

XML Search and XQuery Full-Text

Sihem Amer-Yahia

Yahoo! Research
Community Systems Group

Stanford guest lecture

Feb. 12th, 2007



Outline

- Motivation
- Challenges
- Languages
 - XQuery Full-Text
 - INEX
- Research overview



Motivation

- XML is able to represent a mix of structured and text information:
 - XML applications: *digital libraries, content management.*
 - XML repositories: *IEEE INEX collection, SIGMOD Record in XML, LexisNexis, the Library of Congress collection, HL7, MPEG7.*
- Need for a language to *search XML documents*

109TH CONGRESS
1ST SESSION

H. R. 2739

To address rising college tuition by strengthening the compact between the States, the Federal Government, and institutions of higher education to make college more affordable.

IN THE HOUSE OF REPRESENTATIVES

MAY 26, 2005

Mr. TIERNEY (for himself, Ms. MCCOLLUM of Minnesota, Mr. GEORGE MILLER of California, Mr. KILDEE, Mr. EMANUEL, Mr. BISHOP of New York, Mr. PAYNE, Ms. WOOLSEY, Mrs. MCCARTHY, Mr. WU, Mr. DAVIS of Illinois, Mr. GRIJALVA, Mr. MEEHAN, Mr. BECERRA, Mr. REYES, Mr. GONZALEZ, Ms. LINDA T. SÁNCHEZ of California, Mr. MCGOVERN, Ms. DELAURO, Mr. OWENS, Mr. HINOJOSA, Mr. KUCINICH, Mr. HOLT, Mr. CASE, Mr. VAN HOLLEN, Mr. ORTIZ, Mr. GUTIERREZ, Mr. CARDOZA, Mrs. JONES of Ohio, Ms. BALDWIN, Mr. WEXLER, Mr. BARROW, Mr. JEFFERSON, Mr. RYAN of Ohio, Ms. SOLIS, Ms. VELÁZQUEZ, and Ms. SCHAKOWSKY) introduced the following bill; which was referred to the Committee on Education and the Workforce

A BILL

To address rising college tuition by strengthening the compact between the States, the Federal Government, and institutions of higher education to make college more affordable.

Be it enacted by the Senate and House of Representatives of the United States of America in Congress assembled,



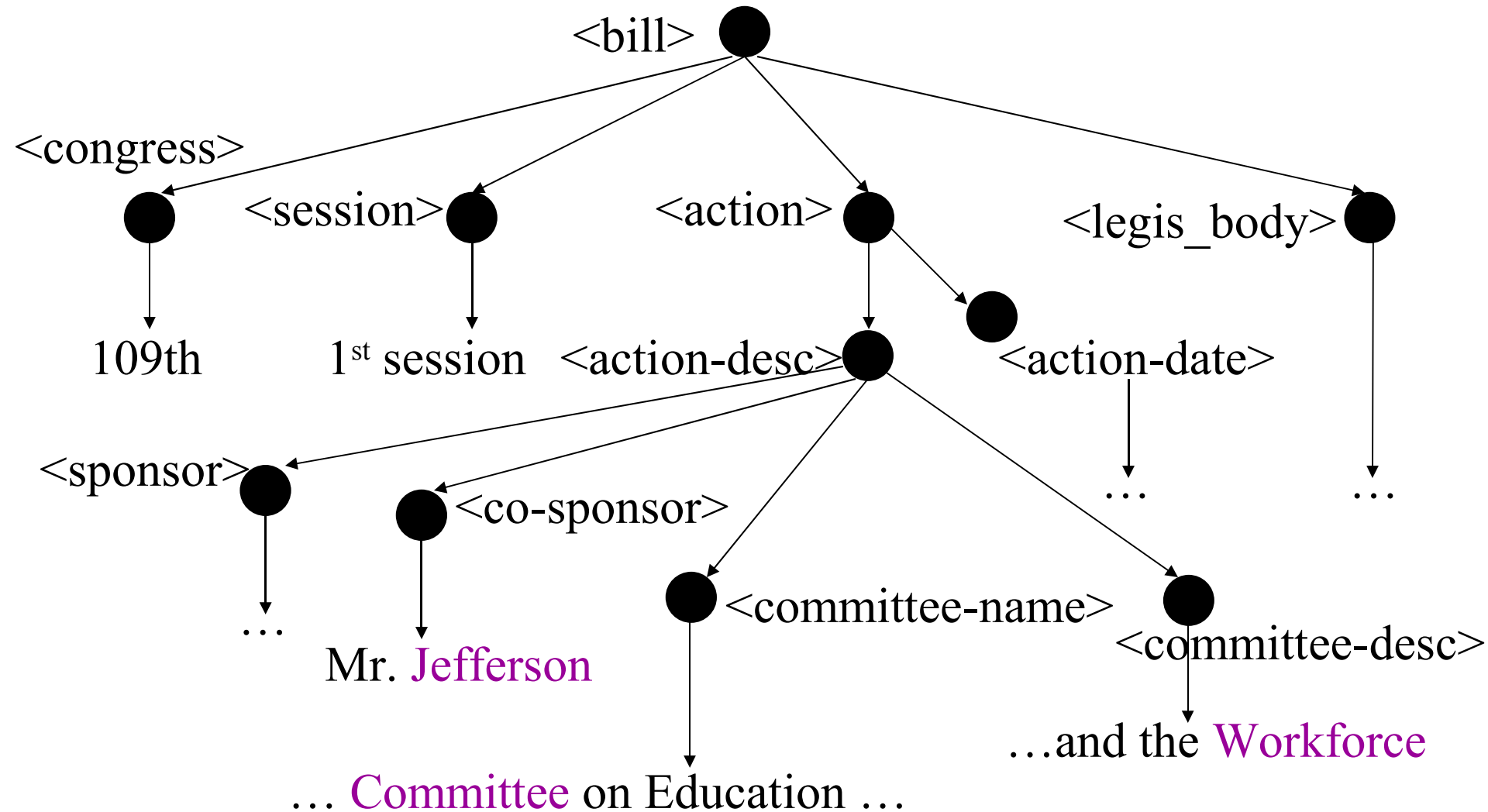
LoC XML Document

http://thomas.loc.gov/home/gpoxmlc109/h2739_ih.xml

```
<bill bill-stage = "Introduced-in-House">
  <congress> 109th CONGRESS </congress>
  <session> 1st Session </session>
  <legis-num> H. R. 2739 </legis-num>
  <current-chamber> IN THE HOUSE OF REPRESENTATIVES </current-chamber>
  <action>
    <action-date date = "20050526"> May 26, 2005 </action-date>
    <action-desc><sponsor name-id = "T000266"> Mr. Tierney </sponsor> (for
      himself, and <cosponsor name-id = "M001143"> Ms. McCollum of Minnesota
      </cosponsor>, <cosponsor name-id = "M000725"> Mr. George Miller of
      California </cosponsor>) introduced the following bill; which was referred to the
      <committee-name committee-id = "HED00"> Committee on Education and the
      Workforce </committee-name>
    </action-desc>
  </action>
  ...
</bill>
```



LoC Document Example





THOMAS Search Engine

Search Full Text of the Congressional Record - 109th Congress - Microsoft Internet Explorer

File Edit View Favorites Tools Help



[Thomas Home](#)
[Library of Congress](#)

Bill Text

109th Congress (2005-2006)

Select Congress: [109](#) | [108](#) | [107](#) | [106](#) | [105](#) | [104](#) | [103](#) | [102](#) | [101](#) [HELP](#)

Bill Number: [\[Help\]](#)

Examples: h.r. 1425, S. 896, h.j.res. 125, scores 24

[View](#) Complete List of Bills in this Congress by Type and Bill Number

The following fields can be used singly or in combination:

Word/Phrase: [\[Help\]](#)

- All Bills
- Bills with **floor action**
- Enrolled bills** sent to the President
- Both House and Senate Bills
- House Bills only
- Senate Bills only
- Exact word(s)
- Word variants (plurals, etc.)

Date/Session: [\[Help\]](#)

Format: mm/dd/yyyy or mm-dd-yyyy

From Through

Words in the Index:

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000

Internet

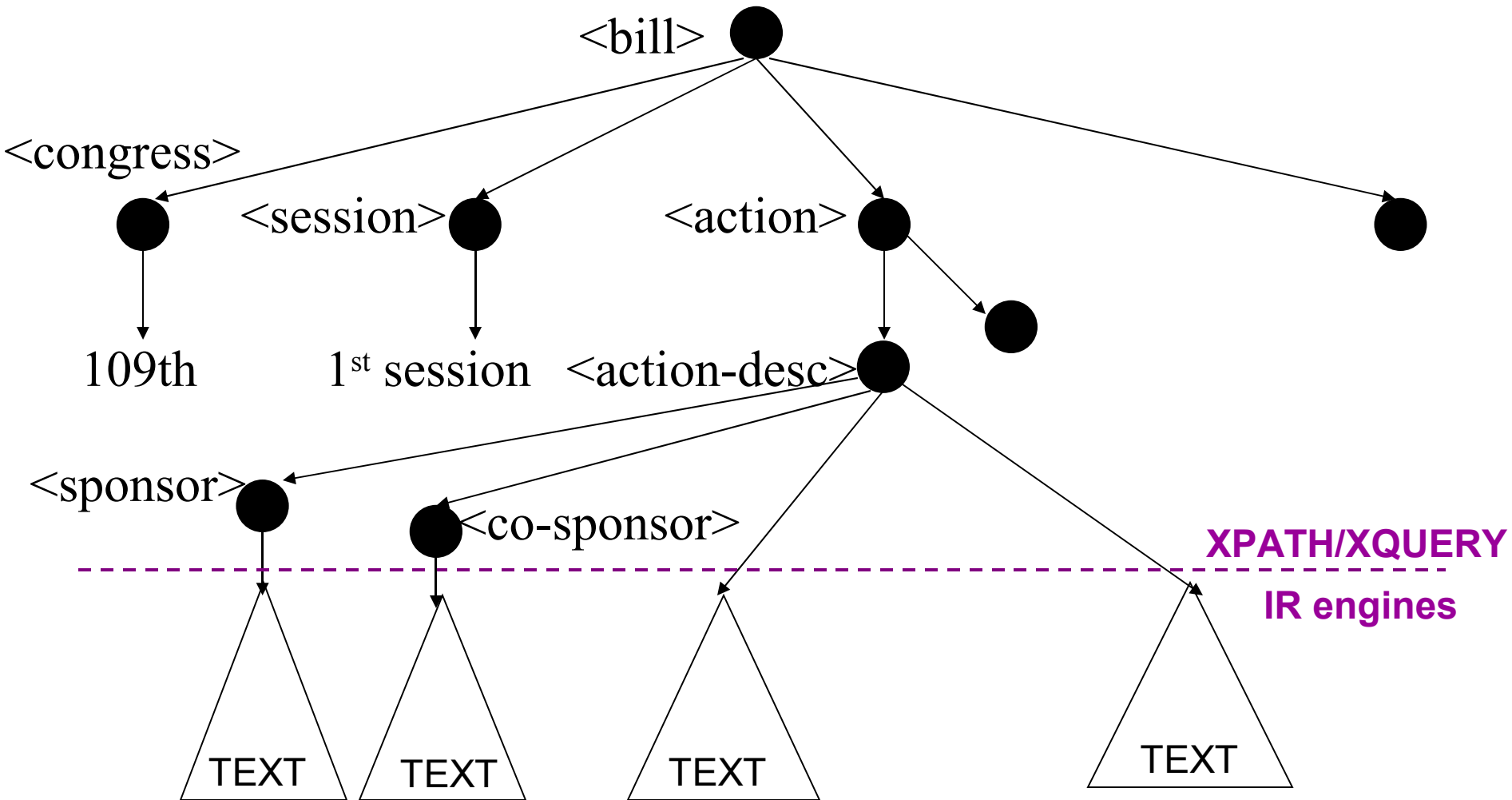


Outline

- Motivation
- Challenges
- Languages
- Research overview



Challenges: DB and IR





Challenges

- Searching over Structure+Text
 - *express complex full-text searches and combine them with structural searches.*
 - *specify a search context and return context.*
- Scores and Ranking
 - *specify a scoring condition,*
 - *possibly over both full-text and structured predicates*
 - *obtain k best results based on query relevance scores*



Motivation

- Current XML query languages are mostly “database” languages
 - Examples: XQuery, XPath
- Provide very rudimentary text/IR support
 - `fn:contains(e, keywords)`
 - Returns true iff element e contains keywords
- No support for complex IR queries
 - Distance predicates, stemming, ...
- No scoring

-
- Full-Text Task Force (FTTF) started in Fall 2002 to extend XQuery with full-text search capabilities: IBM, Microsoft, Oracle, the US Library of Congress.
 - First FTTF documents published on February 14, 2004. (public comments are welcome!):
<http://www.w3.org/TR/xmlquery-full-text-use-cases/>
<http://www.w3.org/TR/xmlquery-full-text-requirements/>
 - XQuery Full-Text highly influenced by TeXQuery.
 - Published a working draft describing the syntax and semantics of XQuery Full-Text on July 9, 2004. Latest version on May 1st 2006 :
<http://www.w3.org/TR/xquery-full-text/>



Example Queries

- From XQuery Full-Text Use Cases Document
 - Find the titles of the books that contain the **phrases** “Usability” and “Web site” in this **order**, in the **same paragraph**, using **stemming** if necessary to match the tokens
 - Find the titles of the books that contain “Usability” and “testing” within a **window of 3 words**, and return them in **score order**
- Such queries are used, e.g. in legal applications



Related Work in IR

- XSEarch, XIRQL, JuruXML, XXL, ELIXIR
 - Not integrated with a powerful language for structured search, such as XQuery
 - Lack expressive power
 - No fully composable
 - Not easily extensible



XML FT Search Definition

- *Context expression*: XML elements searched:
 - pre-defined XML elements.
 - XPath/XQuery queries.
- *Return expression*: XML fragments returned:
 - pre-defined meaningful XML fragments.
 - XPath/XQuery to build answers.
- *Search expression*: FT search conditions:
 - Boolean keyword search.
 - proximity distance, scoping, thesaurus, stop words, stemming.
- *Score expression*:
 - system-defined scoring function.
 - user-defined scoring function.
 - query-dependent keyword weights.



Outline

- Motivation
- Challenges
- Languages
 - XQuery Full-Text
 - INEX
- Research overview



Four Classes of Languages

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



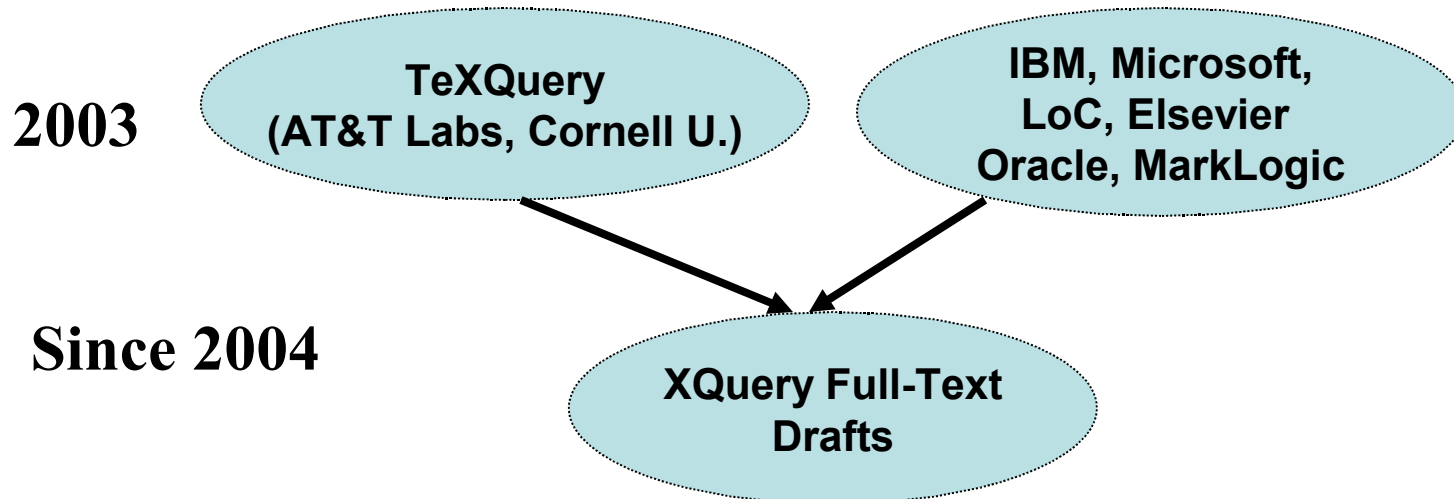
XML Search Languages

- **Keyword-only**
 - Nearest concept (Schmidt, Kersten, Windhouwer, ICDE 2002)
 - XRank (Guo, Botev, Shanmugasundaram, SIGMOD 2003)
 - Schema-free XQuery (Li, Yu, Jagadish, VLDB 2003)
 - INEX Content-Only queries (Trotman, Sigurbjornsson, INEX 2004)
 - XKSearch (Xu & Papakonstantinou, SIGMOD 2005)
- **Tag+Keyword**
 - XSEarch (Cohen, Mamou, Kanza, Sagiv, VLDB 2003)
- **Path+Keyword**
 - XPath 2.0 (<http://www.w3.org/TR/xpath20/>)
 - XIRQL (Fuhr, Großjohann, SIGIR 2001)
 - XXL (Theobald, Weikum, EDBT 2002)
 - NEXI (Trotman, Sigurbjornsson, INEX 2004)



TeXQuery and XQuery Full-Text

- Extends XPath/XQuery with fully composable full-text primitives.
- Scoring and ranking on all predicates.



<http://www.w3.org/TR/xquery-full-text/>



XQuery in a Nutshell

- Functional language. Compositional.
- Input/Output: *sequence of items*
 - *atomic types, elements, attributes, processing instructions, comments,...*
- XPath core navigation language.
- Variable binding.
- Element construction.

return books on XML indexing and ranking sorted by price:

```
for      $item in //books/book
let      $pval := $item//price
where    fn:contains($item/title, "XML")
         and fn:contains($item, "indexing")
         and fn:contains($item, "ranking")
         and $item/price < 50

order by $pval
return  <result>
        {$item/title, $item//authors}
        </result>
```

- sub-string operations: *fn:start-with()*, *fn:end-with()*
- No relevance ranking.



Syntax Overview

Two new XQuery constructs

- **FTContainsExpr**
 - Expresses “Boolean” full-text search predicates
 - Seamlessly composes with other XQuery expressions
- **FTScore**
 - Extension to FLWOR expression
 - Can score FTContainsExpr *and* other expressions



FTContainsExpr and FTScore

- FTContainsExpr := FTWord | FTAnd | FTOr | FTNot | FTMildNot | FTOrder | FTWindow | FTDistance | FTScope | FTTimes | FTSelection (*FTMatchOptions*)*

```
books//section [ . ftcontains (“usability” with stemming occurs 4 times && “Software” case sensitive) window at most 3 ordered with stopwords ]
```

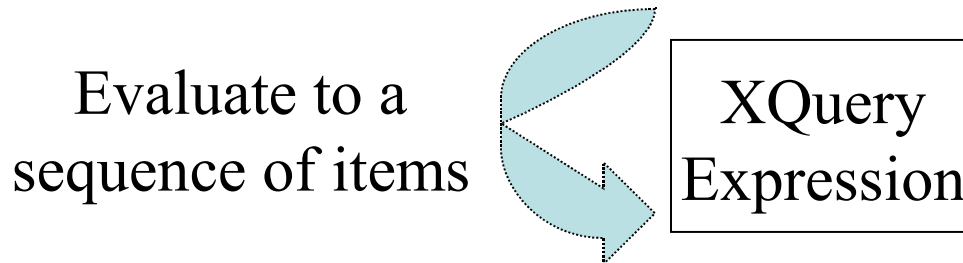
- FTScore

```
for $b SCORE $s in FUZZY  
  //books [ ./title ftcontains “XML” 0.4 and ./section  
    ftcontains (“indexing” with stemming &&  
    “ranking” with thesaurus “synonyms”  
    distance 5 and ./price < 50 ]  
  
order by $s  
return <result score="{$s}"> {$b/title, $b//authors} </result>
```



FTContainsExpr

- Like other XQuery expressions
 - Takes in sequences of items (nodes) as input
 - Produces a sequence of items (nodes) as output



- Can seamlessly compose with other XQuery expressions



FTContainsExpr

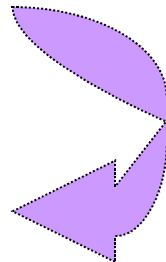
ContextExpr **ftcontains** FTSelection

- ContextExpr (any XQuery expression) is **context spec**
- FTSelection is **search spec**
- Returns true iff at least one node in ContextExpr satisfies the FTSelection
- Examples
 - `//book ftcontains 'Usability' && 'testing' distance 5`
 - `//book[./content ftcontains 'Usability' with stems]/title`
 - `//book ftcontains /article[author='Dawkins']/title`

- Encapsulates all full-text conditions in FTContainsExpr
- Works in a new data model called **AllMatch**
 - Operates on positions within XML nodes (more fine grained than XQuery data model)
 - Fully composable; similar to composition of relational (and XML) operators!



FTSelection



Evaluate to
AllMatch



FTSelection Composability

- 'Usability'
- /book[author='Dawkins']/title
- 'Usability' && /book[author='Dawkins']/title
- ('Usability' && /book[author='Dawkins']/title)
same sentence
- ('Usability' && /book[author='Dawkins']/title)
same sentence window 5
- All of these evaluate to an AllMatch!
 - Allows arbitrary composition of full-text primitives



FTMatchOption

- Can be applied on any FTSelection to specify aspects such as stemming, thesauri, case, etc.
 - Fully composable with other context modifiers and FTSelections
- Examples
 - ‘Usability’ && ‘testing’ **with stems**
 - ‘Usability’ && ‘testing’ with stems window 5 **without stop words**
 - ‘Usability’ && ‘testing’ with stems window 5 without stop words **case insensitive**



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

Example

```
FOR $b SCORE $s in  
    /pub/book[. ftcontains "Usability" && "testing"]  
ORDER BY $s  
RETURN <result score={$s}> $b </result>
```



FTScoreExpr

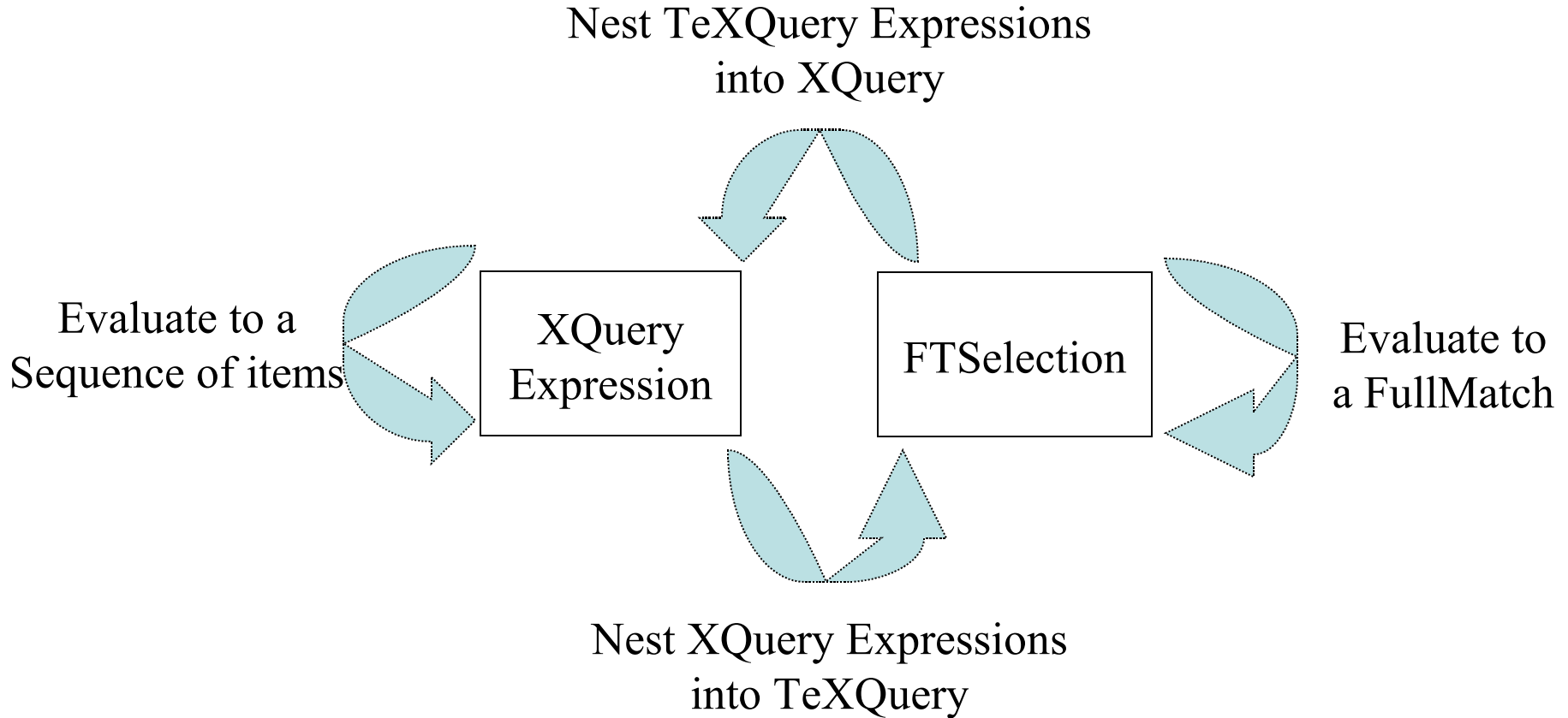
In any order {
FOR \$v [SCORE \$s]? [AT \$i]? IN [FUZZY] Expr
LET ...
WHERE ...
ORDER BY ...
RETURN

Example

```
FOR $b SCORE $s in FUZZY  
    /pub/book[. ftcontains "Usability" && "testing"]  
ORDER BY $s  
RETURN <result score={ $s }> $b </result>
```



Semantics Issues



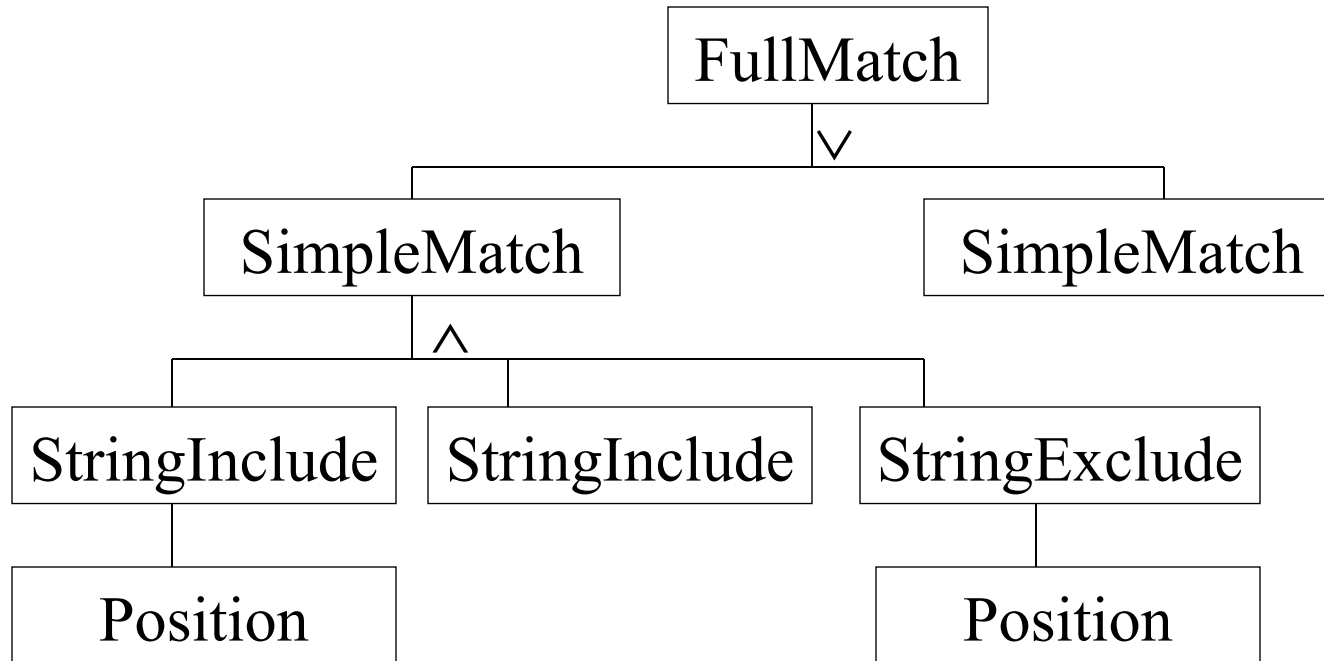


FullMatch Overview

- FTSelections are fully composable
- Extensible with respect to new FTSelections
 - Only have to define semantics w.r.t. FullMatch
- Clean way to specify semantics of FTSelections
 - Like specifying semantics of relational operators
- Provides basis for optimizing complex queries



FullMatch



- FullMatch can be interpreted as a propositional formula over word positions in DNF



Sample Document

```
<book(1) id(2)='`1000(3)''>
  <author(4)>Elina(5) Rose(6)</author(7)>
  <content(8)>
    <p(9)> The(10) usability(11) of(12) software(13)
      measures(14) how(15) well(16) the(17)
      software(18) provides(19) support(20) for(21)
      quickly(22) achieving(23) specified(24)
      goals(25). </p(26)>
    <p(27)>The(28) users(29) must(30) not(31) only(32)
      be(33) well-served(34), but(35) must(36)
      feel(37) well-served(38).</p(39)>
  </content(40)>
</book(41)>
```



Sample Query

```
$doc ftcontains  
( 'usability' with stems &&  
  'Rose' )  
window at most 10
```



Sample FTSelection

```
('usability' with stems &&  
'Rose')  
window at most 10
```

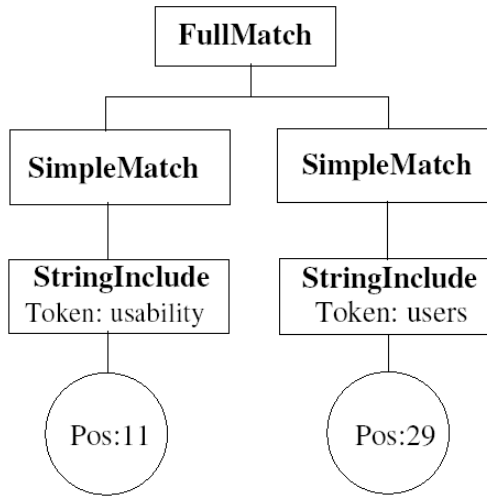


Semantics of FTStringSelection

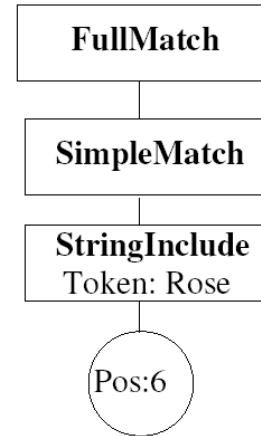
```
<book(1) id(2)='`1000(3)''>
  <author (4)>Elina(5) Rose(6)</author(7)>
  <content(8)>
    <p(9)> The(10) usability(11) of(12) software(13)
      measures(14) how(15) well(16) the(17)
      software(18) provides(19) support(20) for(21)
      quickly(22) achieving(23) specified(24)
      goals(25). </p(26)>
    <p(27)>The(28) users(29) must(30) not(31) only(32)
      be(33) well-served(34), but(35) must(36)
      feel(37) well-served(38).</p(39)>
  </content(40)>
</book(41)>
```



Semantics of FStringSelection



'usability' with stems



'rose'

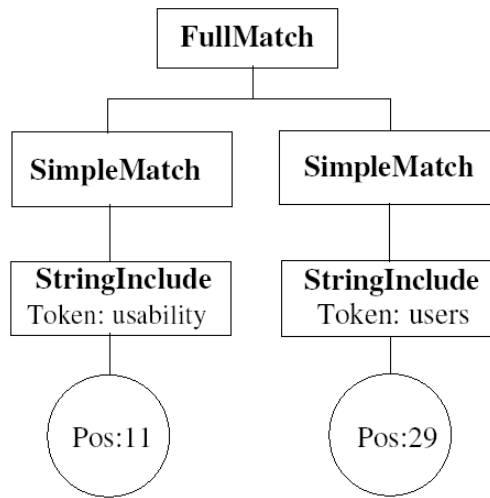


Sample FTSelection

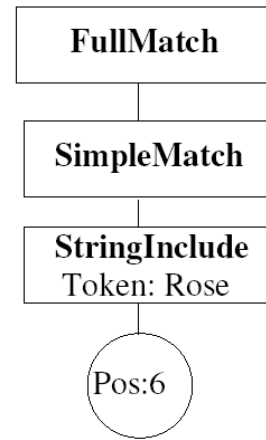
```
('usability' with stems &&  
'Rose')  
window at most 10_
```



Semantics of FTAndConnective



×

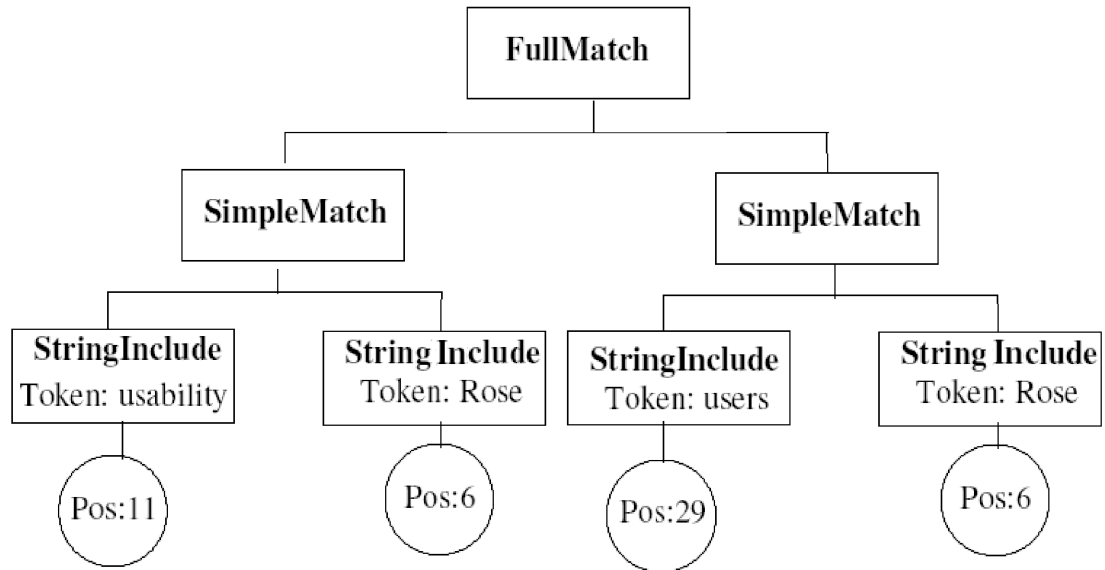


'usability' with stems

'Rose'



Semantics of FTAndConnective



'usability' with stems && 'Rose'



Sample FTSelection

```
('usability' with stems &&  
'Rose')  
window at most 10
```

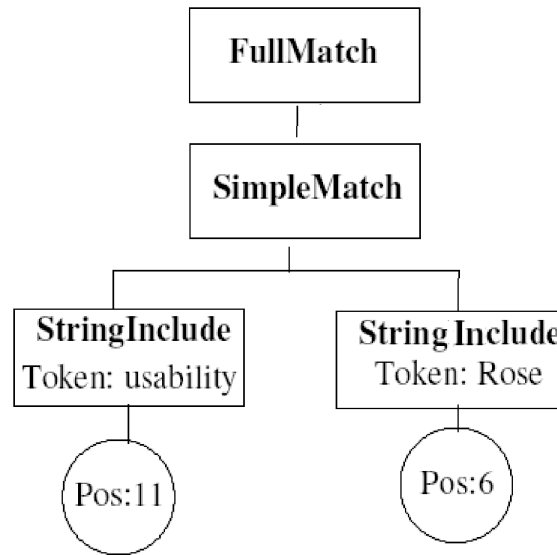


Semantics of FTWindowSelection

```
<book(1) id(2)='`1000(3)''>
  <author (4)>Elina(5) Rose(6)</author(7)>
  <content(8)>
    <p(9)> The(10) usability(11) of(12) software(13)
      measures(14) how(15) well(16) the(17)
      software(18) provides(19) support(20) for(21)
      quickly(22) achieving(23) specified(24)
      goals(25). </p(26)>
    <p(27)>The(28) users(29) must(30) not(31) only(32)
      be(33) well-served(34), but(35) must(36)
      feel(37) well-served(38).</p(39)>
  </content(40)>
</book(41)>
```



Semantics of FTWindowSelection



`('usability' with stems && 'Rose')`
`window at most 10`



FullMatch Benefits

- FullMatch has a hierarchical structure
- Thus FullMatch can be represented as XML
- Semantics of FTSelections can be specified as transformation from input XML FullMatches to the output XML FullMatch
- Thus, semantics of FTSelections can be specified in XQuery itself!
- Full-text conditions and structural conditions represented in the same framework
 - Enables joint optimization and evaluation



- GalaTeX Overview:**
- [About XQuery Full-Text](#)
 - [About GalaTeX](#)
 - [The GalaTeX Team](#)
 - [Thanks!](#)

- Documentation:**
- [W3C XQuery 1.0 / XPath 2.0 Full-Text](#)
 - [XQuery 1.0 / XPath 2.0 Full-Text](#)
 - [W3C Use Cases](#)
 - [Full-Text Requirements](#)

- GalaTeX Documentation**
- [Technical Reports](#)
 - [Presentations](#)

[Galax / XQuery Documentation](#)



An XML Full Text Search Engine

Welcome to the GalaTeX Website!

GalaTeX is a conformant implementation of the [W3C XQuery 1.0 and XPath 2.0 Full-Text](#) standard proposal, an XML full text search language. GalaTeX is built on top of [Galax](#).

GalaTeX's primary goal is to serve as a reference implementation for the XQuery 1.0 and XPath 2.0 Full-Text extension language.

This site contains information about GalaTeX, some documentation, on-line demos, W3C full-text uses cases and the semantic implementation of the full-text primitives. GalaTeX and this Web site keep evolving so please come back again soon.

For any questions, please contact [Emiran Curtmola](#), [Sihem Amer-Yahia](#) or [Mary Fernández](#).

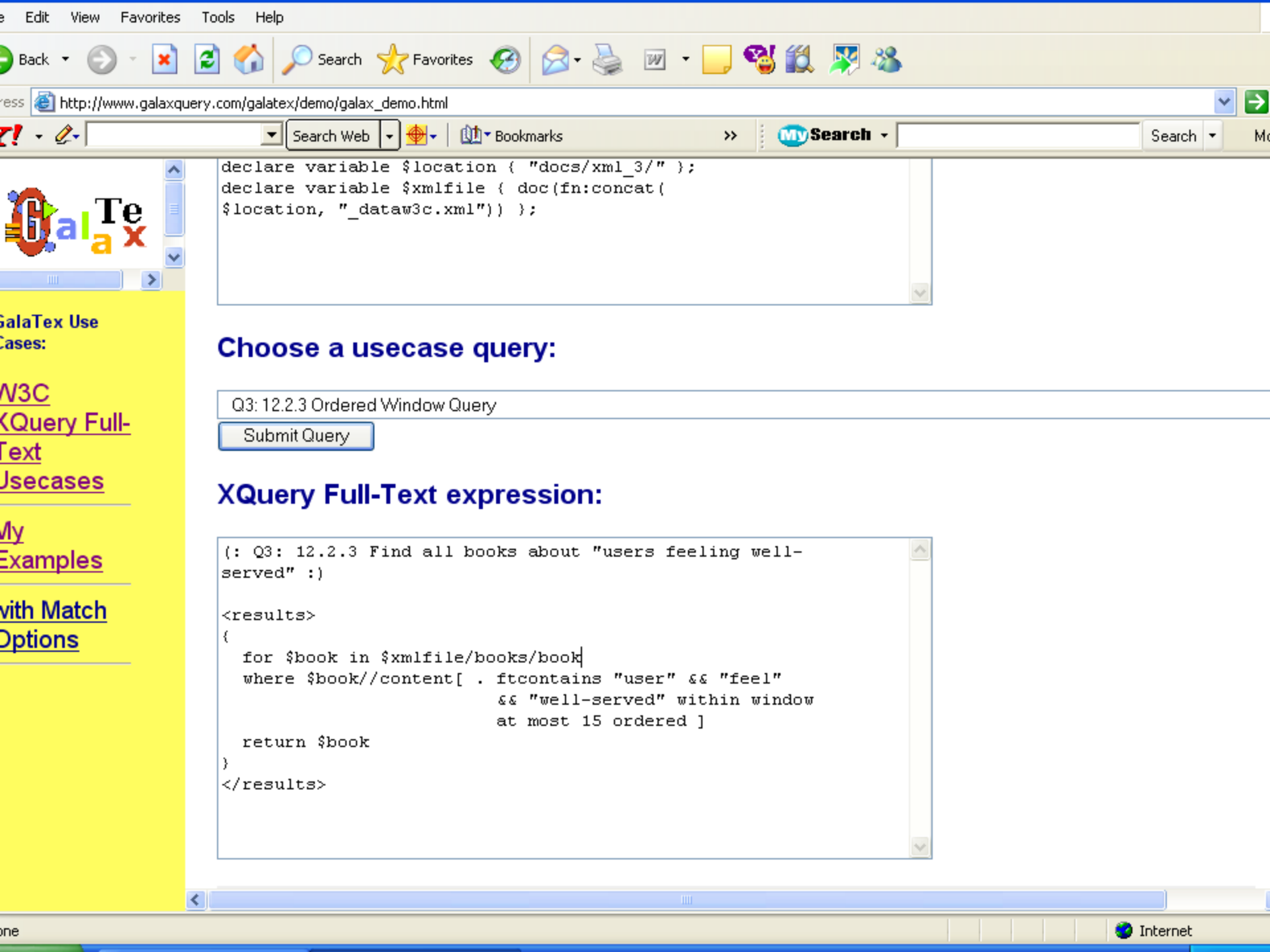
Latest News **Hot**

- **August, 2005.** The GalaTeX pre-release 0.5.1 source-code is available! If you are a GalaTeX user or experimenter, please [drop us a note](#).
- **August, 2004.** The first prototype for [GalaTeX demo](#) is now available! This

- GalaTeX Live!**
- [GalaTeX Demo](#)
 - [\[Galax Demo \]](#)

- Download GalaTeX**
- [GalaTeX 0.5.1](#)
 - [License](#)

- Tech. Center**
- [Research Pointers](#)



- GalaTeX Use Cases:
- [W3C XQuery Full-Text Usecases](#)
- [My Examples with Match Options](#)

```
declare variable $location { "docs/xml_3/" };
declare variable $xmlfile { doc(fn:concat(
$location, "_dataw3c.xml")) };
```

Choose a usecase query:

Q3: 12.2.3 Ordered Window Query

Submit Query

XQuery Full-Text expression:

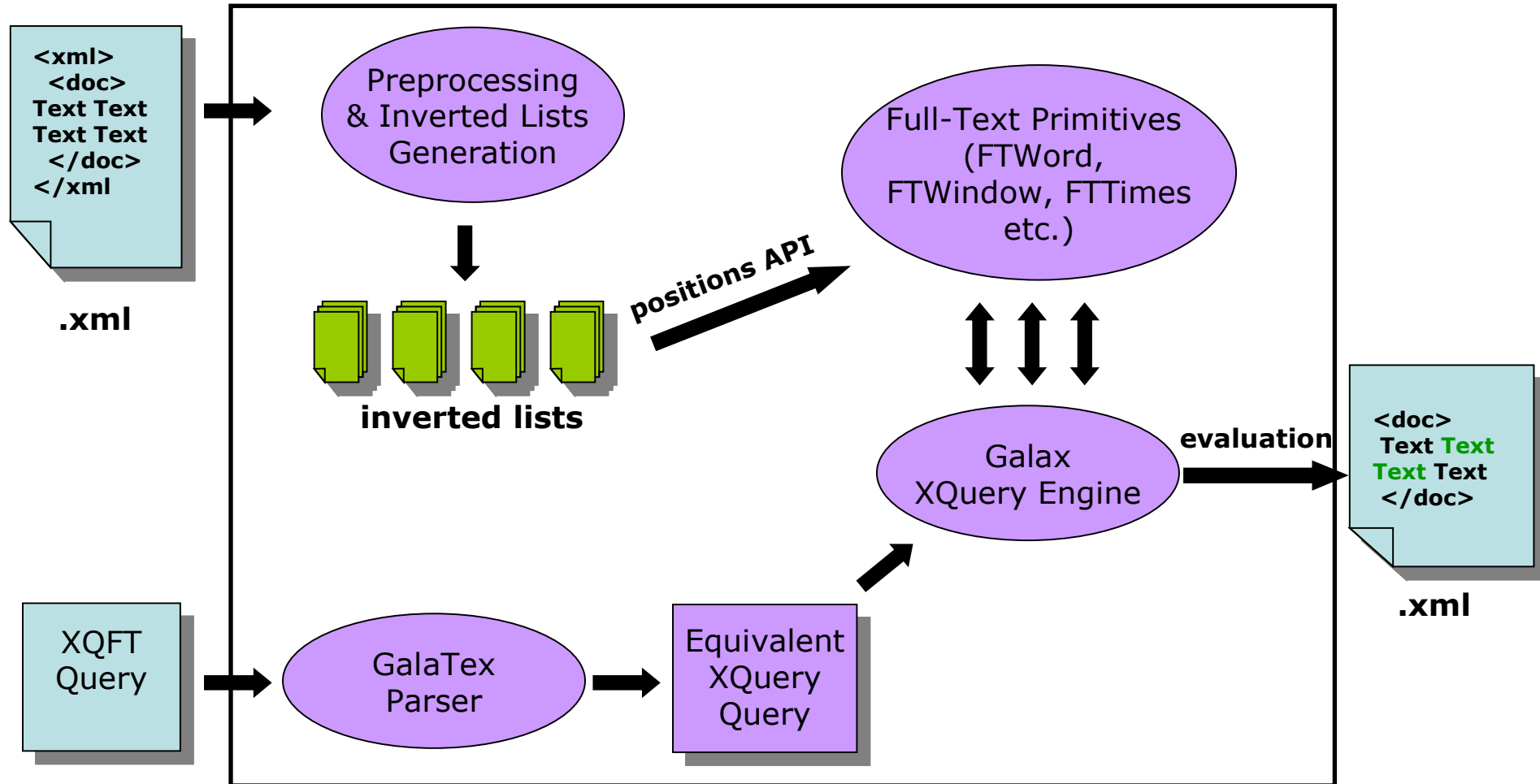
```
(: Q3: 12.2.3 Find all books about "users feeling well-served" :)

<results>
{
  for $book in $xmlfile/books/book
  where $book//content[ . ftcontains "user" && "feel"
                        && "well-served" within window
                        at most 15 ordered ]

  return $book
}
</results>
```



GalaTex <http://www.galaxquery.org/galatex>





Outline

- Motivation
- Challenges
- Languages
 - XQuery Full-Text
 - INEX
- **Research overview**



INitiative for the Evaluation of XML retrieval

- Evaluate effectiveness of content-oriented XML retrieval systems
- Ongoing effort to define:
 - documents
 - queries (topics)
 - relevance assessments
 - metrics



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



INEX document

<article>

<fno>A1002</fno>

<doi>10.1041/A1002s-2004</doi>

<ti>IEEE ANNALS OF THE HISTORY OF COMPUTING</ti>

<issn>1058-6180</issn>

<obi> Published by the IEEE Computer Society</obi>

<mo>JANUARY-MARCH</mo>

<yr>2004</yr>

<bdy>

<sec>

<ip1>Some 25 years ago, 26 if we are to be precise, a small group of computer scientists decided

that their discipline not only had a past, it had a history. A history is a very different thing from a

past. A past is a series of events; some good, bad, pleasing, embarrassing,.....

</ip1>

<p align="left" ind="none">A history, however, looks at the deep trends of modern life and asks where they have been, where they are now, and where they are going. It is a discipline that looks to the future as much as it retells the story of the past. Those of us involved with the <it>Annals</it> believe that the stored program electronic computer helps us understand

almost

</p>



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



NEXI

- Narrowed Extended XPath I
 - Designed for content-oriented XML search (i.e. “aboutness”)
 - query conditions on structure interpreted as hints to find content
- IEEE document collection growth
 - 12,107 to 659,388 documents
 - 8M to 30M elements
 - 494MB to 60GB (total size)

+ontologies -anonyms

//article [about (., ontologies)]

//article [about (., ontologies)]//sec [about (., ontologies case study)]



INEX topic id=202

```
<inex_topic topic_id="202" query_type="CO+S" ct_no="1" >
```

```
<InitialTopicStatement>I'm interested in knowing how ontologies are used to encode knowledge in real world scenarios. I'm writing a report on the use of ontologies. I'm particularly interested in knowing what sort of concepts and relations people use in their ontologies.
```

```
</InitialTopicStatement>
```

```
<title>ontologies case study</title>
```

```
<castitle>//article[about(., ontologies)]//sec[about(., ontologies case study)]</castitle>
```

```
<description>Case studies in the use of ontologies</description>
```

```
<narrative>I'm writing a report on the use of ontologies. I'm interested in knowing how ontologies are used to encode knowledge in real world scenarios. I'm particularly interested in knowing what sort of concepts and relations people use in their ontologies. I'm not interested in general ontology frameworks or technical details about tools for ontology creation or management. An example relevant result contains a description of the real world phenomena described by the ontology and also lists some of the concepts used and relations between concepts.
```

```
</narrative>
```

```
</inex_topic>
```



Relevance

- Precision and recall are not enough:
 - relevance is a binary property (items are relevant or not)
 - relevance of one item independent from other items
 - user spends a constant time on each element
 - user looks at an ordered list and stops at some point
- The problem with retrieving elements:
 - specificity and exhaustiveness matter
 - overlap between elements: return parent (2005) / child (2006)?
 - size of retrieved elements varies => time spent varies
 - near-misses: some elements could be found by browsing



Metrics

- inex-eval (precall)
 - quantisation functions to capture specificity and exhaustivity
 - ignores possible overlap between elements
- inex-eval-ng
 - incorporate overlap and element size in precision and recall
 - consider only increment in text size of elements already seen
- cumulative gain
 - favors specificity
 - computed as the sum of relevance score up to that element
 - favors deeper nodes



Outline

- Motivation
- Challenges
- Languages
- **Research Overview**



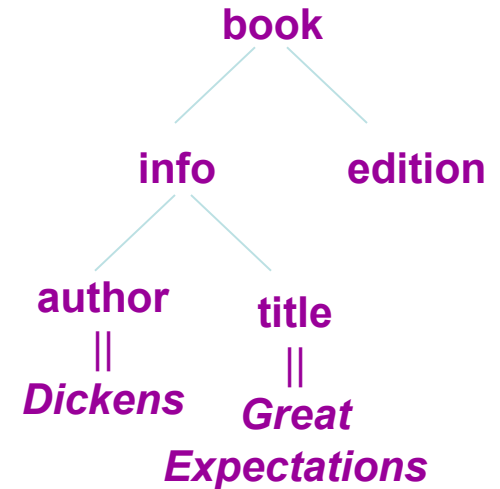
Some papers

- Designed **TeXQuery** (Amer-Yahia, Botev, Shanmugasundaram, WWW 2004), and **XQuery Full-Text**, a full-text extension of XPath/XQuery (Amer-Yahia et al, <http://www.w3.org/TR/xquery-full-text/>, W3C Draft) and developed **GalaTex**, a conformant open-source implementation. (Curtmola, Amer-Yahia, Brown, Fernandez, XIME-P 2005)
- *Beyond DB*: Formalized a query semantics that consistently extends classical XPath semantics to account for XPath relevance ranking. **FleXPath** (Amer-Yahia, Lakshmanan, Pandit, SIGMOD 2004)
- *Beyond IR*: Developed a family of scoring methods for XML on both structure and content that are consistent with $tf*idf$. (Amer-Yahia, Koudas, Marian, Srivastava, Toman, VLDB 2005)
- Developed efficient algorithms for topK processing. **Whirlpool** (Marian, Amer-Yahia, Koudas, Srivastava, ICDE 2005)



Example Query

***//book [./info [./author
ftcontains "Dickens" and ./title
ftcontains "Great Expectations"]
and ./edition]***





Some XPath Relaxations

- Examples of atomic relaxations:
 - Leaf node deletion
 - Edge generalization
 - Subtree promotion
 - ...



Data



Y! Query Representation

```
//book [ ./info [ ./author  
ftcontains "Dickens" and ./title  
ftcontains "Great Expectations" ]  
and ./edition ]
```



*pc(\$1,\$2) and pc(\$2,\$3) and pc(\$2,\$4) and pc(\$1,\$5) and
(\$1.tag = book) and (\$2.tag = info) and (\$3.tag = author) and
(\$4.tag = title) and (\$5.tag = edition) and contains(\$3, "Dickens") and
contains(\$4, "Great Expectations")*



XPath Relaxation Algorithm

- Logical representation of query using predicates on structure and content.
- Compute *query closure* using inference rules below:
 - $pc(\$x, \$y)$ implies $ad(\$x, \$y)$
 - $ad(\$x, \$y), ad(\$y, \$z)$ implies $ad(\$x, \$z)$
 - $ad(\$x, \$y), contains(\$y, FTExp)$ implies $contains(\$x, FTExp)$
 - ...
- Drop predicates.
- Compute *query core* (unique).



Example of XPath Relaxation

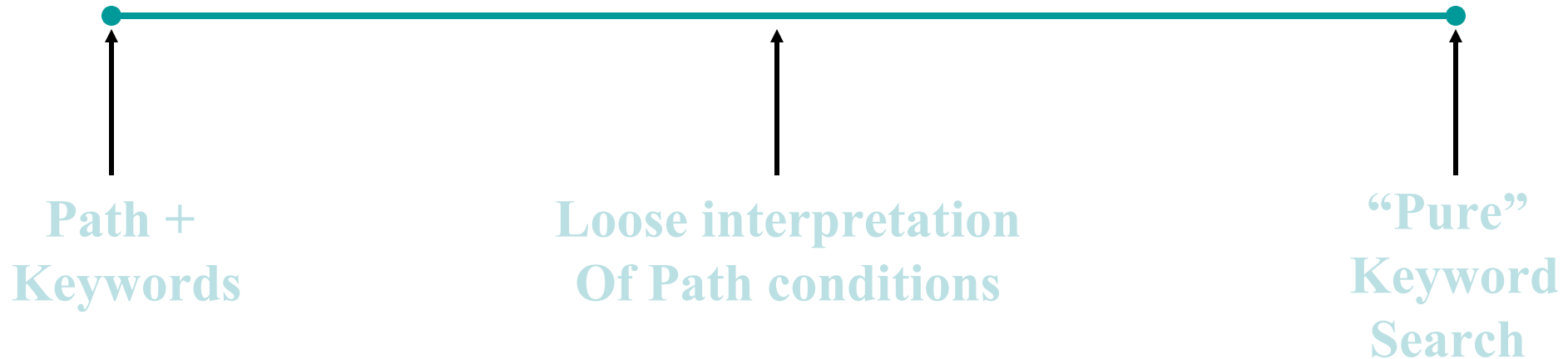


*pc(\$1,\$2) and pc(\$2,\$3) and
pc(\$2,\$4) and pc(\$1,\$5) and
(\$1.tag = article) and (\$2.tag = info)
and (\$3.tag = author) and
(\$4.tag = title) and (\$5.tag = edition)
and contains(\$3, "Dickens") and
contains(\$4, "Great Expectations")*

*pc(\$1,\$2) and ad(\$1,\$3) and
pc(\$2,\$4) and ad(\$1,\$5) and
(\$1.tag = article) and (\$2.tag = info)
and (\$3.tag = author) and
(\$4.tag = title) and (\$5.tag = edition)
and contains(\$3, "Dickens") and
contains(\$4, "Great Expectations")*



Spanning XPath Relaxations



- Framework for defining new relaxations.
- Orthogonal to approximation on content.
- Answers to relaxed query contain answers to exact query.
- Score of answer to relaxed query should be no higher than score of answer to more exact query.



Adaptation of $tf*idf$ to XML

Document Retrieval	XML Retrieval
Document	XML fragment (result is a subtree rooted at an element with a given tag and satisfying content and structure in query)
Keyword	Path + Keyword
<i>idf</i> (inverse document frequency) is a function of the fraction of documents that contain the keyword	<i>idf</i> is a function of the fraction of returned fragments that match the query tree pattern
<i>tf</i> (term frequency) is a function of the number of occurrences of the keyword in the document	<i>tf</i> is a function of the number of ways the query tree pattern matches the returned fragment



A Family of Scoring Methods

- Binary scoring
 - Low quality
 - Fast computation
- Path scoring
- Twig scoring
 - High quality
 - Expensive computation





What does XML mean anyway?

- EDS: Encyclopedia of Database Systems
- Alphabetical organization of ~ 1000 entries
 - definitions and illustrations of basic terminology, concepts, methods, and algorithms,
 - references to literature, and cross-references to other entries and journal articles.
 - Not a textbook
 - <http://refworks.springer.com/database-systems>
- April 15: Initial list of entries for XML
- Send to sihem@yahoo-inc.com