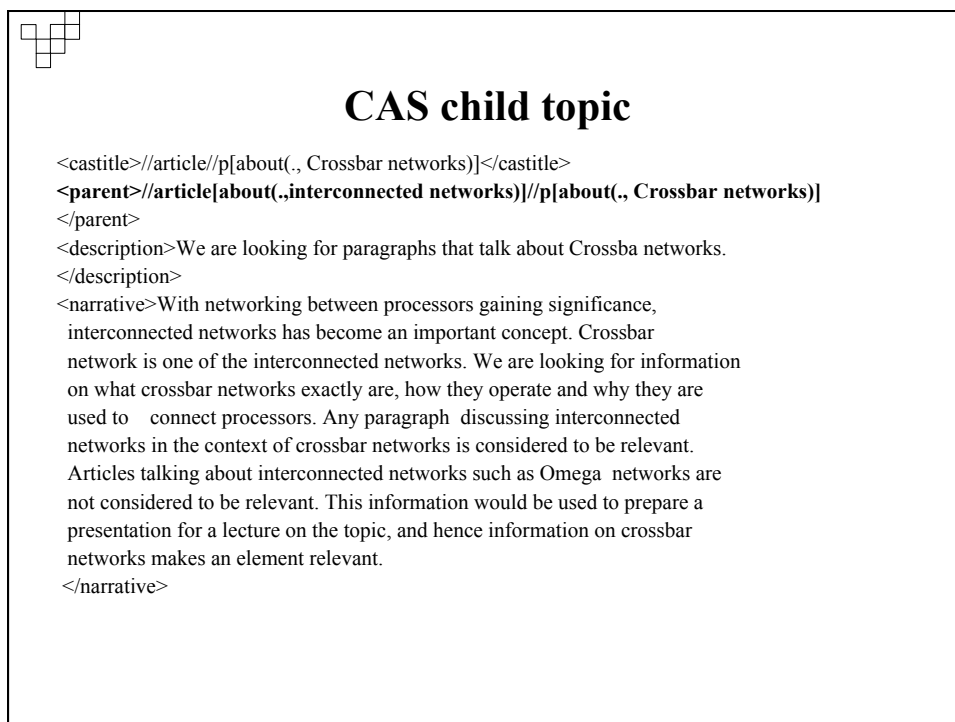


The diagram illustrates the XML structure for a CAS parent topic. It features a central box titled "CAS parent topic" with two labels: "support element" on the left and "target element" on the right. Arrows point from these labels to the corresponding parts of the XML code below. The XML code is as follows:

```

<title></title>
<castitle>//article[about(.,interconnected networks)]//p[about(., Crossbar networks)]</castitle>
<description>We are looking for paragraphs that talk about Crossbar
networks from articles that talk about interconnected networks.</description>
<narrative>With networking between processors gaining significance,
interconnected networks has become an important concept. Crossbar
network is one of the interconnected networks. We are looking for
information on what crossbar networks exactly are, how they operate
and why they are used to connect processors. Any article discussing
interconnected networks in the context of crossbar networks is considered to
be relevant. Articles talking about interconnected networks such as Omega
networks are not considered to be relevant. This information would be used
to prepare a presentation for a lecture on the topic, and hence information on
crossbar networks makes an element relevant.
</narrative>

```



The diagram illustrates the XML structure for a CAS child topic. It features a central box titled "CAS child topic" with XML code below. The XML code is as follows:

```

<castitle>//article//p[about(., Crossbar networks)]</castitle>
<parent>//article[about(.,interconnected networks)]//p[about(., Crossbar networks)]
</parent>
<description>We are looking for paragraphs that talk about Crossba networks.
</description>
<narrative>With networking between processors gaining significance,
interconnected networks has become an important concept. Crossbar
network is one of the interconnected networks. We are looking for information
on what crossbar networks exactly are, how they operate and why they are
used to connect processors. Any paragraph discussing interconnected
networks in the context of crossbar networks is considered to be relevant.
Articles talking about interconnected networks such as Omega networks are
not considered to be relevant. This information would be used to prepare a
presentation for a lecture on the topic, and hence information on crossbar
networks makes an element relevant.
</narrative>

```



## CO+S topics

```

<title>markov chains in graph related algorithms</title>
<castitle>//article//sec[about(.,+"markov chains" +algorithm +graphs)] </castitle>
<description>Retrieve information about the use of markov chains in
graph theory and in graphs-related algorithms.
</description>
<narrative>I have just finished my Msc. in mathematics, in the field
of stochastic processes. My research was in a subject related to
Markov chains. My aim is to find possible implementations of my
knowledge in current research. I'm mainly interested in
applications in graph theory, that is, algorithms related to graphs
that use the theory of markov chains. I'm interested in at
least a short specification of the nature of implementation (e.g.
what is the exact theory used, and to which purpose), hence the
relevant elements should be sections, paragraphs or even abstracts
of documents, but in any case, should be part of the content of the
document (as opposed to, say, vt, or bib).
</narrative>
    
```



## Impact of structural constraints in queries

- Make use of CO+S topics: <castitle>
- **Structural hints:**
  - “Upon discovering that his/her <title> query returned many irrelevant elements, a user might decide to add structural hints, i.e. to write his/her initial CO query as a CAS query”

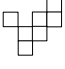
**open standards for digital video in distance learning**



```
//article//sec[about(.,open standards for digital video in distance
learning)]
```

- Results show that processing the structure (from the query) does not consistently lead to any significant improvement in retrieval effectiveness (apart for maybe at very early ranks)

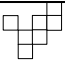
Trotman & Lalmas, SIGIR Poster 2006)



## Retrieval tasks

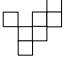
- **Ad hoc retrieval:**
  - “a simulation of how a library might be used and involves the searching of a static set of XML documents using a new set of topics”
  - CO topics
  - CAS (+S) topics
- **Several retrieval strategies**
  - Thorough retrieval
  - Focused retrieval
  - Relevant in context
  - Best in context

*Presentation of results*



## Retrieval strategies


- **Thorough strategy**
  - core system's task underlying most XML retrieval strategies
  - estimate the relevance of all potentially retrievable elements in the collection
  - challenge is to rank elements appropriately, often for further processing
- **Focused strategy:**
  - find the most relevant 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
  - overlapping (redundant) results

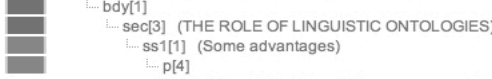



## Fetch & Browse

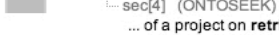
- Document ranking, and in each document, element ranking or set (called Relevant in Context in 2006)

**OntoSeek: Content-Based Access to the Web**  
 Nicola Guarino, Claudio Masolo, Guido Vetere

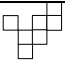

  
 ... and precision of content-based **retrieval**. Our OntoSeek system adopts ... large ontology based on **WordNet** for content matching.


  
 The **retrieval** quality improves considerably if ... linguistic ontology such as **WordNet**. For example, let's add **WordNet** to a simple matching ...


  
 ... linguistic ontologies such as **WordNet** and structured representation formalisms can help an **information-retrieval** system to


  
 ... of a project on **retrieval** and reuse of object-oriented ... system designed for content-based **information retrieval** from online yellow pages ... mostly resulting from merging **WordNet's** thesaurus into the Penman ... broad ontology endowed with **WordNet's** powerful lexical interface, which ...

(Courtesy of Sigurbjörnsson)



## Best in context

- Document ranking, and in each document, return the best entry point
  - Element from where to start reading
  - Analysis:
    - Mostly not the beginning of the document
    - Often the element that is part of the first relevant fragment

(Kamp et al, SIGIR 2007 Poster)





## CO retrieval task

- **Specification:**
  - make use of the CO topics
  - retrieves the most specific elements and only those, which are relevant to the topic
  - no structural constraints regarding the appropriate granularity
  - must identify the most appropriate XML elements to return to the user
  
- **Two main strategies**
  - Thorough strategy
  - Focused strategy



## Thorough strategy

- **Specification:**
  - “core system's task underlying most XML retrieval strategies, which is to estimate the relevance of potentially retrievable elements in the collection”
  
  - overlap problem viewed as an interface and presentation issues
  - challenge is to rank elements appropriately
  
- **Task that most XML approaches performed up to 2004 in INEX.**



## Focused strategy

### ■ Specification:

“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”

- no overlapping elements
- return parent / child if same estimated relevance between parent and child elements
- preference for specificity over exhaustivity



## CAS retrieval task

### ■ Strict content-and-structure:

- retrieve relevant elements that exactly match the structure specified in the query (2002, 2003)

### ■ Vague content-and-structure:

- retrieve relevant elements that may not be the same as the target elements, but are structurally similar (2003)
- retrieve relevant elements even if do not exactly meet the structural conditions; treat structure specification as hints as to where to look (since 2004)

**CAS (+S) retrieval task**

- Make use of CO+S topics: <castitle>
- **Structural hints:**
  - “Upon discovering that his/her <title> query returned many irrelevant elements, a user might decide to add structural hints, i.e. to write his/her initial CO query as a CAS query”

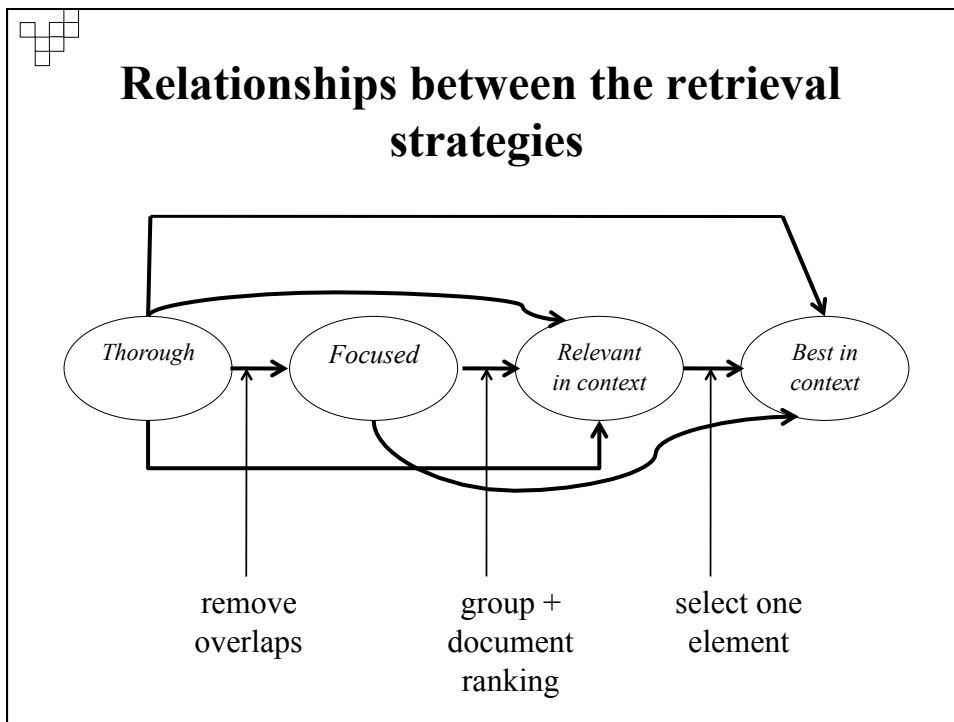
**open standards for digital video in distance learning**

↓

**//article//sec[about(.,open standards for digital video in distance learning)]**

- Two strategies (as for CO retrieval task):
  - *Focussed strategy*
  - *Thorough strategy*

(Trotman & Lalmas, SIGIR Poster 2006)



**Relevance in XML retrieval**

- A document is **relevant** if it “has significant and demonstrable bearing on the matter at hand”.
- Common assumptions in laboratory experimentation:
  - **Objectivity**
  - **Topicality**
  - **Binary nature**
  - **Independence**

(Borlund, JASIST 2003)  
(Goevert et al, JIR 2006)

```

graph TD
    article((article)) --- s1((s1))
    article --- s2((s2))
    article --- s3((s3))
    s3 --- ss1((ss1))
    s3 --- ss2((ss2))
    
```

**Relevance in XML retrieval**

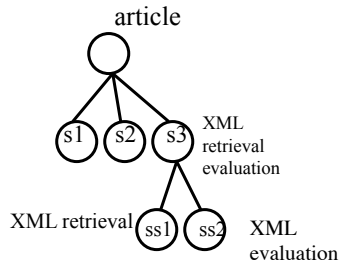
- **Topicality not enough**
- **Binary nature not enough**
- **Independence is wrong**

**specificity:** extent to which a document component is focused on the information need, while being an relevant unit.



(based on Chiramella et al, FERMI fetch and browse model 1996)

## Relevance in XML retrieval: INEX 2003 - 2004



- **Topicality not enough**
- **Binary nature not enough**
- **Independence is wrong**

- **Relevance = (0,0) (1,1) (1,2) (1,3) (2,1) (2,2) (2,3) (3,1) (3,2) (3,3)**

exhaustivity = how much the section discusses the query: 0, 1, 2, 3

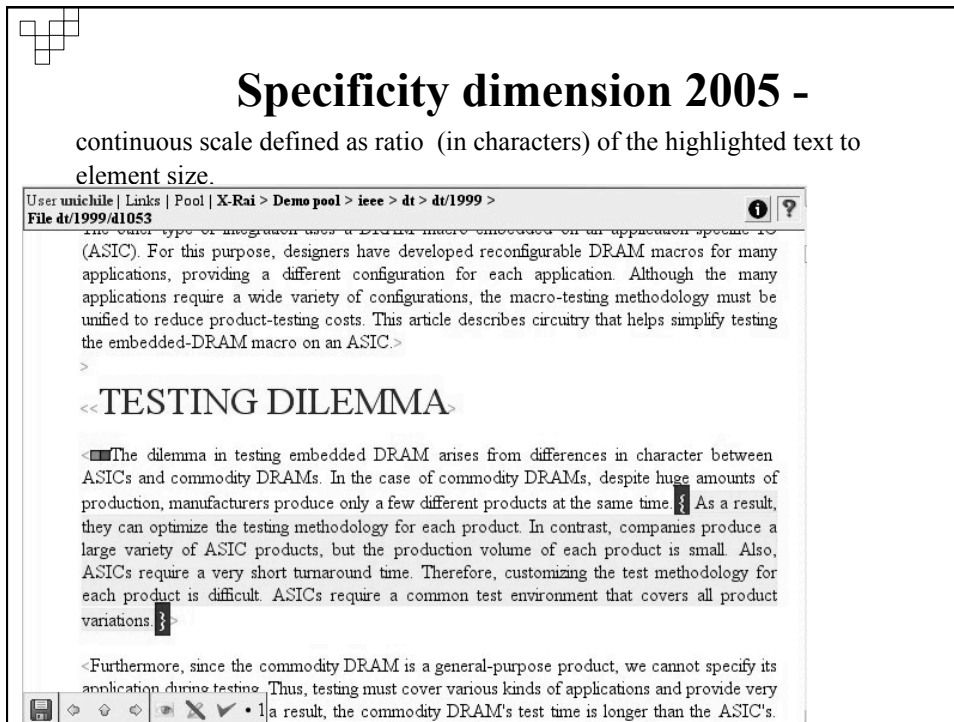
specificity = how focused the section is on the query: 0, 1, 2, 3

- **If a subsection is relevant so must be its enclosing section, ...**



## Relevance - to recap

- find smallest component (→ specificity) that is highly relevant (→ exhaustivity)
- **specificity**: extent to which a document component is focused on the information need, while being an informative unit.
- **exhaustivity**: extent to which the information contained in a document component satisfies the information need.



**Specificity dimension 2005 -**  
continuous scale defined as ratio (in characters) of the highlighted text to element size.

User: unichile | Links | Pool | X-Rai > Demo pool > ieeec > dt > dt/1999 >  
File: dt/1999/d1053

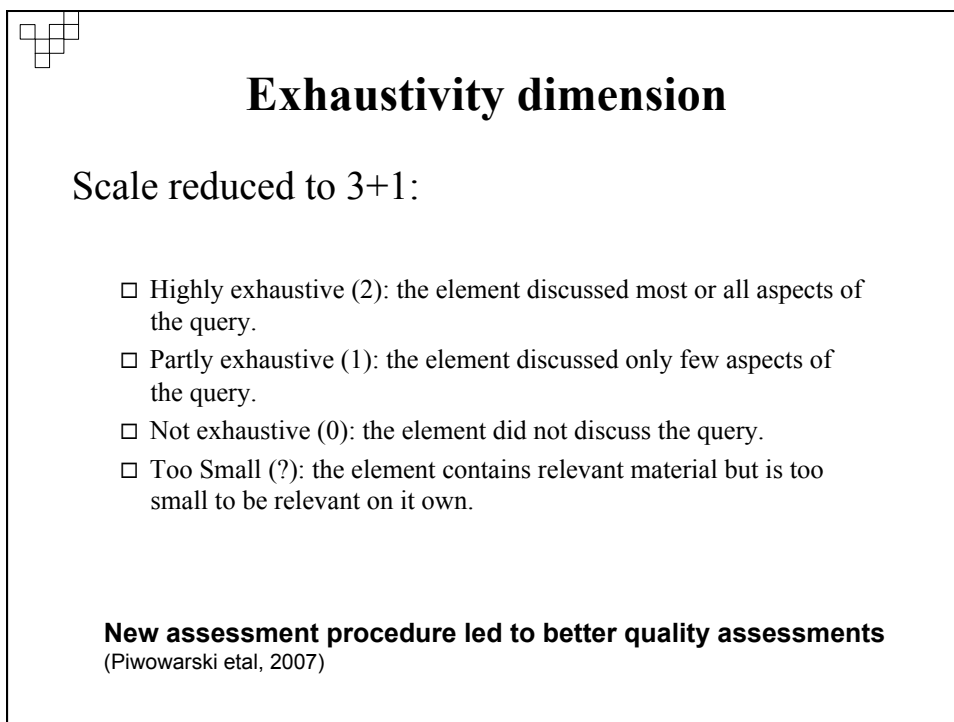
The other type of integration uses a DRAM macro embedded on an application specific IC (ASIC). For this purpose, designers have developed reconfigurable DRAM macros for many applications, providing a different configuration for each application. Although the many applications require a wide variety of configurations, the macro-testing methodology must be unified to reduce product-testing costs. This article describes circuitry that helps simplify testing the embedded-DRAM macro on an ASIC.>

>

<< **TESTING DILEMMA**

<¶¶ The dilemma in testing embedded DRAM arises from differences in character between ASICs and commodity DRAMs. In the case of commodity DRAMs, despite huge amounts of production, manufacturers produce only a few different products at the same time. ¶¶ As a result, they can optimize the testing methodology for each product. In contrast, companies produce a large variety of ASIC products, but the production volume of each product is small. Also, ASICs require a very short turnaround time. Therefore, customizing the test methodology for each product is difficult. ASICs require a common test environment that covers all product variations. ¶¶

<Furthermore, since the commodity DRAM is a general-purpose product, we cannot specify its application during testing. Thus, testing must cover various kinds of applications and provide very ¶¶ a result, the commodity DRAM's test time is longer than the ASIC's.



**Exhaustivity dimension**

Scale reduced to 3+1:

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

**New assessment procedure led to better quality assessments**  
(Piwowarski et al, 2007)



## Further simplification

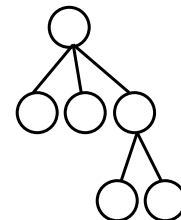
- **Statistical analysis on the INEX 2005 data:**
  - The exhaustivity 3+1 scale is not needed in most scenarios to compare XML retrieval approaches
  - *The two small maybe simulated by some threshold length*
- **INEX 2006-7 use only the specificity dimension to “measure” relevance**
  - The same highlighting approach is used
  - *Some investigation being done regarding the two small elements*

(Ogilvie & Lalmas, 2006)



## Relevance assessment task

- Topics are assessed by the INEX participants
- Pooling technique (~500 elements on runs of 1500 elements)
- **Completeness**
  - **Rules that force assessors to assess related elements**
  - E.g. element assessed relevant → its parent element and children elements must also be assessed
  - ...
- **Consistency**
  - **Rules to enforce consistent assessments**
  - E.g. Parent of a relevant element must also be relevant, although to a different extent
  - E.g. Exhaustivity increases going up; specificity increases going down
  - ...



(Piwowarski & Lalmas, CIKM 2004)



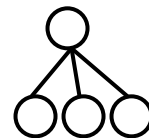
## Quality of assessments - Recap

- Very laborious assessment task, eventually impacting on the quality of assessments (Trotman, Glasgow IR festival 2005)
  - binary document agreement is 27% (compared to TREC 6 (33%) and TREC 4 (42.049%))
  - exact element agreement is 16%
- Interactive study shows that assessors agreement levels are high only at extreme ends of the relevance scale (very vs. not relevant) (Pehcevski *et al*, Glasgow IR festival 2005)
- Statistical analysis in 2004 data showed that comparisons of approaches would lead to same outcomes using a reduced scale (Ogilvie & Lalmas, 2006)
- A simplified assessment procedure based on highlighting (Clarke, Glasgow IR festival 2005)

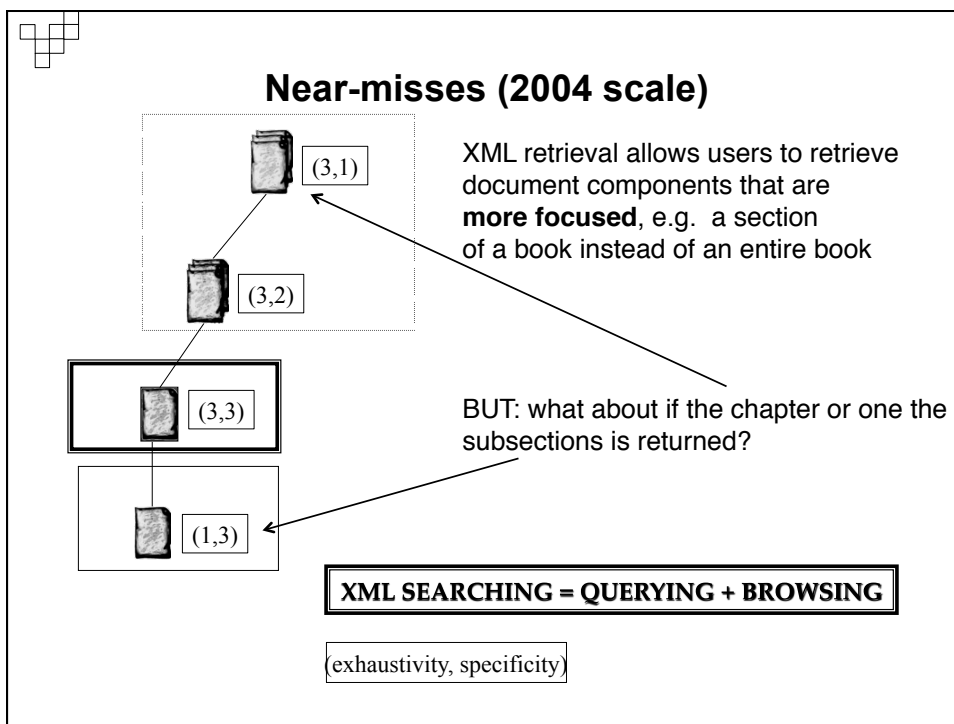
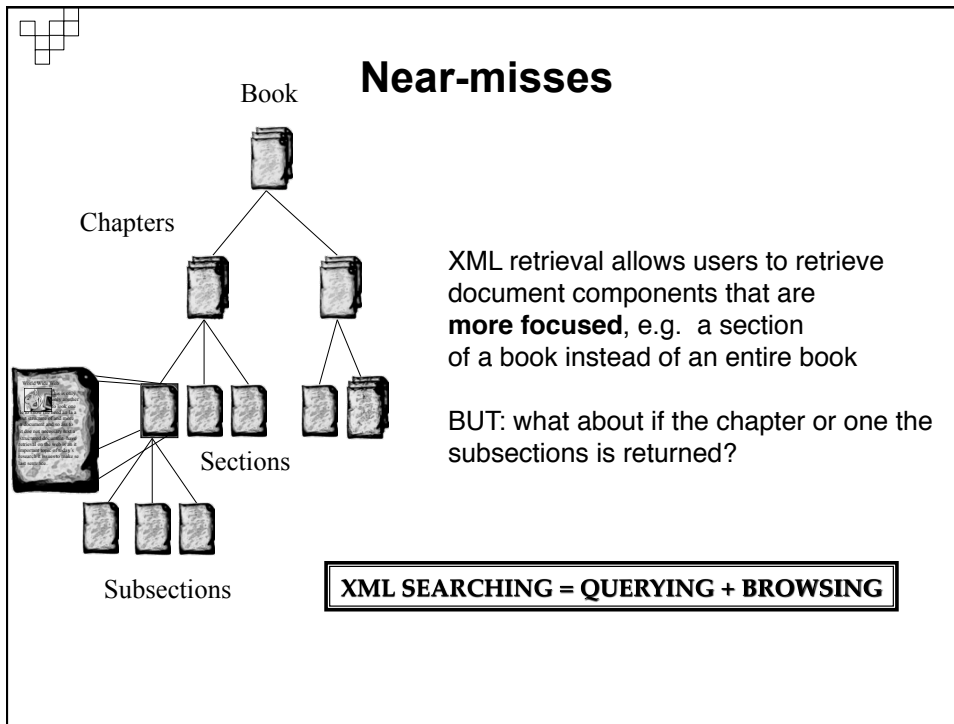


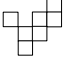
## Measuring effectiveness: Metrics

- Need to consider:
  - Multi-graded dimensions of relevance
  - Near-misses
- Metrics
  - inex\_eval (also known as inex2002) (Goevert & Kazai, INEX 2002)  
official INEX metric 2002-2004
  - inex\_eval\_ng (also known as inex2003) (Goevert *et al*, JIR 2006)
  - ERR (expected ratio of relevant units) (Piwowarski & Gallinari, INEX 2003)
  - xCG (XML cumulative gain) (Kazai & Lalmas, TOIS 2006)  
official INEX metric 2005-2006
  - t2i (tolerance to irrelevance) (de Vries *et al*, RIAO 2004)
  - EPRUM (Expected Precision Recall with User Modelling) (Piwowarski & Dupret, SIGIR 2006)
  - HiXEval (Highlighting XML Retrieval Evaluation) (Pehcevski & Thom, INEX 2005)
    - Variant of it is now official INEX metric 2007- (Kamps *et al*, INEX 2007)
  - ...



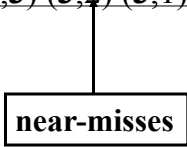


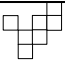




**Retrieve the best XML elements according to content and structure criteria (2004 scale):**

- Most exhaustive and the most specific = (3,3)
- Near misses = (3,3) + (2,3) (1,3) ← specific
- Near misses = (3, 3) + (3,2) (3,1) ← exhaustive
- Near misses = (3, 3) + (2,3) (1,3) (3,2) (3,1) (1,2) ...





**Two multi-graded dimensions of relevance**

- **Several “user models”**
  - **Expert and impatient:** only reward retrieval of highly exhaustive and specific elements (3,3) → *no near-misses*
  - Expert and patient: only reward retrieval of highly specific elements (3,3), (2,3) (1,3) → (2,3) and (1,3) are near-misses
  - ...
  - **Naïve and has lots of time:** reward - to a different extent - the retrieval of any relevant elements; i.e. everything apart (0,0) → *everything apart (3,3) is a near-miss*
- Use a quantisation function for each “user model”



## Examples of quantization functions

### Expert and impatient

$$\text{quant}_{\text{expert}}(e,s) = \begin{cases} 1 & \text{if } (e,s) = (3,3) \\ 0 & \text{otherwise} \end{cases}$$

### Naïve and has a lot of time

$$\text{quant}_{\text{naive}}(e,s) = \begin{cases} 1.00 & \text{if } (e,s) = (3,3) \\ 0.75 & \text{if } (e,s) \in \{(2,3), (3,2), (3,1)\} \\ 0.50 & \text{if } (e,s) \in \{(1,3), (2,2), (2,1)\} \\ 0.25 & \text{if } (e,s) \in \{(1,1), (1,2)\} \\ 0.00 & \text{if } (e,s) = (0,0) \end{cases}$$



## Using “standard” precision/recall

### Simulated runs (Piwowarski & Gallinari, INEX 2003)

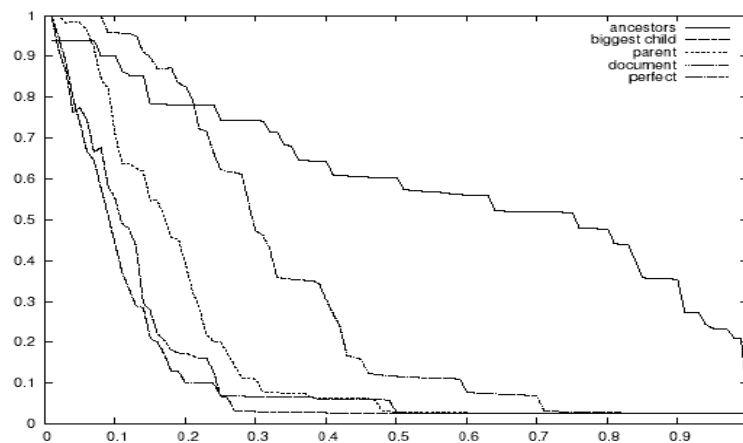
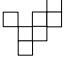


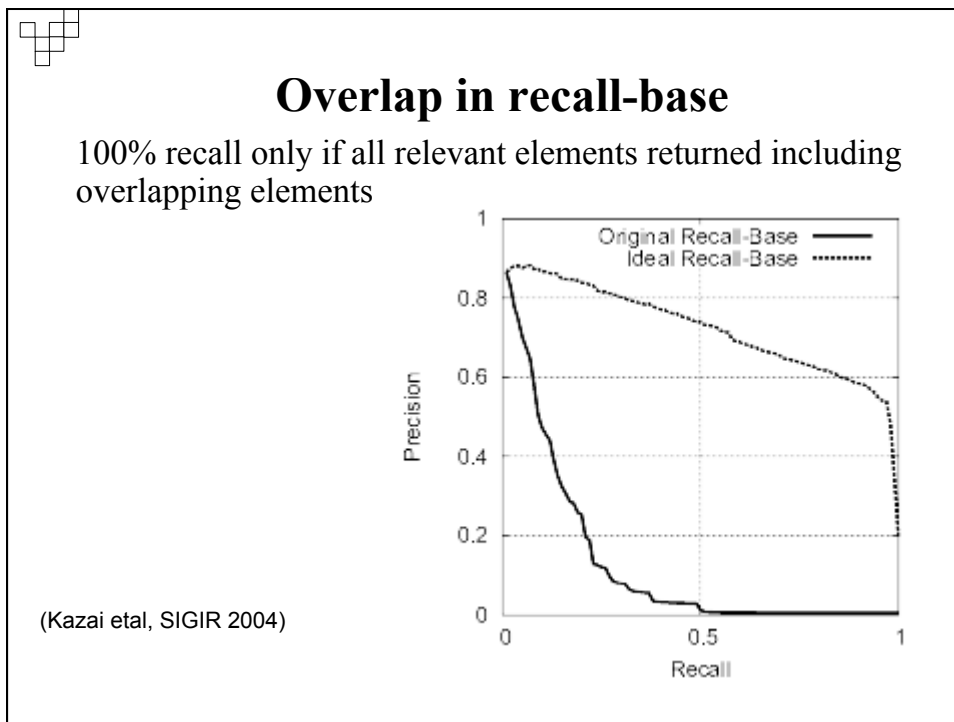
Figure 2. Generalised precision-recall. The axis of abscissas represents recall and the axis of ordinate the precision. Precision are averaged over the queries.




## Overlap in results

### Official INEX 2004 Results for CO topics

Rank	Systems (runs)	Avg Prec	<b>Overlap</b>
1.	IBM Haifa Research Lab (CO-0.5-LAREFIENMENT)	0.1437	<b>80.89</b>
2.	IBM Haifa Research Lab (CO-0.5)	0.1340	<b>81.46</b>
3.	University of Waterloo (Waterloo-Baseline)	0.1267	<b>76.32</b>
4.	University of Amsterdam (UAmS-CO-T-FBack)	0.1174	<b>81.85</b>
5.	University of Waterloo (Waterloo-Expanded)	0.1173	<b>75.62</b>
6.	Queensland University of Technology (CO_PS_Stop50K)	0.1073	<b>75.89</b>
7.	Queensland University of Technology (CO_PS_099_049)	0.1072	<b>76.81</b>
8.	IBM Haifa Research Lab (CO-0.5-Clustering)	0.1043	<b>81.10</b>
9.	University of Amsterdam (UAmS-CO-T)	0.1030	<b>71.96</b>
10.	LIP6 (simple)	0.0921	<b>64.29</b>



Relevance propagates up!



- ~26,000 relevant elements on  
~14,000 relevant paths
- **Propagated assessments: ~45%**
- **Increase in size of recall-base: ~182%**  
(INEX 2004 data)


(Kazai et al, SIGIR 2004)

**XCG: XML cumulated gain measures**

- Based on cumulated gain measure for IR (Kekäläinen and Järvelin, TOIS 2002)
- Accumulate gain obtained by retrieving elements up to a given rank; thus not based on precision and recall → **user-oriented measures**
- Extended to include a precision/recall behaviour → **system-oriented measures**
- Require the construction of
  - an ideal recall-base to separate what should be retrieved and what are near-misses
  - an associated ideal run, which contains what should be retrieved
- with which retrieval runs are compared, which include what is being retrieved, **including near-misses**.

(Kazai & Lalmas, TOIS 2006)

**HiXEval - Generalized precision and recall based on amount of highlighted content**



For each element, we derive:

- rsize: number of highlighted characters
- size: number of characters

For each topic, we derive:

- Trel: number of highlighted characters in collection

**HiXEval - Generalized precision and recall based on amount of highlighted content**

- Precision at rank  $r$ 

$$P(r) = \frac{\sum_{i=1}^r \text{rsize}(e_i)}{\sum_{i=1}^r \text{size}(e_i)}$$
- Recall at rank  $r$ 

$$R(r) = \frac{1}{T_{\text{rel}}} \cdot \sum_{i=1}^r \text{rsize}(e_i)$$
- F-measure at rank  $r$ , average precision, MAP, etc

(Pehcevski & Thom, INEX 2005; Kamps *et al*, INEX 2007)



## Evaluation and INEX - Recap

- Larger and more realistic collection with Wikipedia
- Better understanding of information needs and retrieval scenarios
- Better understanding of how to measure effectiveness
  - Near-misses and overlaps
  - Application to other IR problems
- Who are the real users?
  - Larsen et al, SIGIR 2006 poster; Betsi et al, SIGIR 2006 poster; Pharo & Trotman, SIGIR Forum 2007.
- Book search track at INEX



## Other INEX tracks

- Interactive (iTrack)
- Relevance feedback
- Natural language query processing
- Heterogeneous collection
- Multimedia track
- Document mining together with PASCAL network
- Use-case studies
- Entity ranking
- Link-the-Wiki
- Web service discovery
- Question Answering (QA@INEX)
- Efficiency
- Data-Centric



## Conclusions

- Major advances in XML search (ranking) approaches made possible with INEX
- Evaluating XML retrieval effectiveness itself a research problem
- Many open problems for research



## Areas for Open Problems

- DB and IR
  - Interaction between traditional DB query optimization (query rewriting) and ranking
- “Old” vs. new IR models
  - Combination of evidence problem
  - What evidence to use?
- Simple/succinct vs. complex/verbose QL
  - Define an XQuery core?





## Areas for Open Problems

- **Indexing & searching**
  - Efficient algorithms
  
- **INEX test collection and effectiveness**
  - Too complex?
  - What constitutes a retrieval baseline?
  - Generalisation of the results on other data sets
  
- **Quality evaluation (Web, XML)**
  - Who are the users?
  - What are their information needs?
  - What are the requirements?



## Outline

- Introduction to XML, basics and standards
  
- Document-oriented XML retrieval
  
- Evaluating XML retrieval effectiveness
  
- **Going beyond XML retrieval**



## **Beyond XML retrieval**

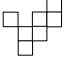
- Focused retrieval
- Aggregated results
- Structural context summarization
- Beyond the logical structure



## **Focused retrieval**

- Best performance obtained using evidence from element, document, and element size, and this whatever the model.
  - How can we apply this to other so-called “focused” retrieval problem?
  - What other evidence, e.g. semantic tags, should be used?
  - What combination formalism should be used?

See workshop on focused retrieval at SIGIR 2007, SIGIR 2008



## Aggregated results

- We know how to retrieve “snippets”.
- We know how to return “snippets” within a document (e.g. heatmap).
  
- How to combine/mix snippets from across documents to return **meaningful** aggregated results?
  - Virtual” documents (from Chiramella)
  - This is not just a data fusion problem!
  - “User experience” in search engine



## An example of aggregated search (1)




[Advanced Search](#)  
[Preferences](#)

**Web** [Show options...](#) Rt

**Cristiano Ronaldo - Wikipedia, the free encyclopedia** [f](#) [x](#) - 10:47am  
 Cristiano Ronaldo dos Santos Aveiro, OIH (Portuguese pronunciation: [kɾiʃˈtiɲu wɐˈnɐɾdu]; born 5 February 1985) is a Portuguese footballer who plays as a ...  
[en.wikipedia.org/wiki/Cristiano\\_Ronaldo](http://en.wikipedia.org/wiki/Cristiano_Ronaldo) - 254k - [Cached](#) - [Similar pages](#) - [f](#) [x](#)

**Cristiano Ronaldo - Manchester United Official Web Site** [f](#) [x](#)  
 Cristiano Ronaldo has exhausted all superlatives, except to say that having developed and matured from an inexperienced, young winger when he left Sporting ...  
[www.manutd.com/default.sps?bioid=91960&pageid=%7B7FE60904B-C2A8-4E60-9B05-700DBBC29BBC%7D...](http://www.manutd.com/default.sps?bioid=91960&pageid=%7B7FE60904B-C2A8-4E60-9B05-700DBBC29BBC%7D...) - 38k - [Cached](#) - [Similar pages](#) - [f](#) [x](#)

**Cristiano Ronaldo 7.com | The Ultimate Cristiano Ronaldo Fan Site** [f](#) [x](#)  
 Photographs, match reports, biography, fixtures, and a shop.  
[www.cristianoronaldo7.com/](http://www.cristianoronaldo7.com/) - 27k - [Cached](#) - [Similar pages](#) - [f](#) [x](#)

**Image results for Cristiano Ronaldo** - [Report images](#)




**Video results for Cristiano Ronaldo**

 <a href="#">cristiano ronaldo</a> 7 min <a href="http://www.youtube.com">www.youtube.com</a>	 <a href="#">Cristiano Ronaldo - The Perfect Player 2008</a> 5 min <a href="http://www.youtube.com">www.youtube.com</a>
---	---

**Cristiano Ronaldo Fan | News, Photos, Blog, Pics, Videos, Wallpapers** [f](#) [x](#)  
 Cristiano Ronaldo - Fan site dedicated to Cristiano Ronaldo. All the latest news , videos,

## An example of aggregated search (2)




Hiya Guest | Sign in  
Customise my search

Also try: [heligate london](#), [london time](#), [london hotels](#), [lauren london](#), [london eye](#) more...

Web Search by Yahoo! Search 782,000,000 results

- Visit **London** - London's official website - **London Tourist** ... Visit **London** with ease using **London Tourist**, Entertainment and Travel Guides. ... those tiaras and feather boas – the world's campest musical is back in **London!** ... <http://www.visitlondon.com/>
- LondonTown.com | The Number One Internet Site for **London** England ... Sightseeing, attraction and other tourist information. Theater tickets, **London** maps and more. ... Christmas Lights in **London** 2007 ... <http://www.londontown.com/>
- London** Entertainment Guide from The Evening Standard | This is **London** ... tickets booking for **London** covering Film, Restaurants, ... **London** 'safest place' ... This is **London** Magazine Ltd. Daily Mail Mail on Sunday Travel Mail ... <http://www.thisislondon.co.uk/>
- London** - Wikipedia, the free encyclopedia  
An important settlement for two millennia, **London** is one of the world's leading ... **London** is a major tourist destination with four world heritage sites, several ... <http://en.wikipedia.org/wiki/London>
- Mayor of **London**, the **London** Assembly and the Greater **London** Authority  
The official website for the Mayor of **London**, the **London** Assembly and the Greater **London** Authority ... **London** Assembly members, Richard Barnes and Tony Arbour, ... <http://www.london.gov.uk/>
- Time Out **London** - Discover Things To Do in **London** Including ... Highlights this month's events in **London** and provides details of accommodation, ... **London** Jazz Festival ... **London** ice-skating rinks. Start practising your ... <http://www.timeout.com/london/>
- BBC - **London** - **London** Homepage  
**London** homepage ... **London** seeks throw-away bags ban. **London**'s boroughs approve a bill seeking a city-wide ban of free disposable shopping bags. ...




Photos by Flickr



More results | Prev | Next

---

Yahoo! Answers

-  London!!!!!!!?  
...northwest of England and i love going down to London a few times a yr, & each time i go...  
<http://au.answers.yahoo.com/question/?qid=2007C>
-  London???????  
I am going London tommorow and im a kid what places are main...  
<http://au.answers.yahoo.com/question/?qid=2007C>
-  london.....?  
...of beautiful people down the place i live in london getting on and off the bus...do any...  
<http://au.answers.yahoo.com/question/?qid=2006f>

More results | Prev | Next

## An example of aggregated search (3)



로그인 | 내아바 | 알림 | 카페 | 블로그 | 더보기

통합검색 사이트 | 홈연서 | 지식서 | 블로그 | 카페 | 이미지 | 동영상 | 사진 | 뉴스 | 더보기

### '대한민국' 국가정보



대한민국

위치: 아시아 동북 동북 끝  
수도: 서울 (Seoul)  
언어: 한국어  
기후: 대륙성 / 온순기후  
면적: 9만 9538㎢(남한)  
역사: 원시, 고대, 중세, 근대, 현대사...more  
인구: 4842만 2000명(2005)  
주요도시: 서울

대한정보

관광명소: 서울 - 한강, 궁궐, 한계천, 공원, 시장...  
부산 - 해운대, 동백섬, 오륙도...  
대전 - 엑스포과학공원, 유성온천...  
강원도 - 설악산국립공원, 정선교...  
제주도 - 서귀포시, 제주시, 남제주군...  
대한민국 여행정보 더보기

대외정보

한국의 새: 선명 사는 새, 울며 사는 새...more  
한국의 말: 삼국시대 말, 통일신라 말...more  
한국의 사람: 왕양자, 박경자, 박석자...more

### 연관검색어

korean  
japan  
china  
한국  
대한민국  
corea

도움말

---

### 국가 검색순위

오스트레일리아	↑ 2
대한민국	↓ 1
일본	↓ 1
중국	- 0
미국	- 0
영국	↑ 1
중국	↓ 1
프랑스	↑ 5
캐나다	↑ 6
타이	- 0

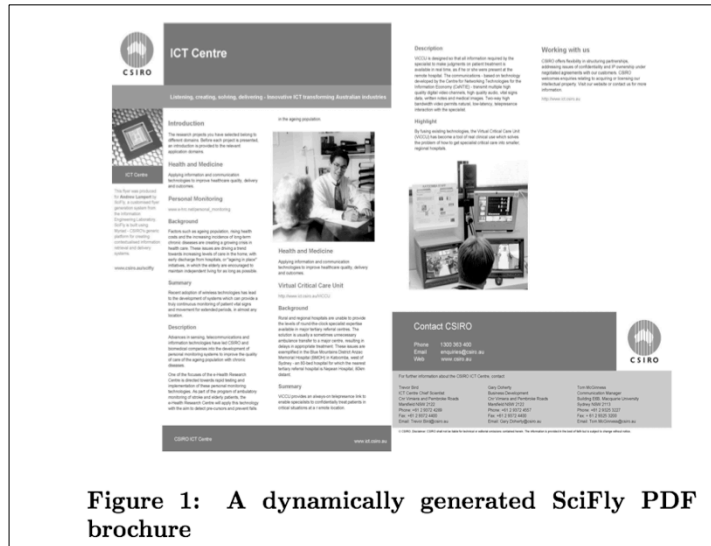
도움말 2008.09.21

사전

백과사전 **대한민국 [Korea]**  
한국(韓國) 또는 남한(South Korea)이라고도 부르고, 북한에서는 남조선(南朝鮮)이라고 부른다. 남북으로 길게 뻗은 반도와 3200여...

영어사전 **Korea**  
\*\*\* Korea [ka:ri:ka: - |ka:ri:] [Kor. 「고려」가 변한 말] n. 한국 (공식 명칭 the Republic of Korea) 동칭 South Korea:수도 Seoul(漢 PDK) > Korjān a.

## An example of aggregated search (4)




**Figure 1: A dynamically generated SciFly PDF brochure**

(Courtesy of Cecile Paris, CSIRO, Sydney, Australia)


## Structural context summarization

- Users require document context when viewing an elements result
- We know how to summarize the structure (ToC) of a document (depth, relevance, etc)
- How can we summarize the structure of the search results, to provide context for the whole search.
  - Not just clusters
  - Refer to Amazon property list for a given product




---

dbdk\_training in **Baseline System**



query was: text classification naive bayes  
Results **1 - 10** of **100**.  
Result pages: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [next](#)



---

### Search Result

- 1: (0.247) **Scalable Feature Mining for Sequential Data**  
*Neal Lesh Mitsubishi Electric Research Lab Mohammed J. Zaki Rensselaer Polytechnic Institute Mitsunori Ogihara University of Rochester*  
Result path: /article[1]/bdy[4]/sec[5]
- 2: (0.204) **Probability and Agents**  
*Marco G. Valtorta University of South Carolina mgv@csce.sc.edu Michael N. Huhns University of South Carolina huhns@sc.edu*  
Result path: /article[1]/bdy[4]/sec[3]
- 3: (0.176) **Combining Image Compression and Classification Using Vector Quantization**  
*Karen L. Oehler Member IEEE Robert M. Gray Fellow IEEE*  
Result path: /article[1]/bdy[4]/sec[4]/ss1[2]/ss2[4]
- 4: (0.175) **Text-Learning and Related Intelligent Agents: A Survey**  
*Dunja Mladenic J. Stefan Institute*  
Result path: /article[1]/hm[5]/app[4]/sec[5]
- 5: (0.175) **Detecting Faces in Images: A Survey**  
*Ming-Hsuan Yang Member IEEE David J. Kriegman Senior Member IEEE Narendra Ahuja Fellow IEEE*  
Result path: /article[1]/bdy[4]/sec[2]/ss1[9]/ss2[10]

## Providing context for the element

**Table of contents:**

- [Albert Einstein](#)
  - [Einstein Albert Einstein...](#)
    - [Biography](#)
      - [Youth and college](#)
      - [Work and doctorate](#)
      - [Middle years](#)
        - [The Einstein refrigerator](#)
        - [World War II](#)
      - [Final years](#)
    - [Personality](#)
      - [Albert Einstein was much...](#)
      - [Political views](#)
    - [Popularity and cultural impact](#)
      - [Einstein's popularity has...](#)
      - [Entertainment](#)
        - [Albert Einstein has becom...](#)
      - [Licensing](#)
      - [Honors](#)
    - [References](#)
    - [Works by Albert Einstein](#)
    - [External links](#)

### Entertainment

Albert Einstein has become the subject of a number of novels, films and plays, including Nicolas Roeg's film *Insignificance*, Fred Schepin's film *L.Q.*, Alan Lightman's novel *Einstein's Dreams*, and Steve Martin's comedic play "Picasso at the Lapin Agile". He was the subject of Philip Glass's groundbreaking 1976 opera *Einstein on the Beach*. Since 1978, Einstein's humorous side has been the subject of a live stage presentation *Albert Einstein: The Practical Bohemian*, a one man show performed by actor Ed Metzger.

He is often used as a model for depictions of eccentric scientists in works of fiction, his own character and distinctive hairstyle suggest eccentricity, electricity, or even insanity and are widely copied or exaggerated. TIME magazine writer Frederic Golden referred to Einstein as "a cartoonist's dream come true."

On Einstein's 72nd birthday in 1951, the UPI photographer Arthur Sasse was trying to coax him into smiling for the camera. Having done this for the photographer many times that day, Einstein stuck out his tongue instead. The image has become an icon in pop culture for its contrast of the genius scientist displaying a moment of levity. Yahoo! Serious, an Australian film maker, used the photo as an inspiration for the intentionally anachronistic movie *Young Einstein*.



## Beyond the logical structure

- We know how to exploit the tags representing the logical structure to provide focused retrieval.
- What about other tags, e.g. semantic tags, formatting tags, template tags, etc?



## Acknowledgements

- These slides are based on a number of presentations from the presenters at other events and from other researchers.
  - **Accessing XML content - An information retrieval perspective, Winter School, Champéry, 2008**
  - **M. Consens, R. A. Baeza-Yates, M. Lalmas, S. Amer-Yahia: XML retrieval: DB/IR in theory, web in practice. VLDB 2007**
  - **S. Amer-Yahia, R. Baeza-Yates, M. Consens and M. Lalmas. XML Retrieval: Integrated IR-DB Challenges and Solutions. SIGIR 2007.**
  - **S. Amer-Yahia and M. Lalmas. Accessing XML Content: From DB and IR Perspectives, CIKM 2005.**
  - R. Baeza-Yates and N. Fuhr. XML Retrieval, SIGIR 2004
  - R. Baeza-Yates and M. Consens. The Continued Saga of DB-IR Integration, SIGIR 2005
  - M. Lalmas. Structure/XML retrieval. ESSIR 2005, ESSIR 2007
  - M. de Rijke, J. Kamps and M. Marx. Retrieving Content and Structure, ESSLI 2005
  - B. Sigurbjörnsson, Element Retrieval in Action, QMUL Seminar 2005
  - R. Baeza-Yates and M. Lalmas, XML Information Retrieval, SIGIR 2006



- S. Amer-Yahia & M. Lalmas. XML Search: Languages, INEX and Scoring, Submitted for Publication, 2006.
- E. Amitay, D. Carmel, R. Lempel & A. Soffer. Scaling IR-system evaluation using term relevance sets. SIGIR 2004, pp 10-17.
- P. Arvola, J. Kekäläinen & M. Junkkari. Query Evaluation with Structural Indices. INEX 2005.
- P. Arvola, M. Junkkari & J. Kekäläinen. Generalized contextualization method for XML information retrieval. CIKM 2005.
- R. A. Baeza-Yates, N. Fuhr & Y.S. Maarek. SIGIR XML and Information Retrieval workshop, SIGIR Forum, 36(2):53–57, 2002. 3.
- R. A. Baeza-Yates, Y. S. Maarek, T. Roelleke & A.P. de Vries. SIGIR joint XML & Information Retrieval and Integration of IR and DB workshops, SIGIR Forum, 38(2):24–30, 2004.
- R. Baeza-Yates, D. Carmel, Y.S. Maarek, and A. Sofer (eds). Special issue on XML Retrieval, JASIST, 53, 2002.
- R. Baeza-Yates & G. Navarro, Integrating contents and structure in text retrieval, SIGMOD 25:67-79, 1996.
- R. Baeza-Yates and G. Navarro, XQL and Proximal Nodes, JASIST 53:504-514, 2002.
- S. Betsi, M. Lalmas, A. Tombros & T. Tsirikika. User Expectations from XML Element Retrieval, SIGIR 2006 (Poster).



- Henk M. Blanken, T. Grabs, H.-J. Schek, R. Schenkel & G. Weikum (eds). Intelligent Search on XML Data, Applications, Languages, Models, Implementations, and Benchmarks, 2003.
- P. Borlund. The concept of relevance in IR. JASIS, 54(10):913-925, 2003.
- J.P. Callan. (1994). Passage-level evidence in document retrieval. SIGIR 1994.
- D. Carmel, Y.S. Maarek & A. Soffer. XML and Information Retrieval. SIGIR Forum, 34(1): 31–36, 2000.
- D. Carmel, Y.S. Maarek, M. Mandelbrod, Y. Mass & A. Soffer: Searching XML documents via XML fragments. SIGIR 2003.
- Y. Chiramella, P. Mulhem & F. Fourel. A model for multimedia information retrieval. FERMI Technical report, University of Glasgow, 1996.
- Chinenyanga and Kushmerik, Expressive retrieval from XML documents, SIGIR 2001. 2001.
- C. Clarke. Range results in XML retrieval. INEX 2005 Workshop on Element Retrieval Methodology.
- C. Clarke. Controlling Overlap in Content-Oriented XML Retrieval. SIGIR 2005.
- W.S. Cooper. Expected search length: A single measure of retrieval effectiveness based on weak ordering action of retrieval systems. JASIS, 19:30-41, 1968.
- A. Delgado & R. Baeza-Yates. A Comparison of XML Query Languages, Upgrade 3, 12-25, 2002.
- L. Denoyer & P. Gallinari. The Wikipedia XML Corpus. SIGIR Forum, 40(1), 2006.





- A. de Vries, G. Kazai & M. Lalmas. Tolerance to Irrelevance: A User-effort Oriented Evaluation of Retrieval Systems without Predefined Retrieval Unit, RIAO 2004.
- N. Fuhr, N. Goevert, G. Kazai & M. Lalmas (eds). INitiative for the Evaluation of XML Retrieval (INEX 2002): Proceedings of the First INEX Workshop. ERCIM Workshop Proceedings, 2003.
- N. Fuhr & Kai Großjohann. XIRQL: A Query Language for Information Retrieval in XML Documents. SIGIR 2001.
- N. Fuhr, M. Lalmas & S. Malik (eds). INitiative for the Evaluation of XML Retrieval (INEX 2003). Proceedings of the Second INEX Workshop, 2004.
- N. Fuhr, M. Lalmas, S. Malik & G. Kazai (eds). Advances in XML Information Retrieval and Evaluation: Fourth Workshop of the INitiative for the Evaluation of XML Retrieval (INEX 2005), LNCS 3977, 2006.
- N. Fuhr, M. Lalmas, S. Malik & Z. Szavik (eds). Advances in XML Information Retrieval, Third International Workshop of the Initiative for the Evaluation of XML Retrieval (INEX 2004), INEX 2004, LNCS 3493, 2005.
- S. Geva. GPX - Gardens Point XML IR at INEX 2005. INEX 2005.
- N. Goevert, N. Fuhr, M. Lalmas, & G. Kazai. Evaluating the effectiveness of content-oriented XML retrieval methods. Journal of Information Retrieval, 2006 (In Press).
- N. Goevert & G. Kazai. Overview of the INitiative for the Evaluation of XML retrieval (INEX) 2002, INEX 2002.



- M. A. Hearst & C. Plaunt. Subtopic structuring for full-length document access. SIGIR 1993.
- K. Järvelin & J. Kekäläinen. Cumulated gain-based evaluation of IR techniques. ACM TOIS 20(4):422-446, 2002.
- J. Kamps, M. de Rijke & B. Sigurbjornsson. Length normalization in XML retrieval. SIGIR 2004.
- J. Kamps, M. de Rijke & B. Sigurbjornsson. The importance of length normalization for XML retrieval. Information Retrieval, 8(4):631-654, 2005.
- G. Kazai & M. Lalmas. eXtended Cumulated Gain Measures for the Evaluation of Content-oriented XML Retrieval. ACM TOIS, 2006 (To appear).
- G. Kazai, M. Lalmas & A. de Vries. The overlap problem in content-oriented xml retrieval evaluation. SIGIR 2004.
- G. Kazai, M. Lalmas, N. Fuhr & N. Gövert. A report on the first year of the INitiative for the evaluation of XML retrieval (INEX 02). JASIST, 54, 2004.
- G. Kazai, M. Lalmas & J. Reid. Construction of a test collection for the focussed retrieval of structured documents. ECIR 2003.
- M. Lalmas & E. Moutogianni. A Dempster-Shafer indexing for the focussed retrieval of a hierarchically structured document space: Implementation and experiments on a web museum collection, RIAO 2000.
- M. Lalmas and I. Ruthven. Representing and Retrieving Structured Documents using the Dempster-Shafer Theory of Evidence: Modelling and Evaluation. JDoc, 54(5):529-565, 1998.



- B. Larsen, A. Tombros & S. Malik . Is XML retrieval meaningful to users? Searcher preferences for full documents vs. elements. SIGIR 2006 (Poster).
- Luk, Leong, Dillon, Chan, Croft & Allan, A Survey on Indexing and Searching XML, "Special Issue on XML and IR", *JASIST*, 2002.
- Mass, Mandelbrod, Amitay, and Soffer, JuruXML - an XML retrieval system at INEX 2002. INEX 2003.
- Y. Mass & M. Mandelbrod. Retrieving the most relevant XML Components. INEX 2004.
- Y. Mass & M. Mandelbrod. Using the INEX environment as a test bed for various user models for XML Retrieval. INEX 2005.
- V. Mihajlovic, G. Ramirez, T. Westerveld, D. Hiemstra, H. E. Blok & A. P. de Vries. TIJAH Scratches INEX 2005: Vague Element Selection, Image Search, Overlap, and Relevance Feedback. INEX 2005.
- Navarro and Baeza-Yates, Proximal Nodes, SIGIR 1995 (journal version in ACM TOIS, 1997).
- P. Ogilvie & M. Lalmas. Investigating the exhaustivity dimension in content-oriented XML element retrieval evaluation. 2006. Submitted for publication.
- Paul Ogilvie & Jamie Callan: Hierarchical Language Models for XML Component Retrieval. INEX 2004.
- Paul Ogilvie & Jamie Callan: Parameter Estimation for a Simple Hierarchical Generative Model for XML Retrieval. INEX 2005.
- J. Pehcevski & J. A. Thom. Hixeval: Highlighting xml retrieval evaluation. INEX 2005.



- J. Pehcevski, J.A. Thom & A.-M. Vercoustre. Hybrid XML Retrieval: Combining Information Retrieval and a Native XML Database. *Journal of Information Retrieval* 8(4): 571-600, 2005.
- J. Pehcevski, J. A. Thom & A.M. Vercoustre. Users and assessors in the context of INEX: Are relevance dimensions relevant? INEX 2005 Workshop on Element Retrieval Methodology.
- B. Piwowarski & G. Dupret. Evaluation in (XML) Information Retrieval: Expected Precision-Recall with User Modelling (EPRUM), SIGIR 2006.
- B. Piwowarski & P. Gallinari. Expected ratio of relevant units: A measure for structured information retrieval. INEX 2003.
- B. Piwowarski & M. Lalmas. Providing consistent and exhaustive relevance assessments for XML retrieval evaluation. CIKM 2004.
- B. Piwowarski, A. Trotman & M. Lalmas. Sound and complete relevance assessments for XML retrieval. 2006. Submitted for publication.
- V.V. Raghavan, P. Bollmann, & G. S. Jung. A critical investigation of recall and precision as measures of retrieval system performance. *ACM TOIS*, 7(3):205-229, 1989.
- G. Ramirez, T. Westerveld & A. P. de Vries. Using structural relationships for focused XML retrieval. FQAS 2006.
- T. Roelleke, M. Lalmas G. Kazai, I Ruthven & S. Quicker. The Accessibility Dimension for Structured Document Retrieval, ECIR 2002.
- G. Salton, J. Allan & C. Buckley. Approaches to Passage Retrieval in Full Text Information Systems, SIGIR 1993.



- K. Sauvagnat, L. Hlaoua & M. Boughanem XFIRM at INEX 2005: ad-hoc and relevance feedback tracks. INEX 2005.
- B. Sigurbjornsson, J. Kamps & M. de Rijke The Importance of Length Normalization for XML Retrieval. *Journal of Information Retrieval*,8(4), 2005.
- B. Sigurbjornsson, J. Kamps & M. de Rijke The Effect of Structured Queries and Selective Indexing on XML Retrieval. INEX 2005.
- B. Sigurbjornsson & A. Trotman. Queries: INEX 2003 working group report. INEX 2003.
- A. Trotman and M. Lalmas. Strict and Vague Interpretation of XML-Retrieval Queries, SIGIR 2006 (Poster).
- M. Theobald, R. Schenkel & G. Weikum. TopX & XXL at INEX 2005. INEX 2005.
- A. Tombros, S. Malik & B. Larsen. Report on the INEX 2004 interactive track. *ACM SIGIR Forum*, 39(1):43–49, 2005.
- A. Trotman. Wanted: Element retrieval users. INEX 2005 Workshop on Element Retrieval Methodology.
- A. Trotman & M. Lalmas. Why Structural Hints in Queries do not Help XML Retrieval, SIGIR 2006 (Poster).
- A. Trotman & B. Sigurbjornsson. NEXI, now and next. INEX 2004.
- A. Trotman & B. Sigurbjornsson. Narrowed extended XPATH I (NEXI). INEX 2004.
- C. J. van Rijsbergen. *Information Retrieval*. Butterworths, 1979.



- K. Sauvagnat, L. Hlaoua & M. Boughanem XFIRM at INEX 2005: ad-hoc and relevance feedback tracks. INEX 2005.
- B. Sigurbjornsson, J. Kamps & M. de Rijke The Importance of Length Normalization for XML Retrieval. *Journal of Information Retrieval*,8(4), 2005.
- B. Sigurbjornsson, J. Kamps & M. de Rijke The Effect of Structured Queries and Selective Indexing on XML Retrieval. INEX 2005.
- B. Sigurbjornsson & A. Trotman. Queries: INEX 2003 working group report. INEX 2003.
- A. Trotman and M. Lalmas. Strict and Vague Interpretation of XML-Retrieval Queries, SIGIR 2006 (Poster).
- M. Theobald, R. Schenkel & G. Weikum. TopX & XXL at INEX 2005. INEX 2005.
- A. Tombros, S. Malik & B. Larsen. Report on the INEX 2004 interactive track. *ACM SIGIR Forum*, 39(1):43–49, 2005.
- A. Trotman. Wanted: Element retrieval users. INEX 2005 Workshop on Element Retrieval Methodology.
- A. Trotman & M. Lalmas. Why Structural Hints in Queries do not Help XML Retrieval, SIGIR 2006 (Poster).
- A. Trotman & B. Sigurbjornsson. NEXI, now and next. INEX 2004.
- A. Trotman & B. Sigurbjornsson. Narrowed extended XPATH I (NEXI). INEX 2004.
- C. J. van Rijsbergen. *Information Retrieval*. Butterworths, 1979.



- J.-N. Vittaut & P. Gallinari. Machine Learning Ranking for Structured Information Retrieval. ECIR 2006.
- R. Wilkinson. Effective Retrieval of Structured Documents. SIGIR 1994.
- A. Woodley & S. Geva. NLPX at INEX 2004. INEX 2004.
- J.-N. Vittaut, B. Piwowarski & P. Gallinari. An Algebra for Structured Queries in Bayesian Networks. INEX 2004.