# Efficient Filtering of XML Documents for Selective Dissemination of Information

By Mehmet Altinel and Michael J. Franklin

Presented by: Weimin Li
October 16, 2002

---

# Outline

- XML-based SDI

- XFilter Structure

- Enhanced Filtering Algorithm
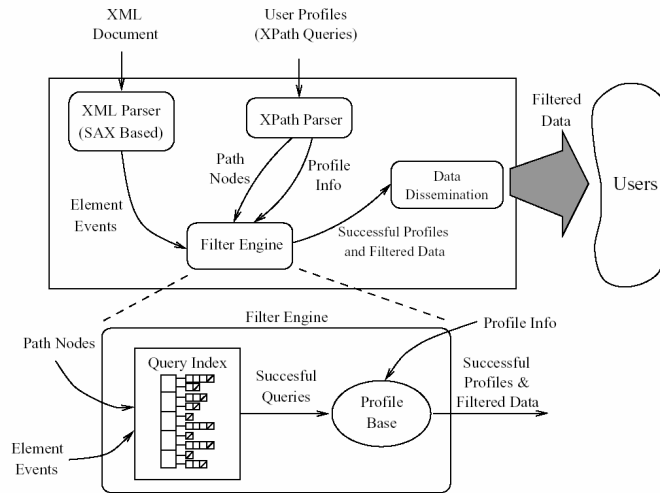
- Performance

- Comments

# XML-based SDI

- Selective Dissemination of Information
  Distribute the right information to users based upon their profiles (interests)

- Approaches in the Information Retrieval (IR) community
  Match keywords: Boolean or Similarity-based

- Approaches in the database community
  Use queries in the context of Continuous Queries (CQ)

---
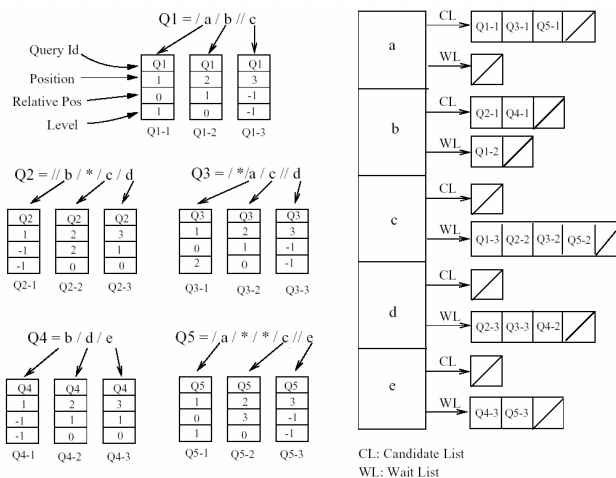
# XML-based SDI

- Why XML?
  eXtensible Markup Language derived from SGML
  - Semi-structured
  - Self-described
  XML becomes a standard format in data exchange
  The cost: the complexity to process XML documents

- XPath as a Profile Language
  A language to navigate or address parts in an XML documents
  "/catalog/product//name"

# XFilter Structure

XML Document     User Profiles (XPath Queries)

XML Parser (SAX Based)

XPath Parser

Path Nodes

Profile Info

Element Events

Filter Engine

Data Dissemination

Filtered Data

Users

Successful Profiles and Filtered Data

Filter Engine

Path Nodes

Query Index

Successful Queries

Profile Base

Profile Info

Successful Profiles & Filtered Data

Element Events

# XFilter

- Finite State Machine (FSM)

$Q1 = / a / b // c$

Query Id
Position
Relative Pos
Level

| Q1 | Q1 | Q1 |
|----|----|----|
| 1 | 2 | 3 |
| 0 | 1 | -1 |
| 1 | 0 | -1 |

Q1-1   Q1-2   Q1-3

$Q2 = // b / * / c / d$

| Q2 | Q2 | Q2 |
|----|----|----|
| 1 | 2 | 3 |
| -1 | 2 | 1 |
| -1 | 0 | 0 |

Q2-1   Q2-2   Q2-3

$Q3 = / * / a / c // d$

| Q3 | Q3 | Q3 |
|----|----|----|
| 1 | 2 | 3 |
| 0 | 1 | -1 |
| 2 | 0 | -1 |

Q3-1   Q3-2   Q3-3

$Q4 = b / d / e$

| Q4 | Q4 | Q4 |
|----|----|----|
| 1 | 2 | 3 |
| -1 | 1 | 1 |
| -1 | 0 | 0 |

Q4-1   Q4-2   Q4-3

$Q5 = / a / * / * / c // e$

| Q5 | Q5 | Q5 |
|----|----|----|
| 1 | 2 | 3 |
| 0 | 3 | -1 |
| 1 | 0 | -1 |

Q5-1   Q5-2   Q5-3

| | CL | | | | |
|---|----|---|---|---|---|
| a | | Q1-1 | Q3-1 | Q5-1 | |
| | WL | | | | |
| b | CL | Q2-1 | Q4-1 | | |
| | WL | Q1-2 | | | |
| c | CL | | | | |
| | WL | Q1-3 | Q2-2 | Q3-2 | Q5-2 |
| d | CL | | | | |
| | WL | Q2-3 | Q3-3 | Q4-2 | |
| e | CL | | | | |
| | WL | Q4-3 | Q5-3 | | |

CL: Candidate List
WL: Wait List

# Enhanced Filtering Algorithms

- List Balancing

  Skewed lengths of the Candidate Lists do not provide little selectivity

  Select a "pivot" as the start element node

- Prefiltering

  The idea: pre-delete the queries that are impossible to match the document

# Performance

- Four policies

  Basic

  Prefiltering+Basic

  List Balance

  Prefiltering+List Balance

- Prefiltering+List Balance works best in nearly all cases

# Comments

- Contributions
  - An modified FSM
    Run and evaluate all the queries in one FSM at the same time

  - Algorithms
    Basic algorithm and List Balance and Prefiltering algorithms enable efficiently filter XML documents

# Future Work

- Adaptive
  - Source data and queries may change with time
  - Need reevaluate and re-balance the Candidate Lists

- Extract parts of an XML document
  - Save bandwidth
  - Need more complex algorithm