# XML Retrieval

# Overview

① Introduction

② Basic XML concepts

③ Challenges in XML IR

④ Vector space model for XML IR

⑤ Evaluation of XML IR

# IR and relational databases

IR systems are often contrasted with relational databases (RDB).

- Traditionally, IR systems retrieve information from *unstructured text* ("raw" text without markup).

- RDB systems are used for querying *relational data*: sets of records that have values for predefined attributes such as employee number, title and salary.

|  | RDB search | unstructured IR |
|---|---|---|
| objects | records | unstructured docs |
| main data structure | table | inverted index |
| model | relational model | vector space & others |
| queries | SQL | free text queries |

Some structured data sources containing text are best modeled as structured documents rather than relational data (Structured retrieval).

# Structured retrieval

Basic setting: queries are structured or unstructured; documents are structured.

## Applications of structured retrieval

Digital libraries, patent databases, blogs, tagged text with entities like persons and locations (named entity tagging)

## Example

- Digital libraries: *give me a full-length article on fast fourier transforms*

- Patents: *give me patens whose claims mention RSA public key encryption and that cite US patent 4,405,829*

- Entity-tagged text: *give me articles about sightseeing tours of the Vatican and the Coliseum*

4

# Why RDB is not suitable in this case

Three main problems

1.  An unranked system (DB) would return a potentially large number of articles that mention the Vatican, the Coliseum and sightseeing tours without ranking them by relevance to query.

2.  Difficult for users to precisely state structural constraints – may not know which structured elements are supported by the system.

    *tours* AND (COUNTRY: *Vatican* OR

    LANDMARK: *Coliseum*)?

    *tours* AND (STATE: *Vatican* OR BUILDING: *Coliseum*)?

3.  Users may be completely unfamiliar with structured search and advanced search interfaces or unwilling to use them.

Solution: adapt ranked retrieval to structured documents to address these problems.

# Structured Retrieval

## RDB search, Unstructured IR, Structured IR

|  | RDB search | unstructured retrieval | structured retrieval |
|---|---|---|---|
| objects | records | unstructured docs | trees with text at leaves |
| main data structure | table | inverted index | ? |
| model | relational model | vector space & others | ? |
| queries | SQL | free text queries | ? |

Standard for encoding structured documents: Extensible Markup Language (XML)

- structured IR → XML IR

- also applicable to other types of markup (HTML, SGML, …)

# XML document

- Ordered, labeled tree
- Each node of the tree is an XML element, written with an opening and closing XML tag (e.g. <title…>, </title…>)
- An element can have one or more XML attributes (e.g. number)
- Attributes can have values (e.g. vii)
- Attributes can have child elements (e.g. title, verse)

<play>

<author>Shakespeare</author>

<title>Macbeth</title>

<act number="I">

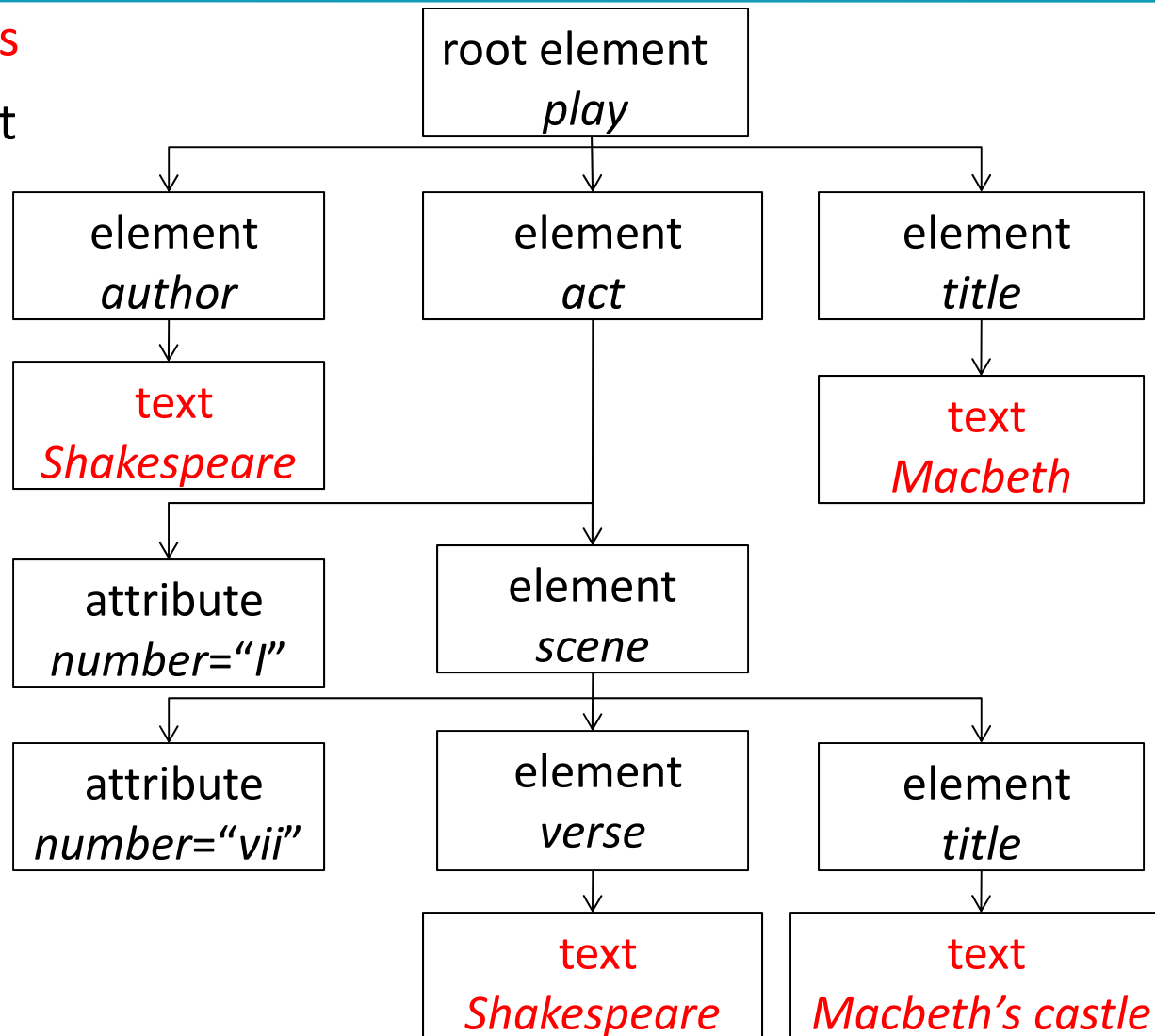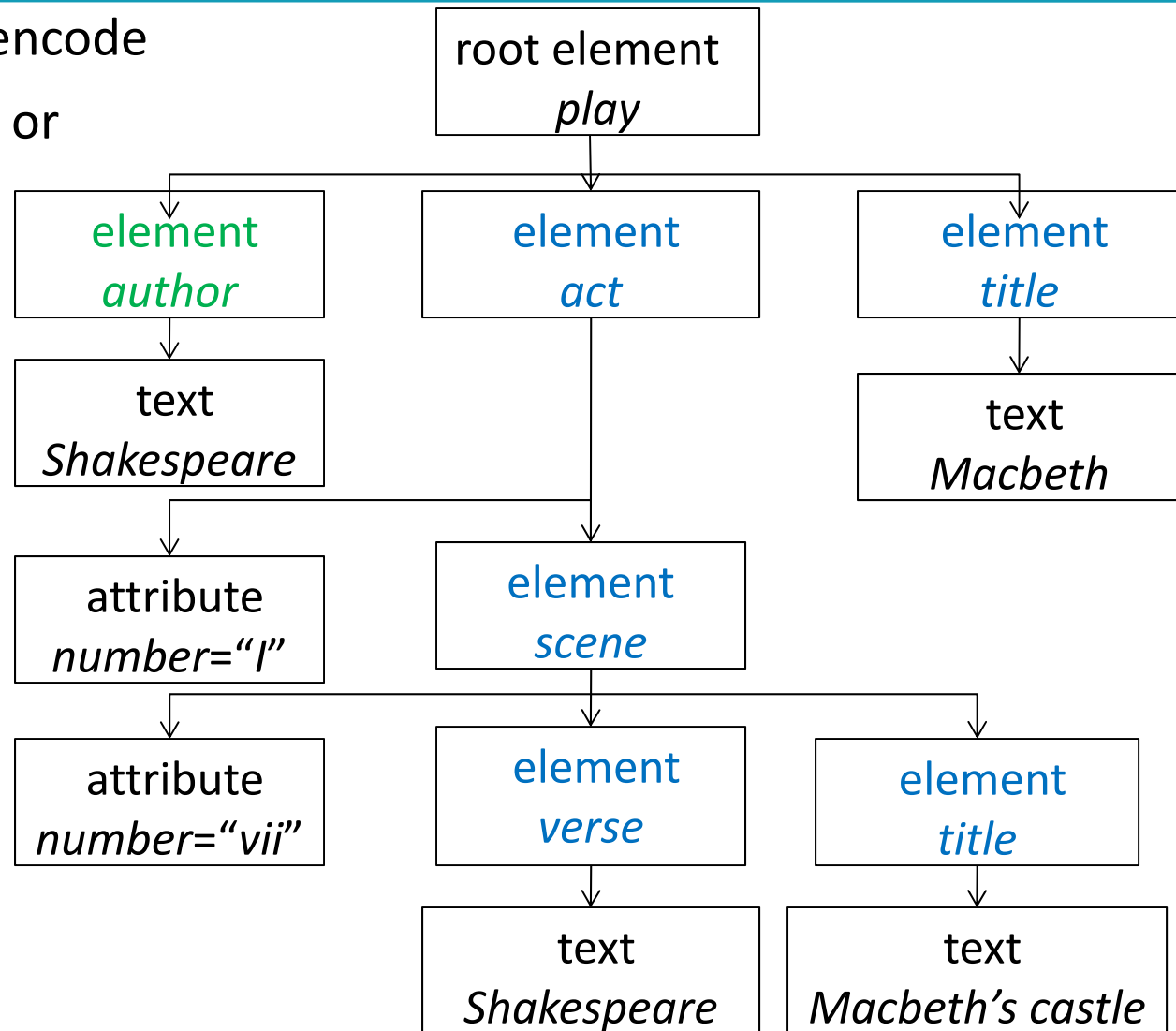<scene  number=""vii">

<title>Macbeth's castle</title>

<verse>Will I with wine

…</verse>

</scene>

</act>

</play>

# XML document



```
                    root element
                        play
        |               |               |
    element         element         element
     author            act             title
        |               |               |
      text                            text
   Shakespeare                       Macbeth
        |               |
    attribute       element
  number="I"          scene
        |               |               |
    attribute       element         element
  number="vii"        verse            title
                        |               |
                      text            text
                   Shakespeare    Macbeth's castle
```

8

# XML document

The leaf nodes

consist of text

root element
*play*

element
*author*

element
*act*

element
*title*

text
*Shakespeare*

text
*Macbeth*

attribute
*number="I"*

element
*scene*

attribute
*number="vii"*

element
*verse*

element
*title*

text
*Shakespeare*

text
*Macbeth's castle*

9

# XML document

The internal nodes encode

document structure or

metadata functions

| root element |
| *play* |

| element | element | element |
| *author* | *act* | *title* |

| text | | text |
| *Shakespeare* | | *Macbeth* |

| attribute | element |
| *number="I"* | *scene* |

| attribute | element | element |
| *number="vii"* | *verse* | *title* |

| text | text |
| *Shakespeare* | *Macbeth's castle* |

10

# XML basics

- **XML Documents Object Model (XML DOM)**: standard for accessing and processing XML documents

  - The DOM represents elements, attributes and text within elements as nodes in a tree.

  - With a DOM API, we can process an XML documents by starting at the root element and then descending down the tree from parents to children.

- **XPath**: standard for enumerating path in an XML document collection.

  - We will also refer to paths as XML contexts or simply contexts

- **Schema**: puts constraints on the structure of allowable XML documents. E.g. a schema for Shakespeare's plays: scenes can occur as children of acts.

  - Two standards for schemas for XML documents are: XML DTD (document type definition) and XML Schema.

11

# First challenge: document parts to retrieve

Structured or XML retrieval: users want us to return parts of documents (i.e., XML elements), not entire documents as IR systems usually do in unstructured retrieval.

## Example

If we query Shakespeare's plays for *Macbeth's castle*, should we return the scene, the act or the entire play?

- In this case, the user is probably looking for the scene.
- However, an otherwise unspecified search for *Macbeth* should return the play of this name, not a subunit.

Solution: structured document retrieval principle

# Structured document retrieval principle

**Structured document retrieval principle**

One criterion for selecting the most appropriate part of a document:
*A system should always retrieve the most specific part of a document answering the query.*

- Motivates a retrieval strategy that returns the smallest unit that contains the information sought, but does not go below this level.

- Hard to implement this principle algorithmically. E.g. query: title:*Macbeth* can match both the title of the tragedy, *Macbeth*, and the title of Act I, Scene vii, *Macbeth's castle*.

  - But in this case, the title of the tragedy (higher node) is preferred.

  - Difficult to decide which level of the tree satisfies the query.

# Second challenge: document parts to index

Central notion for indexing and ranking in IR: documents unit or **indexing unit**.

- In unstructured retrieval, usually straightforward: files on your desktop, email massages, web pages on the web etc.

- In structured retrieval, there are four main different approaches to defining the indexing unit
  1. non-overlapping pseudodocuments
  2. top down
  3. bottom up
  4. all

# XML indexing unit: approach 1

Group nodes into non-overlapping pseudodocuments.



Indexing units: books, chapters, section, but without overlap.

Disadvantage: pseudodocuments may not make sense to the user

because they are not coherent units.

# XML indexing unit: approach 2

Top down (2-stage process):

1. Start with one of the latest elements as the indexing unit, e.g. the book element in a collection of books

2. Then, postprocess search results to find for each book the subelement that is the best hit.

This two-stage retrieval process often fails to return the best subelement because the relevance of a whole book is often not a good predictor of the relevance of small subelements within it.

# XML indexing unit: approach 3

Bottom up:

Instead of retrieving large units and identifying subelements (top down), we can search all leaves, select the most relevant ones and then extend them to larger units in postprocessing.

Similar problem as top down: the relevance of a leaf element is often not a good predictor of the relevance of elements it is contained in.

# XML indexing unit: approach 4

Index all elements: the least restrictive approach. Also problematic:

- Many XML elements are not meaningful search results, e.g., an ISBN number.

- Indexing all elements means that search results will be highly redundant.

## Example

For the query *Macbeth's castle* we would return all of the *play*, *act*, *scene* and *title* elements on the path between the root node and *Macbeth's castle*. The leaf node would then occur 4 times in the result set: 1 directly and 3 as part of other elements.

We call elements that are contained within each other **nested elements**. Returning redundant nested elements in a list of returned hits is not very user-friendly.

# Third challenge: nested elements

Because of the redundancy caused by the nested elements it is common to restrict the set of elements eligible for retrieval. Restriction strategies include:

- discard all small elements

- discard all element types that users do not look at (working XML retrieval system logs)

- discard all element types that assessors generally do not judge to be relevant (if relevance assessments are available)

- only keep element types that a system designer or librarian has deemed to be useful search results

In most of these approaches, result sets will still contain nested elements.

# Third challenge: nested elements

Further techniques:

- remove nested elements in a postprocessing step to reduce redundancy.

- collapse several nested elements in the results list and use highlighting of query terms to draw the user's attention to the relevant passages.

## Highlighting

- Gain 1: enables users to scan medium-sized elements (e.g., a section); thus, if the section and the paragraph both occur in the results list, it is sufficient to show the section.

- Gain 2: paragraphs are presented in-context (i.e., their embedding section). This context may be helpful in interpreting the paragraph.

# Nested elements and term statistics

Further challenge related to nesting: we may need to distinguish different contexts of a term when we compute term statistics for ranking, in particular inverse document frequency (idf ).

## Example

The term *Gates* under the node *author* is unrelated to an occurrence under a content node like *section* if used to refer to the plural of *gate*. It makes little sense to compute a single document frequency for *Gates* in this example.

Solution: compute idf for XML-context term pairs.

- sparse data problems (many XML-context pairs occur too rarely to reliably estimate df)

- compromise: consider the parent node x of the term and not the rest of the path from the root to x to distinguish contexts.

# Main idea: lexicalized subtrees

Aim: to have each dimension of the vector space encode a word together with its position within the XML tree.
How: Map XML documents to lexicalized subtrees.

# Main idea: lexicalized subtrees

1. Take each text node (leaf) and break it into multiple nodes, one for each word. E.g. split *Bill Gates* into *Bill* and *Gates*

2. Define the dimensions of the vector space to be lexicalized subtrees of documents – subtrees that contain at least one vocabulary term.

# Lexicalized subtrees

We can now represent queries and documents as vectors in this space of lexicalized subtrees and compute matches between them, e.g. using the vector space formalism.

**Vector space formalism in unstructured VS. structured IR**

The main difference is that the dimensions of vector space in unstructured retrieval are vocabulary terms whereas they are lexicalized subtrees in XML retrieval.

# Structural term

There is a tradeoff between the dimensionality of the space and the accuracy of query results.

▪ If we restrict dimensions to vocabulary terms, then we have a standard vector space retrieval system that will retrieve many documents that do not match the structure of the query (e.g., *Gates* in the title as opposed to the author element).

▪ If we create a separate dimension for each lexicalized subtree occurring in the collection, the dimensionality of the space becomes too large.

**Compromise:** index all paths that end in a single vocabulary term, in other words all XML-context term pairs. We call such an XML-context term pair a structural term and denote it by *<c, t>:* a pair of XML-context $c$ and vocabulary term $t$.

# Context resemblance

A simple measure of the similarity of a path $c_q$ in a query and a path $c_q$ in a document is the following *context resemblance* function CR:

$$\text{CR}(c_q, c_d) = \begin{cases} \frac{1+|c_q|}{1+|c_d|} & \text{if } c_q \text{ matches } c_d \\ 0 & \text{if } c_q \text{ does not match } c_d \end{cases}$$

$|c_q|$ and $|c_d|$ are the number of nodes in the query path and document path, resp.

$c_q$ matches $c_d$ iff we can transform $c_q$ into $c_d$ by inserting additional nodes.

# Context resemblance example



$$\mathrm{CR}(c_q, c_d) = \begin{cases} \dfrac{1+|c_q|}{1+|c_d|} & \text{if } c_q \text{ matches } c_d \\ 0 & \text{if } c_q \text{ does not match } c_d \end{cases}$$

$\mathrm{CR}(c_q, c_d) = 3/4 = 0.75$. The value of $\mathrm{CR}(c_q, c_d)$ is 1.0 if $q$ and $d$ are identical.

# Context resemblance example



$$\mathrm{CR}(c_q, c_d) = \begin{cases} \frac{1+|c_q|}{1+|c_d|} & \text{if } c_q \text{ matches } c_d \\ 0 & \text{if } c_q \text{ does not match } c_d \end{cases}$$

$\mathrm{CR}(c_q, c_d) = ?$ $\mathrm{CR}(c_q, c_d) = 3/5 = 0.6.$

# Document similarity measure

The final score for a document is computed as a variant of the cosine measure, which we call SIMNOMERGE.

SIMNOMERGE($q$, $d$) =

$$\sum_{c_k \in B} \sum_{c_l \in B} \mathrm{CR}(c_k, c_l) \sum_{t \in V} \mathrm{weight}(q, t, c_k) \frac{\mathrm{weight}(d, t, c_l)}{\sqrt{\sum_{c \in B, t \in V} \mathrm{weight}^2(d, t, c)}}$$

- *V* is the vocabulary of non-structural terms

- *B* is the set of all XML contexts

- weight ($q$, $t$, $c$), weight($d$, $t$, $c$) are the weights of term t in XML context $c$ in query $q$ and document $d$, resp. (standard weighting e.g. $\mathrm{idf}_t$ x $\mathrm{wf}_{t,d}$, where $\mathrm{idf}_t$ depends on which elements we use to compute $\mathrm{df}_t$.)

SIMNOMERGE($q$, $d$) is not a true cosine measure since its value can be larger than 1.0.

# SIMNOMERGE algorithm

SCOREDOCUMENTSWITHSIMNOMERGE(*q, B, V, N, normalizer*)

```
1    for n ← 1 to N
2    do score[n] ← 0
3    for  each ⟨cq, t⟩ ∈ q
4    do wq ← WEIGHT(q, t, cq)
5        for  each c ∈ B
6        do if CR(cq, c) > 0
7            then postings ← GETPOSTINGS(⟨c, t⟩)
8                for  each posting ∈ postings
9                do x ← CR(cq, c) * wq * weight(posting)
10                   score[docID(posting)]+ = x
11   for n ← 1 to N
12   do score[n] ← score[n]/normalizer[n]
13   return score
```

# Initiative for the Evaluation of XML retrieval (INEX)

INEX: standard benchmark evaluation (yearly) that has produced test collections (documents, sets of queries, and relevance judgments).
Based on IEEE journal collection (since 2006 INEX uses the much larger English Wikipedia test collection).
The relevance of documents is judged by human assessors.

| INEX 2002 collection statistics | |
|---|---|
| 12,107 | number of documents |
| 494 MB | size |
| 1995—2002 | time of publication of articles |
| 1,532 | average number of XML nodes per document |
| 6.9 | average depth of a node |
| 30 | number of CAS topics |
| 30 | number of CO topics |

# INEX topics

Two types:

1. content-only or **CO topics**: regular keyword queries as in unstructured information retrieval

2. content-and-structure or **CAS topics**: have structural constraints in addition to keywords

Since CAS queries have both structural and content criteria, relevance assessments are more complicated than in unstructured retrieval

# INEX relevance assessments

INEX 2002 defined component coverage and topical relevance as orthogonal dimensions of relevance.

## Component coverage

Evaluates whether the element retrieved is "structurally" correct, i.e., neither too low nor too high in the tree.

We distinguish four cases:

1. Exact coverage (E): The information sought is the main topic of the component and the component is a meaningful unit of information.

2. Too small (S): The information sought is the main topic of the component, but the component is not a meaningful (self-contained) unit of information.

3. Too large (L): The information sought is present in the component, but is not the main topic.

4. No coverage (N): The information sought is not a topic of the component.

# INEX relevance assessments

The **topical relevance** dimension also has four levels: highly relevant (3), fairly relevant (2), marginally relevant (1) and nonrelevant (0).

## Combining the relevance dimensions

Components are judged on both dimensions and the judgments are then combined into a digit-letter code, e.g. 2S is a fairly relevant component that is too small. In theory, there are 16 combinations of coverage and relevance, but many cannot occur. For example, a nonrelevant component cannot have exact coverage, so the combination 3N is not possible.

# INEX relevance assessments

The relevance-coverage combinations are quantized as follows:

$$\mathbf{Q}(\textit{rel}, \textit{cov}) = \begin{cases} 1.00 & \text{if} \quad (\textit{rel}, \textit{cov}) = 3\text{E} \\ 0.75 & \text{if} \quad (\textit{rel}, \textit{cov}) \in \{2\text{E}, 3\text{L}\} \\ 0.50 & \text{if} \quad (\textit{rel}, \textit{cov}) \in \{1\text{E}, 2\text{L}, 2\text{S}\} \\ 0.25 & \text{if} \quad (\textit{rel}, \textit{cov}) \in \{1\text{S}, 1\text{L}\} \\ 0.00 & \text{if} \quad (\textit{rel}, \textit{cov}) = 0\text{N} \end{cases}$$

This evaluation scheme takes account of the fact that binary relevance judgments, which are standard in unstructured IR, are not appropriate for XML retrieval. The quantization function **Q** does not impose a binary choice relevant/nonrelevant and instead allows us to grade the component as partially relevant. The number of relevant components in a retrieved set *A* of components can then be computed as:

$$\#(\text{relevant items retrieved}) = \sum_{c \in A} \mathbf{Q}(\textit{rel}(c), \textit{cov}(c))$$

# INEX evaluation measures

As an approximation, the standard definitions of precision and recall can be applied to this modified definition of relevant items retrieved, with some subtleties because we sum graded as opposed to binary relevance assessments.

**Drawback**

Overlap is not accounted for. Accentuated by the problem of multiple nested elements occurring in a search result.

Recent INEX focus: develop algorithms and evaluation measures that return non-redundant results lists and evaluate them properly.

# Recap

- Structured or XML IR: effort to port unstructured (standard) IR know-how onto a scenario that uses structured (DB-like) data

- Specialized applications (e.g. patents, digital libraries)

- A decade old, unsolved problem

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

# A Data Mashup Language for the Data Web

Mustafa Jarrar, Marios D. Dikaiakos

University of Cyprus

LDOW 2009, April 20, 2009

Edited & Presentation by Sangkeun Lee, IDS Lab

Original Slides : http://www.cs.ucy.ac.cy/~mjarrar/Internal/MashQL.V07.ppt

# Outline

- Introduction & Motivation

- The MashQL Language

- The Notion of Query Pipes

- Implementation

- Use cases

- Discussion and Future Directions

# Introduction & Motivation

- ## We are witnessing

  – A rapid emergence of the Data Web

  – Many companies started to make their content freely accessible through APIs

    - E.g. Google Base, eBay, Flickr, eBay

  – Many accessible data in RDF, RDFa

# Web 2.0 and the phenomena of APIs

# Web 2.0 and the phenomena of APIs



Wikipedia in RDF

API

# Web 2.0 and the phenomena of APIs

# Web 2.0 and the phenomena of APIs

# Web 2.0 and the phenomena of APIs



Also supports microformats/RDFa

API

# Web 2.0 and the phenomena of APIs

# Web 2.0 and the phenomena of APIs

# Web 2.0 and the phenomena of APIs



**And many, many others APIs**

API

# Web 2.0 and the phenomena of APIs

Moving to the **Data Web**, in parallel to the web of documents.

# Introduction & Motivation

- ## A Mashup?

  - A Web application that consumes data originated from third parties and retrieved via APIs

  - Problem

    - Building mashups is an art that is limited to skilled programmers

    - Some mashup editors have been proposed by Web 2.0 communities, but…?



(A puzzle of APIs)

$(API_1 + API_2) + API_3 = money$

# How to Build a Mashup?

What do you want to do?

Which data you need? APIs/RSS available?
How is your programming skills?

Geek

Semi-Technical Skills

Sign up for a developer token

**http://aws.amazon.com/**
**http://www.google.com/apis/maps/**
**http://api.search.yahoo.com/webservices/re**

Use mashup editors

Microsoft Popfly
Yahoo! Pipes
QEDWiki by IBM
Google Mashup Editor (Coming)
Serena Business Mashups
Dapper
JackBe Presto Wires

Start coding

Start Configuring

# Mashup Editors

# Mashup Editors

# Mashup Editors

# Mashup Editors

# Mashup Editors

# Limitations of Mashup Editors

- Focus only on providing encapsulated access to (some) public APIs and feeds (rather than querying data sources).

- Still require programming skills.

- Cannot play the role of a *general-purpose data retrieval*, as mashups are sophisticated applications.

- Lacks a formal framework for pipelining mashups.

# Vision

- Position
  - The author propose to regard the web as a database
  - Mashup is seen as a query over one or multiple sources

- So, instead of developing a mashup as an application that access structured data through APIs,

- We regard *a mashup* as *a query*

- *Challenges*
  - ☐ But the problem then is: users need to know the schema and technical details of the data sources they want to query.

# Vision and Challenges

How a user can query a source without knowing its schema, structure, and vocabulary?

**DateSources**

```
SELECT S.Title FROM Google.Scholar S
Where (S.Author='Hacker')
Union
SELECT P.PattentTitle FROM Ggoogle.Patent P
Where (P.Inventor ='Hacker')
Union
SELECT A.Title FROM Citeseer A
Where (P.Author ='Hacker')
```

# Vision and Challenges

How a user can query a source without knowing
its schema, structure, and vocabulary?



```
SELECT S.Title FROM Google.Scholar S
Where (S.Author='Hacker')
Union
SELECT P.PattentTitle FROM Ggoogle.Patent P
Where (P.Inventor ='Hacker')
Union
SELECT A.Title FROM Citeseer A
Where (P.Author ='Hacker')
```
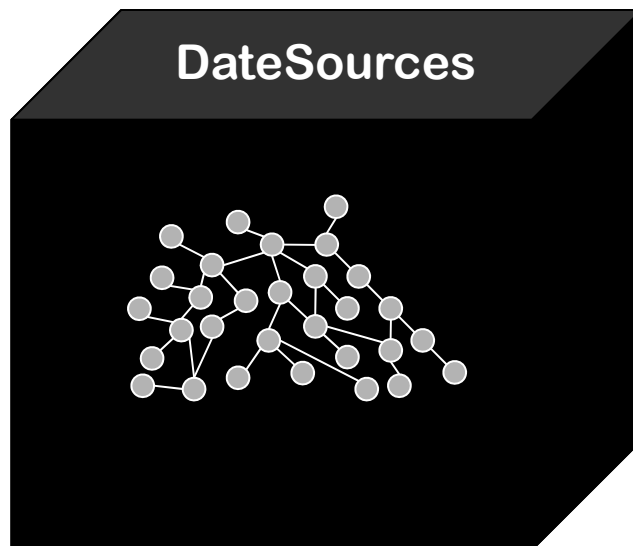
# MashQL

- A simple query language for the Data Web, in a mashup style.

- MashQL allows <span style="color:red">querying a dataspace(s) without any prior knowledge about its schema, vocabulary or technical details</span> (a source may not have a schema al all).   Explore unknown graph

- <span style="color:red">Does not assume any knowledge about RDF, SPARQL, XML, or any technology</span>, to get started.

- Users only use drop-lists to formulate queries (query-by-diagram/interaction).

# MashQL Example 1

`http:www.site1.com/rdf`

```
<:a1>   <:Title>       "Web 2.0"
<:a1>   <:Author>      "Hacker B."
<:a1>   <:Year>        2007
<:a1>   <:Publisher>   "Springer"
<:a2>   <:Title>       "Web 3.0"
<:a2>   <:Author>      "Smith B."
<:a2>   <:Cites>       <:a1>
```

`http:www.site2.com/rdf`

```
<:4>   <:Title>     "Semantic Web"
<:4>   <:Author>    "Tom Lara"
<:4>   <:PubYear>   2005
<:5>   <:Title>     "Web Services"
<:5>   <:Author>    "Bob Hacker"
```

## Hacker's Articles after 2000?

**RDF Input** [?][−][✕]

From:
- http://www.site1.com/rdf
- http://www.site2.com/rdf

**MashQL** [?][−][✕]

○ *Everything*

　　○ Title ☑*ArticleTitle*

　　○ Author "^Hacker"

　　○ Year\PubYear > 2000

# MashQL Example 1

http:www.site1.com/rdf
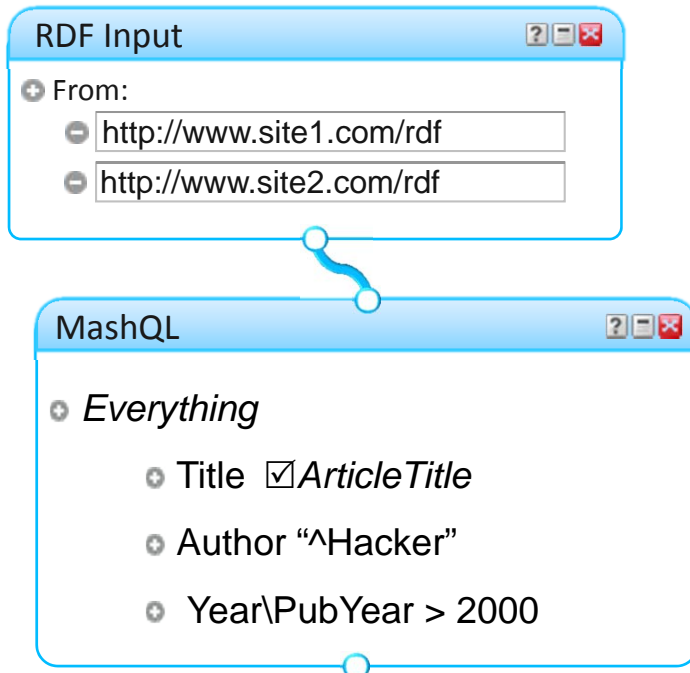
```
<:a1>  <:Title>      "Web 2.0"
<:a1>  <:Author>     "Hacker B."
<:a1>  <:Year>       2007
<:a1>  <:Publisher>  "Springer"
<:a2>  <:Title>      "Web 3.0"
<:a2>  <:Author>     "Smith B."
<:a2>  <:Cites>      <:a1>
```

http:www.site2.com/rdf

```
<:4>  <:Title>    "Semantic Web"
<:4>  <:Author>   "Tom Lara"
<:4>  <:PubYear>  2005
<:5>  <:Title>    "Web Services"
<:5>  <:Author>   "Bob Hacker"
```

Hacker's Articles after 2000?

### RDF Input  ? ▢ ☒

⊕ From:
   ⊖ http://www.site1.com/rdf
   ⊖ http://www.site2.com/rdf

### MashQL  ? ▢ ☒

⊕ Everything ▾
   a1
   a2
   4
   5
   ☐ Types  ☑ Instances

Interactive
query formulation

# MashQL Example 1

`http:www.site1.com/rdf`

```
<:a1>   <:Title>       "Web 2.0"
<:a1>   <:Author>      "Hacker B."
<:a1>   <:Year>        2007
<:a1>   <:Publisher>   "Springer"
<:a2>   <:Title>       "Web 3.0"
<:a2>   <:Author>      "Smith B."
<:a2>   <:Cites>       <:a1>
```

`http:www.site2.com/rdf`

```
<:4>   <:Title>     "Semantic Web"
<:4>   <:Author>    "Tom Lara"
<:4>   <:PubYear>   2005
<:5>   <:Title>     "Web Services"
<:5>   <:Author>    "Bob Hacker"
```

## Hacker's Articles after 2000?

**RDF Input**  ? ▬ ✖

⊕ From:
  ⊖ http://www.site1.com/rdf
  ⊖ http://www.site2.com/rdf

**MashQL**  ? ▬ ✖

⊕ *Everything*
  ⊕ Title ▼      ▼ ☑ ArticleTitle ▼
    Author
    Cites
    Publisher
    PubYear
    Title
    Year

# MashQL Example 1

http:www.site1.com/rdf

```
<:a1>   <:Title>        "Web 2.0"
<:a1>   <:Author>       "Hacker B."
<:a1>   <:Year>         2007
<:a1>   <:Publisher>    "Springer"
<:a2>   <:Title>        "Web 3.0"
<:a2>   <:Author>       "Smith B."
<:a2>   <:Cites>        <:a1>
```
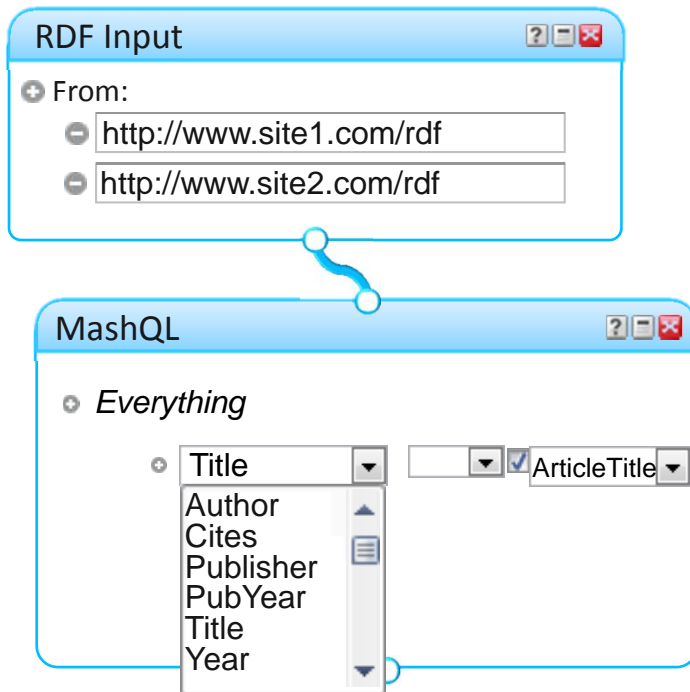
http:www.site2.com/rdf

```
<:4>   <:Title>     "Semantic Web"
<:4>   <:Author>    "Tom Lara"
<:4>   <:PubYear>   2005
<:5>   <:Title>     "Web Services"
<:5>   <:Author>    "Bob Hacker"
```

## Hacker's Articles after 2000?

**RDF Input**   ? ☐ ✖

⊕ From:
  ⊖ http://www.site1.com/rdf
  ⊖ http://www.site2.com/rdf

**MashQL**   ? ☐ ✖

⊕ *Everything*

  ⊕ Title   ☑*Article title*

  ⊕ | Author ▼ |   | Con ▼ | Hacker ▼ |

```
Author    ▲        Equals      ▲
Cites     ▤        Contains    ▤
Publisher          OneOf
PubYear            Not
Title              Between
Year      ▼        LessThan
                   MoreThan    ▼
```

# MashQL Example 1

`http:www.site1.com/rdf`

```
<:a1>   <:Title>      "Web 2.0"
<:a1>   <:Author>     "Hacker B."
<:a1>   <:Year>       2007
<:a1>   <:Publisher>  "Springer"
<:a2>   <:Title>      "Web 3.0"
<:a2>   <:Author>     "Smith B."
<:a2>   <:Cites>      <:a1>
```

`http:www.site2.com/rdf`

```
<:4>   <:Title>     "Semantic Web"
<:4>   <:Author>    "Tom Lara"
<:4>   <:PubYear>   2005
<:5>   <:Title>     "Web Services"
<:5>   <:Author>    "Bob Hacker"
```

Hacker's Articles after 2000?

**RDF Input** ? _ ✖

◍ From:
⊖ http://www.site1.com/rdf
⊖ http://www.site2.com/rdf

**MashQL** ? _ ✖

◍ *Everything*

  ◍ Title  ☑*Article title*

  ◍ *Author "^Hacker"*

  ◍ [Year \PubYe ▾] [mor ▾] [2000 ▾]

  | Publisher ▲ | | OneOf ▲ |
  | PubYear ▤ | | Not ▤ |
  | Title | | Between |
  | Year | | LessThan |
  | | | MoreThan |

# MashQL Example 1

http:www.site1.com/rdf

```
<:a1>   <:Title>       "Web 2.0"
<:a1>   <:Author>      "Hacker B."
<:a1>   <:Year>        2007
<:a1>   <:Publisher>   "Springer"
<:a2>   <:Title>       "Web 3.0"
<:a2>   <:Author>      "Smith B."
<:a2>   <:Cites>       <:a1>
```

http:www.site2.com/rdf

```
<:4>   <:Title>     "Semantic Web"
<:4>   <:Author>    "Tom Lara"
<:4>   <:PubYear>   2005
<:5>   <:Title>     "Web Services"
<:5>   <:Author>    "Bob Hacker"
```

## Hacker's Articles after 2000?

**RDF Input**  ? ▬ ✖

⊕ From:
  ⊖ http://www.site1.com/rdf
  ⊖ http://www.site2.com/rdf

**MashQL**  ? ▬ ✖

⊕ *Everything*

  ⊕ Title  ☑ *Article title*

  ⊕ Author *"^Hacker"*

  ⊕ Year/PubYear > 2000

```
PREFIX S1: <http://site1.com/rdf>
PREFIX S2: <http://site1.com/rdf>
SELECT ? ArticleTitle
FROM <http://site1.com/rdf>
FROM <http://site2.com/rdf>
WHERE {
 {{?X S1:Title ?ArticleTitle}UNION
 {?X S2:Title ?ArticleTitle}}
 {?X S1:Author ?X1} UNION {?X S2:Author ?X1}
 {?X S1:PubYear ?X2} UNION {?X S2:Year ?X2}
 FILTER regex(?X1, "^Hacker")

 FILTER (?X2 > 2000)}
```
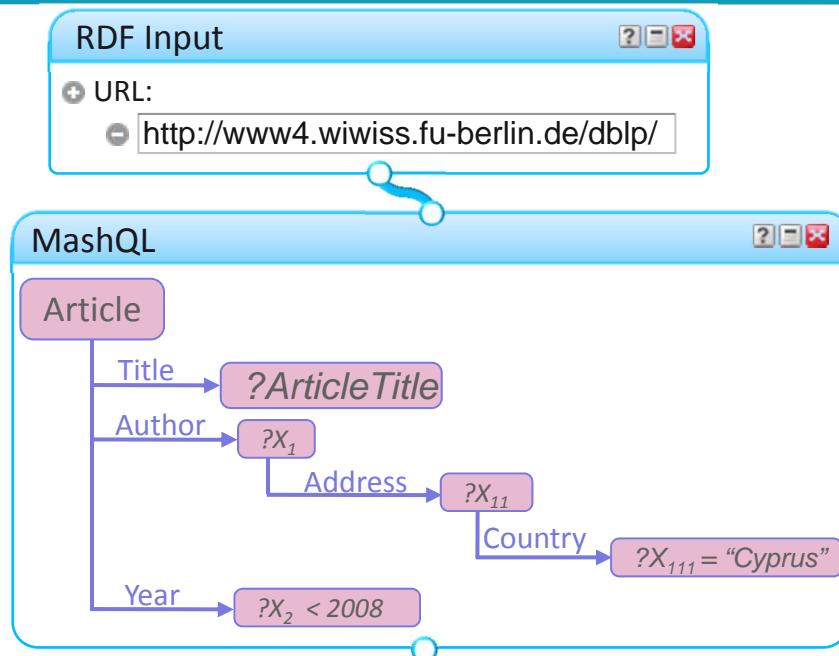
# MashQL Example 2

**RDF Input**  ?  —  ✖

⊕ URL:
   ⊖ http://www4.wiwiss.fu-berlin.de/dblp/

**MashQL**  ?  —  ✖

*Article*
  Title  ☑ *Article Title*
  Author
     Address
       Country "Cyprus"
  Year > 2008

The recent articles from Cyprus

→ Retrieve every Article that has a title, written by an author, who has an address, this address has a country called Cyprus, and the article published after 2008.

# The Intuition of MashQL



**A query is a tree**

- The root is called the query *subject.*

- Each branch is a *restriction*.

- Branches can be expanded, (information path)

- Object value filters

*Def.  A **Query** Q with a subject S, denoted by Q(S), is a set of restrictions on S. Q(S) = $R_1$ AND … AND $R_n$.*

*Dif.  A **Subject** S $\in$ (I $\cup$ V), where I is an identifier and V is a variable.*

*Dif. A **Restriction** R = <$R_x$ , P, $O_f$>, where $R_x$ is an optional restriction prefix that can be (maybe | without), P is a predicate (P $\in$ I $\cup$ V), and $O_f$ is an object filter.*

# The Intuition of MashQL

**RDF Input**  ? □ ✖

⊕ URL:
⊖ http://www4.wiwiss.fu-berlin.de/dblp/
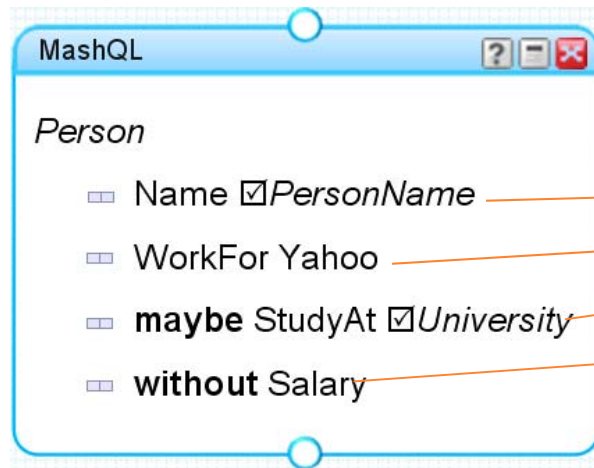
**MashQL**  ? □ ✖

*Article*

▭ Title  ☑ *ArticleTitle*

▭ Author

▭ Address

▭ Country "Cyprus"

▭ Year > 2008

**An Object filter** *is one of :*
- *Equals*
- *Contains*
- *MoreThan*
- *LessThan*
- *Between*
- *one of*
- *Not(f)*
- *Information Path (sub query)*

# More MashQL Constructs

➢ Resection Operators {Required, Maybe, or Without}

  All restriction are required (i.e. AND), unless they are prefixed with "maybe" or "without"



```
SELECT ?PersonName, ?University
WHERE {
    ?Person :Name ?PersonName.
    ?Person :WorkFor :Yahoo.
    OPTIONAL{?Person :StudyAt ?University}
    OPTIONAL{?Person :Salary ?X1}
    FILTER (!Bound(?X1))}
}
```

# More MashQL Constructs

➢ Union operator (denoted as "\") between Objects, Predicates, Subjects and Queries



```
SELECT ?Person
WHERE {
    ?Person :WorkFor :Google
    UNION
    ?Person WorkFor :Yahoo}
```

```
SELECT ?FName
WHERE {
    ?Person :Surname ?FName
    UNION
    ?Person :Firstname ?FName}
```

```
SELECT ?AgentName, ?AgentPhone
WHERE {
    {?Person rdf:type :Person.
    ?Person :Name ?AgentName.
    ?Person :Phone ?AgentPhone}
UNION
    {?Company rdf:type :Company.
    ?Company :Name ?AgentName.
    ?Company :Phone ?AgentPhone}}
```

# MashQL Queries

- In the background, MashQL queries are translated into and executed as SPARQL queries.

- At the moment, we focus on RDF (/RDFa) as a data format, and SPARQL (/Oracle's SPARQL) as a backend query language. However, MashQL can be easily mappable to other query languages.

# MashQL Compilation

Depending on the pipeline structure, MashQL generates either SELECT or CONSTRUCT queries:

- SELECT returns the results in a tabular form
  (e.g. ArticleTitle, Author)

- CONSTRUCT returns the results in a triple form
  (e.g. Subject, Predicate, Object).

```
…
CONSTRUCT *
WHERE{?Job :JobIndustry ?X1.
  ?Job :Type ?X2.
  ?Job :Currency ?X3.
  ?Job :Salary ?X4.
  FILTER(?X1="Education"|| ?X1="HealthCare")
  FILTER(?X2="Full-Time"|| ?X2="Fulltime")||
                           ?X2="Contract")
  FILTER(?X3="^Euro"|| ?X3="^€")
  FILTER(?X4>=75000||  ?X4<=120000)}
```

```
…
SELECT ?Job ?Firm
WHERE
   {?Job :Location ?X1. ?X1 :Country ?X2.
   FILTER (?X2="Italy"||?X2="Spain")||
            ?X2="Greece"||?X2="Cyprus")}
   OPTIONAL{{?job :Organization ?Firm} UNION
{?job :Employer ?Firm}}
```

# MashQL Editor
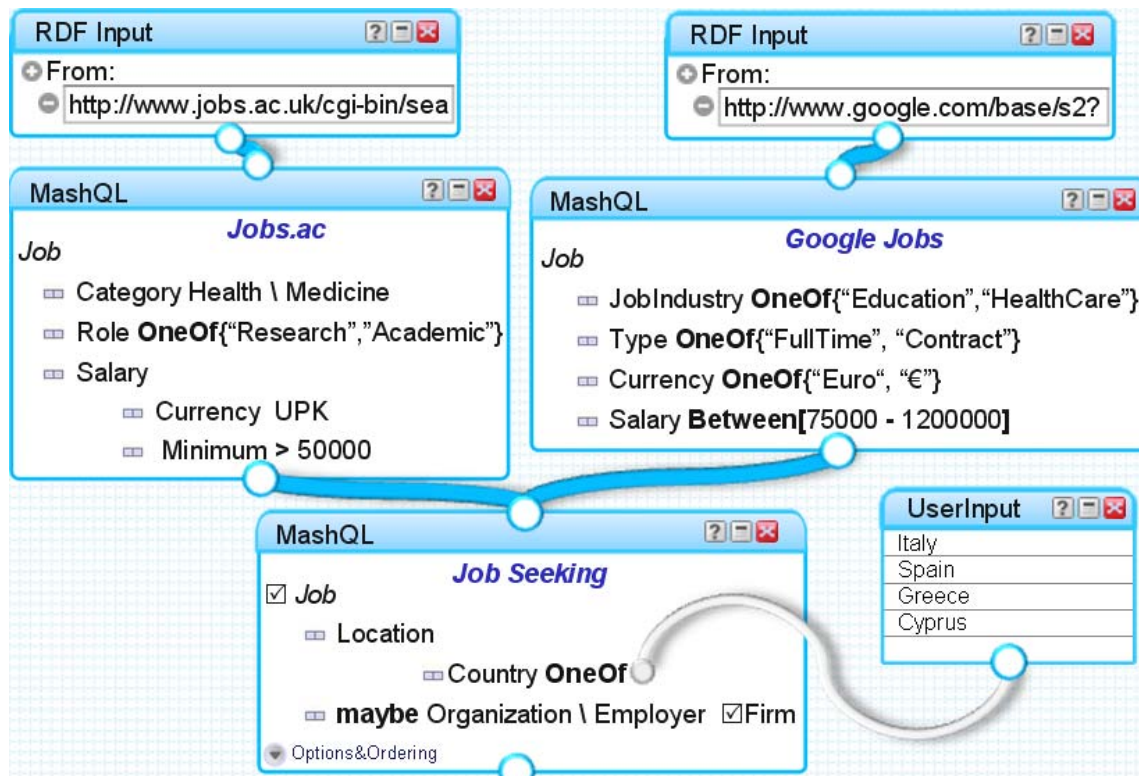
# MashQL Firefox Add-On (Light-mashups @ your browser)

# Use Case: Job Seeking

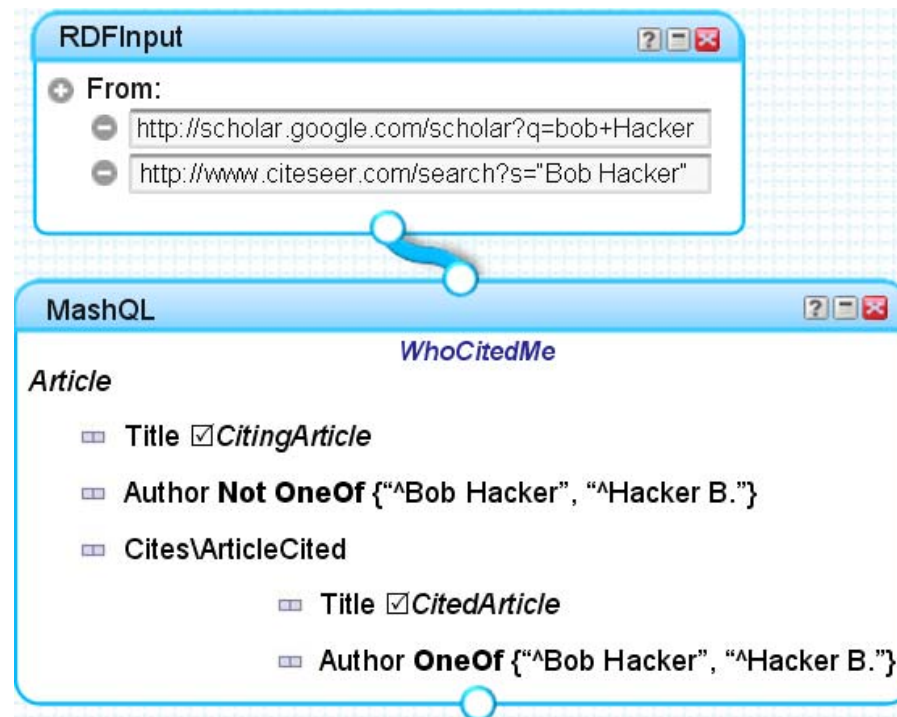A mashup of job vacancies based on Google Base and on Jobs.ac.uk.

# Use Case: My Citations

A mashup of cited Hacker's articles (but no self citations), over Scholar and Siteseer

# Evaluation

Query Execution :

- The performance of executing a MashQL query is bounded to the performance to executing its backend language (i.e. SPARQL/SQL).

- A query with medium size complexity takes one or few seconds (Oracle's SPARQL, [Chong et al 2007]).

# Conclusions

- A formal but yet simple query language for the Data Web, in a mashup and declarative style.

- Allows people to discover and navigate unknown data spaces(/graphs) without prior knowledge about the schema or technical details.

- Can be use as a general purpose data retrieval and filtering