

Database and Information Systems

Homework 4

Due on November 7, 2007

For this homework, you should test your answers using SAXON, an XSLT and XQuery processor. If you would like to work from your own machine, you may download SAXON and its documentation from <http://saxon.sourceforge.net/>. Alternatively, you can ssh to *eniac.seas.upenn.edu* and use the resources in *~zives/saxon/* to test your query source file. You might need to add some lines to your .profile file first:

```
export JAVA_HOME=/usr/java/jdk1.6.0_02
export PATH=${PATH}: ${JAVA_HOME}/bin
export CLASSPATH=${CLASSPATH}: ~zives/saxon/saxon8.jar
```

Source your .profile. (This assumes your shell is bash. If not, simply type “bash” and then “source profile.”) There’s an example XQuery file “*example.xq*,” which uses the XML data file “*book.xml*”. Type in the command:

```
java net.sf.saxon.Query -t -s book.xml example.xq
```

Problem 1 [70 points]: The XML data files used for this problem are **proc.xml** (list of conference proceedings) and **inproc.xml** (list of articles in proceedings) in *~zives/saxon/* on eniac. Based on the two XML files, write the following queries in XQuery, with the query output embedded within tags `<answer> ... </answer>`. To test, you may type in the command:

```
java net.sf.saxon.Query -t -s proc.xml -s inproc.xml
/yourdirectory/hw4_1_1.xq    >/yourdirectory/hw4_1_1.xml
```

And view your results in the *hw4_1_1.xml* file in your directory.

1. For each paper authored by Jeffrey D. Ullman, return a **paper** tag, with the following contents: the paper title (in the appropriate element) and Ullman’s **coauthors**. Note that Ullman is not his own coauthor (coauthorship is not reflexive).
2. Find the **authors** whose papers were most cited in the 1990s. Note that you can test the crossref to the proceedings key attribute (which ends with the year) to determine the publication year. The result should be nested in an **author** tag.

3. Find the **authors** who published the most papers in the proceedings of SIGMOD (any year). You may wish to use the XQuery function fn:contains() to match substrings. (See <http://saxon.sourceforge.net/saxon7.6/functions.html>.)
4. Return, for each paper in the proceedings of VLDB 2003, the **title** and the *second author*. If the paper only has one author, omit the **author** tag entirely (but still return the title). The result should be nested within a **paper** tag.
5. Find the distinct **titles** of proceedings in **proc** that don't have any paper information in **inproc**. The result should be nested in **title** tag.
6. Write a query that *maps* information for the year 2002 from **inproc.xml** and **proc.xml** into a schema that follows the pattern:

```

<publications>
  <author> *
    <name> .. </name>
    <in> *
      <forum> .. </forum>
      <article> .. </article>
    </in>
  </author>
</publications>

```

Again, you may wish to use fn:contains() or fn:ends-with().

7. Sort papers in **inproc** in ascending alphabetical order of first author. When two papers were published by the same first author, further sort them in ascending alphabetical order of second author, etc. Results should be in the following format:

```

<answer>
  <paper>
    sorted papers here
  </paper>
</answer>

```

Since your XQuery handout does not cover the “order by” clause of XQuery (instead it uses the older sortby), see http://www.w3schools.com/xquery/xquery_select.asp for an example of sorting.

Problem 2 [30 points]: Design an XML schema for the PBAY system.

Sellers(*sellerID*: int, *rating*: char, *email*: string)
Items(*itemID*: int, *type*: string)
Buyers(*buyerID*: int, *email*: string, *city*: string, *state*: string)
Stock(*itemID*: int, *sellerID*: int, *startBid*: float, *quantity*: int, *endTime*: int)
Purchases(*itemID*: int, *buyerID*: int, *sellerID*: int, *price*: float, *purchaseQuantity*: int, *bidTime*: int)

The entity sets are **Sellers**, **Items** and **Buyers**. Make up some appropriate attributes for each entity set. Each of the entity sets should have 7 attributes. Specify keys and keyrefs for each complexType. Explain how to encode relationship sets like **Stock** and **Purchases** in the hierarchy of the XML schema.