

# Rule-based Query Tree Evaluation Over Fragmented XML Documents

Ron G. McFadyen<sup>+</sup>, Yangjun Chen<sup>\*</sup>

Department of Applied Computer Science, University of Winnipeg

515 Portage Avenue, Winnipeg, Manitoba, Canada, R3B 2E9

{r.mcfadyen, ychen2}@uwinnipeg.ca

<sup>+</sup>Supported by NSERC 105709-03 (139988) (Natural Sciences and Engineering Council of Canada)

<sup>\*</sup>Supported by NSERC 239074-05 (242523) (Natural Sciences and Engineering Council of Canada)

## Abstract

XML documents are used to hold information and to make exchanges between systems. In this paper, we consider documents that embed knowledge and rules, which may contain considerable redundancy. To control redundancy, documents are decomposed into fragments that are stored separately. Then, to materialize documents for end-users, they need to be dynamically constructed from their sources (separately stored fragments) by evaluating rules, which requires database queries to be executed according to the document structure.

Keywords: XML, business rule, query tree, query evaluation, document fragmentation and assembly.

## 1 Introduction

XML and related technologies are used extensively for storing, managing, and exchanging information. Our work is motivated by the consideration of documents describing requirements or business rules to be met to achieve some designation or status. According to [1, 2], a business rule is a statement that defines or constrains some aspect of the business and is used to control the behavior of the business. For example, in a university calendar, we have rules that specify the valid collection of courses for a student to obtain a certain university degree.

As shown in [3], such rules can be organized into a special structure, the so-called *synthesized query tree* (SQT), to govern their execution. An observation shows that such trees may heavily overlap. Then, the separate storage of SQTs not only leads to redundant space, but also causes repeated evaluation of rules. For this reason, we fragment individual SQTs and organize the sub-SQTs into a hyper tree structure, in which a leaf node can be a simple query or a subtree itself. This kind of organization is similar to document fragmentation [4, 5]. However, in our case, a fragment is a subtree representing a set of rules, which produces a piece of a document dynamically. Our method also shares the flavour of active XML documents [6, 7, 8], by which parts of the contents are generated by invoking a

program or a web service; but differs from these in that we are concerning with the processing of business rules, which are evaluated along a tree structure in a bottom-up way.

The remainder of this paper is organized as follows. Next in section 2, we present the background information on requirement documents and synthesized query trees. In section 3, we present the SQT fragmentation based on the so-call *virtual SQTs*, which enable us to efficiently evaluate queries. Section 4 presents a short conclusion and directions for further work.

## 2 Requirements Document and SQT Structure

In Figure 1, we show some typical majors found in a university calendar. Their XML representation is illustrated in Figure 2. In Figure 1, each document is a requirement for graduation from a certain major.

For instance, in Figure 1(a), we specify that to graduate with a 3-Year BSc in Geography a student must satisfy all of: *Graduation Requirement* (completion of 90 credit hours), *Residence Requirement* (completion of a minimum of 30 credit hours at the university and a minimum of 18 credit hours in Geography at the university), a *General Requirement* (12 credit hours in Humanities subjects and 6 credit hours in Science subjects), and a *Major Requirement* for certain courses in Geography. The fact that all four of these requirements must be met simultaneously is indicated by the attribute: *combining="AND"*, associated with the element GeographyRule in the document shown in Figure 2.

As discussed in [3], we can construct a synthesized query tree over such a document. [3] presents two such tree structures: the boolean and the general SQT; for simplicity, here we only discuss the boolean SQT.

**Definition 1:** a *boolean synthesized query tree* (BSQT) is a tree where each leaf node  $v$  is associated with a boolean query  $Q(v)$ , and each internal node  $v$  is labelled with a tag

$T(v)$ , and an operator  $\theta = \vee$  or  $\wedge$ ; and each node  $v$  is assigned a boolean value,  $V(v)$ , determined as follows:

- a) for a leaf node,  $V(v)$  is *true* if the return value of  $Q(v)$  is not empty; otherwise, it is *false*, and
- b) for an internal node, with children  $v_1, \dots, v_n$ ,

$$V(v) = V(v_1)\theta V(v_2)\dots\theta V(v_n).$$

For instance, for the graduation requirement of Geography, we will construct a tree structure as shown in Figure 3, which represents the rules, queries and relationships corresponding to the requirements shown in Figure 2. One of the defining characteristics of a query tree is that database queries are only present in leaf nodes.

### 3-Year BSc (Geography)

- Graduation Requirement
  - 90 credit hours
- Residence Requirement
  - Degree: minimum 30 credit hours
  - Major: minimum 18 credit hours
- General Degree Requirement
  - Humanities: 12 credit hours
  - Science: 6 credit hours
- Major Requirement
  - Minimum 30 credit hours
  - Maximum 48 credit hours
  - Required Courses
    - 23.202 Intro Geography I
    - 23.203 Intro Geography II
    - 23.331 Advanced Geography
  - Choice
    - 23.205 Atmos Sci or 23.206 Earth Sci

#### a) Graduation Requirements - Geography

### 3-Year BSc (Physics)

- Graduation Requirement
  - 90 credit hours
- Residence Requirement
  - Degree: minimum 30 credit hours
  - Major: minimum 18 credit hours
- General Degree Requirement
  - Humanities: 12 credit hours
  - Science: 6 credit hours
- Major Requirement
  - Minimum 36 credit hours
  - Maximum 54 credit hours
  - Required Courses
    - 44.101 Intro Physics
    - 44.203 Mechanics
    - 44.331 Relativity

#### b) Graduation Requirements - Physics

**Figure 1. Graduation Requirements**

To determine if a student can graduate, it is necessary to evaluate the appropriate SQT and its queries in the context of the student. Thus, to do this for a number of students, we have to traverse the SQT trees repeatedly, each time for a single student.

An observation of the queries present in the leaves shows that a slight modification will facilitate set processing. To

see this, let's have a look at the query  $Q_1$  shown in Figure 3. If we remove the condition  $studentNum = x$  from the where-clause, the execution of the query will find all the students with grade point  $\geq 1$  and more than 90 credit hours.

```
<GeographyRule title="Degree Requirement for 3-Year BSc (Geography)",
  combining="AND">
  <GraduationRule title="Graduation Requirement",
    display="90 credit hours", query="..." , ...>
  </GraduationRule>
  <ResidenceRule title="Residence Requirement",
    combining="AND">
    <DegreeRule title="Degree",
      display="minimum 30 credit hours",
      query="SELECT sum(creditHours)
        FROM studentHistory
        WHERE
          studentNumber=parameterValue",
      expected="30",
      comparison=">=">
    </DegreeRule>
    <MajorRule title="Major", display="minimum 18 credit hours",
      query="..." , ...>
    </MajorRule>
  </ResidenceRule>
  <GeneralRule title="General Degree Requirement,
    combining="AND">
    <HumanitiesRule title="Humanities",
      display="12 credit hours", query="..." , ...>
    </HumanitiesRule>
    <ScienceRule title="Science", display="6 credit hours", query="..." , ...>
    </ScienceRule>
  </GeneralRule>
  <MajorRule title="Major Requirement"
    combining="AND">
    <MinMaxRule
      display="Minimum 30 credit hours, Maximum 48 credit hours",
      query="..." , ...>
    </MinMaxRule>
    <ReqCoursesRule title="Required Courses",
      combining="AND">
    <Course
      display="23.02 Intro Geography I", query="..." , ...> </Course>
    <Course
      display="23.203 Intro Geography II", query="..." , ...> </Course>
    <Course
      display="23.331 Advanced Geography", query="..." , ...> </Course>
    </ReqCoursesRule>
    <ChoiceRule title="Choice",
      display="23.205 Atmos Sci or 23.206 Earth Sci"
      combining="OR">
    <Course query="..." , ...></Course>
    <Course query="..." , ...></Course>
    </ChoiceRule>
  </MajorRule>
</GeographyRule>
```

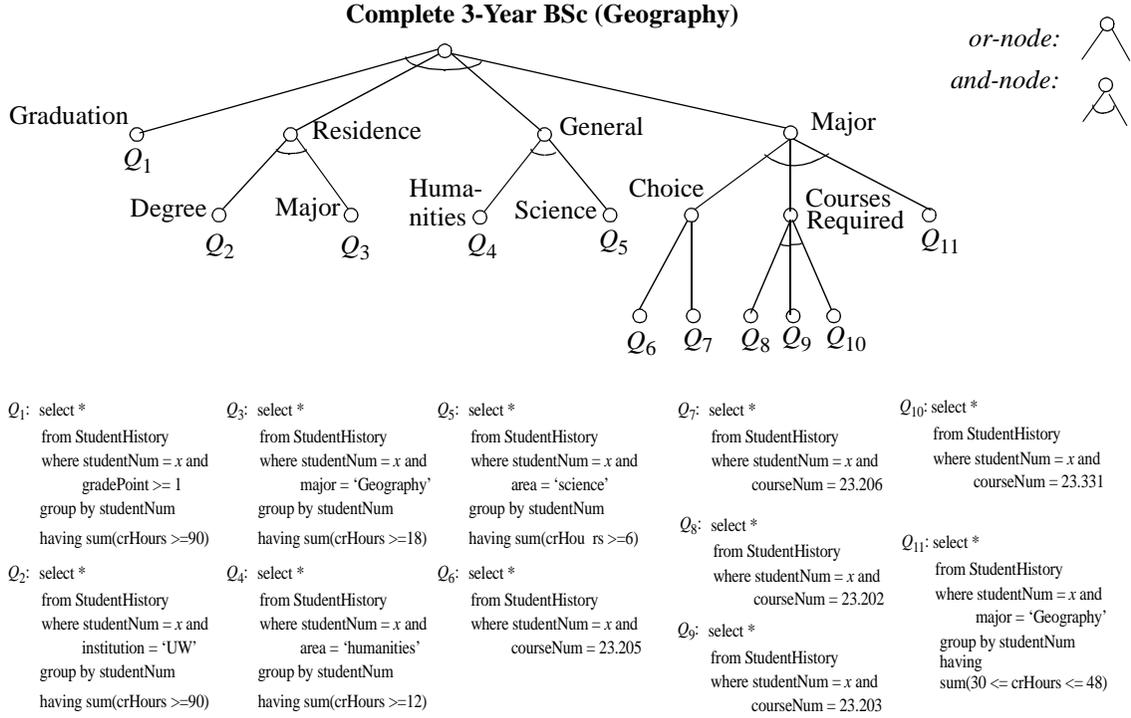
**Figure 2. XML for 3-Year BSc (Geography)**

Now, for a given set of students, to check whether they are eligible to get a degree, we traverse the query tree bottom-up. During this process, all the queries attached with the leaf nodes are evaluated against the student records and the results are transferred to the internal nodes for further checking the specified logic conditions.

For students of a different major, a different SQT will be instantiated and traversed. Obviously, if two SQTs share a common subtree, this subtree will be traversed two times. For example, the Residence requirement in the document shown in Figure 1(a) is completely the same as the Residence requirement shown in Figure 1(b). So it is better for us to remove the condition  $major = 'Geography'$  from the where-clause in the query  $Q_3$  shown in Figure 3 so that the

execution of the corresponding subtree will find all the students satisfying the subtree without considering their majors. The returned result is then separated according to their majors and transferred to the parent in the respective

SQTs. In this way, the common subtree is evaluated only once. This observation leads to the SQT fragments discussed in the next section.



**Figure 3: A Boolean SQT for graduation requirements**

### 3 Documents: fragmentation and evaluation

In this section, we discuss SQT fragmentation. First, we show how to fragment SQTs in 3.1. Then, we discuss how to evaluate fragmented SQTs in 3.2.

#### 3.1 SQT fragmentation

In [3], two kinds of SQTs: binary SQTs and general SQTs are defined. Both can be fragmented to speed up query evaluation. For simplicity, however, we show only how to fragment the binary SQTs and the general SQTs can be handled in a similar way.

First, we introduce the concept of *virtual boolean synthesized query trees*, based on which the boolean SQT fragmentation is conducted.

**Definition 2:** a *virtual boolean synthesized query tree (VBSQT)* is a tree where a leaf node  $v$  is either

- a) associated with a boolean query  $Q(v)$ , or
- b) specifies a fragment that is another VBSQT (such a leaf node is called a virtual leaf node),

and each internal node  $v$  is labelled with a tag  $T(v)$ , and an operator  $\theta = \vee$  or  $\wedge$ ; and each node  $v$  is assigned a boolean value,  $V(v)$ , determined as follows:

- a) for a leaf node that is a query,  $V(v)$  is *true* if the return value of  $Q(v)$  is not empty; otherwise, it is *false*, and
- b) for a leaf node that specifies a fragment,  $V(v)$  is the value of the fragment, and
- c) for an internal node, with children  $v_1, \dots, v_n$ ,  

$$V(v) = V(v_1)\theta V(v_2)\theta \dots \theta V(v_n).$$

For instance, the tree shown in Figure 4(a) is a virtual version of the tree shown in Figure 3, in which the leaf nodes labeled with  $v_1, v_2$ , and  $v_3$  represent the three trees shown in Figure 4(b), respectively. They are singled out since they also belong to other SQTs. To see this, examine the tree shown in Figure 4(c), which is a VBSQT for the 3-Year BSc in Physics and where  $v_1, v_2$ , and  $v_3$  are three of its leaf nodes, too.

### Complete 3-Year BSc (Geography)

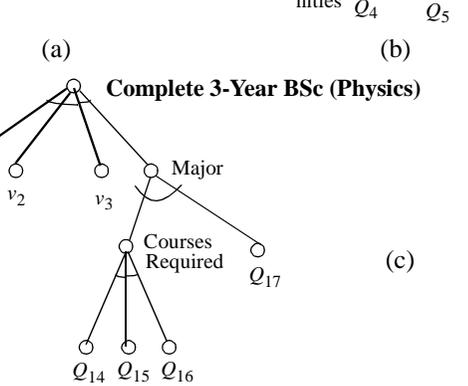
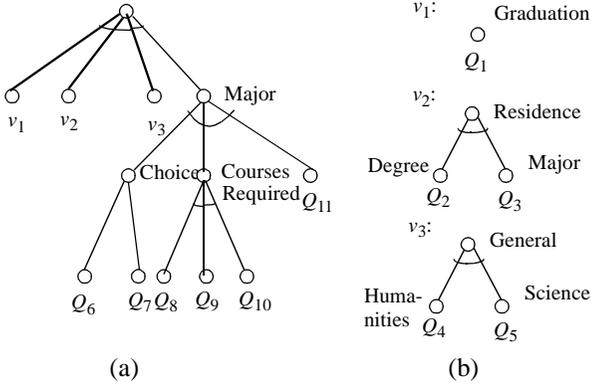


Figure 4: Illustration for SQT fragmentation

When more than two SQTs are involved, more complicated SQT fragmentation has to be considered as illustrated in Figure 5.

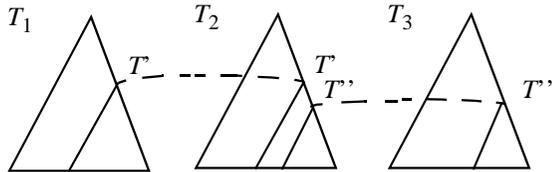


Figure 5: Fragmentation of three SQTs

In this figure, we show three SQTs:  $T_1$ ,  $T_2$ , and  $T_3$ . Among them,  $T_1$  and  $T_2$  share a common subtree  $T$ ; and  $T_2$  and  $T_3$  share a different common subtree  $T'$ . Furthermore,  $T'$  itself is a subtree of  $T$ . In such a case, we will generate five VBSQTs. They are  $T_1/T$  (which represents the tree obtained by replacing  $T$  with a virtual leaf node in  $T_1$ ),  $T_2/T$ ,  $T_3/T'$ ,  $T/T'$ , and  $T'$ .

As an example, consider a 3-Year BA (English) major where the Humanities requirement is specified, but no specification for the Science requirement. Thus, the subtree representing the Humanities requirement in the SQT for the English major is a proper subtree of the General requirement in the Geography major as illustrated in Fig-

ure 6(a), in which  $v_4$  represents the subtree shown in Figure 6(b). Accordingly, the virtual SQT for the English major will be of the form shown in Figure 7.

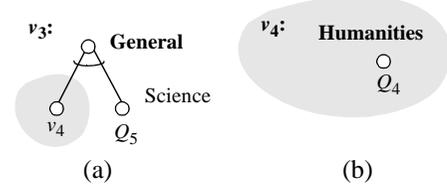


Figure 6: Humanities as a fragment of General

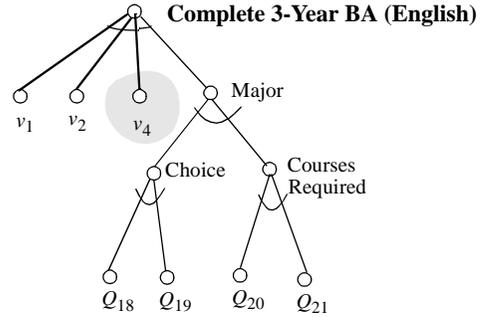


Figure 7: VBSQT for English

Note that the situation gets more complicated if we allow for some other major that specifies the Science requirement but not the Humanities. In such a case, the General requirement would then have two subtrees common to some different SQTs.

In general, we have the following algorithm to fragment any number of SQTs.

#### Algorithm SQT-fragmentation

1. Let  $T_1, T_2, \dots, T_n$  be SQT trees;
2. Let  $\Delta_{ij}^1, \dots, \Delta_{ij}^k$  be all the subtrees shared by  $T_i$  and  $T_j$  ( $i \neq j$ );
3. Repeat until fragments have been created for all  $\Delta_{ij}^l$  ( $i \neq j$ ).
  - a) From unmarked subtrees, select and mark  $\Delta_{ij}^l$  if  $\Delta_{st}^k \not\subset \Delta_{ij}^l$  for all  $\Delta_{st}^k$ , ( $i, j \neq s, t$ ). Generate a fragment for  $\Delta_{ij}^l$ . Mark any  $\Delta_{uv}^w$ , if  $\Delta_{uv}^w = \Delta_{ij}^l$ .
  - b) For each  $\Delta_{st}^k$ , if  $\Delta_{st}^k \supset \Delta_{ij}^l$ , and there is no any other  $\Delta_{uv}^w$  such that  $\Delta_{st}^k \supset \Delta_{uv}^w \supset \Delta_{ij}^l$ , replace  $\Delta_{st}^k$  by  $\Delta_{st}^k / \Delta_{ij}^l$ .
4. Generate VBSQTs:  $T_1 / \{\text{all } \Delta_{st}^k\}, \dots, T_n / \{\text{all } \Delta_{st}^k\}$ .

To explain the SQT-fragmentation algorithm we consider a more complicated scenario involving Geography Physics, and English, which are considered as  $T_1, T_2, T_3$  respectively in Step 1. Step 2 determines 9 common subtrees as shown below:

For Geography and Physics:

$\Delta_{12}^1 =$  Graduation subtree

$\Delta_{12}^2 =$  Residence subtree

$\Delta_{12}^3 =$  General subtree

For Geography and English:

$\Delta_{13}^1 =$  Graduation subtree

$\Delta_{13}^2 =$  Residence subtree

$\Delta_{13}^3 =$  Humanities subtree

For Physics and English:

$\Delta_{23}^1 =$  Graduation subtree

$\Delta_{23}^2 =$  Residence subtree

$\Delta_{23}^3 =$  Humanities subtree

In Step 3, we first generate 3 fragments for Graduation requirement  $\Delta_{12}^1$  ( $= \Delta_{13}^1 = \Delta_{23}^1$ ), Residence requirement  $\Delta_{12}^2$  ( $= \Delta_{13}^2 = \Delta_{23}^2$ ), and Humanities requirement  $\Delta_{13}^3$  ( $= \Delta_{23}^3$ ) as these do not contain any identified subtrees (see (a) in Step 3). Then, the fragment for  $\Delta_{12}^3/\Delta_{13}^3$  will be created (see (b) in Step 3). Finally, Step 4 generates the following VBSQTs:  $T_1/\{\Delta_{12}^1 \cup \Delta_{12}^2 \cup \Delta_{12}^3\}$ ,  $T_2/\{\Delta_{12}^1 \cup \Delta_{12}^2 \cup \Delta_{12}^3, T_3/\{\Delta_{12}^1 \cup \Delta_{12}^2 \cup \Delta_{13}^3\}$ .

```
<GeographyRule title= "Degree Requirement for 3-Year BSc (Geography)",
  combining="AND">
< Xinclude href="GraduationRule.xml" >
< Xinclude href="ResidenceRule.xml" >
< Xinclude href="GeneralRule.xml" >
<MajorRule> title="Major Requirement"
  combining="AND">
  <MinMaxRule
    display="Minimum 30 credit hours, Maximum 48 credit hours",
    query="... " , ...>
  </MinMaxRule>
  <ReqCoursesRule title="Required Courses" ,
    combining="AND">
    <Course
      display="23.02 Intro Geography I", query="..." , ...> </Course>
    <Course
      display="23.203 Intro Geography II", query="..." , ...> </Course>
    <Course
      display="23.331 Advanced Geography", query="..." , ...> </Course>
  </ReqCoursesRule>
  <ChoiceRule title="Choice",
    display="23.205 Atmos Sci or 23.206 Earth Sci"
    combining="OR">
    <Course query="..." , ...></Course>
    <Course query="..." , ...></Course>
  </ChoiceRule>
</MajorRule>
```

**Figure 8. Requirements document referencing fragments using the XML Xinclude feature.**

The fragmentation algorithm creates fragments that, for a given set of documents, controls redundancy present in rules by extracting common rules into separate documents. Dividing a collection of documents into sub-documents

where those sub-documents are common components is a way of structuring documents into manageable pieces that can be considered separately or in combination.

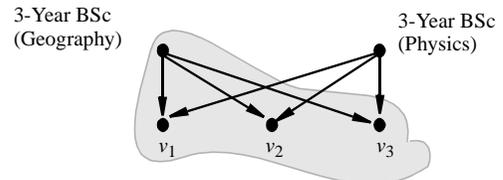
Figure 8 presents the 3-Year BSc Geography document as presented in Figures 4(a) and 4(b). Note the use of the Xinclude feature of XML [10] to link a pair of documents.

The next section discusses the evaluation of these documents which requires the documents be re-assembled in some way.

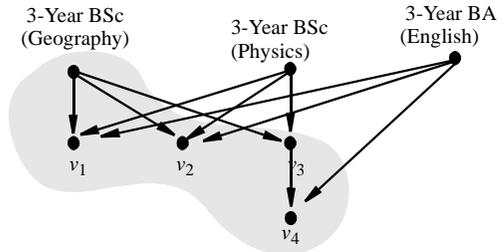
### 3.2 Evaluation of fragmented SQTs

To determine the students who can graduate, we need to evaluate the necessary SQTs for the students in question. If we were to follow the procedure in [3] we would instantiate all SQTs for each degree and every common fragment would be evaluated many times. In contrast, we present another more efficient procedure based on a fragment graph for evaluating graduation status of students.

In order to evaluate a fragmented SQT, we construct a directed graph, called a *fragment graph*, in which each node represents a fragment (or say, a VBSQT), and we have an edge from a node  $frag_a$  to another node  $frag_b$  if  $frag_b \subset frag_a$  and there is not any node  $frag_c$  such that  $frag_b \subset frag_c \subset frag_a$ . For instance, the VBSQTs shown in Figure 4 and Figure 7 can be organized into graphs as shown in Figure 9(a) and (b), respectively.



(a) Fragment graph derived from Figure 4



(b) Fragment graph derived from Figure 6

**Figure 9: Fragment graphs**

To evaluate the status of all students, we evaluate the fragment graph bottom-up. During this process, for each encountered node, we evaluate the VBSQT represented by it and transfer the result obtained to its parents. For exam-

ple, consider Figure 9(a) for students majoring in Geography and Physics. In a bottom-up fashion, we first evaluate the SQTs represented by  $v_1$ ,  $v_2$ , and  $v_3$ . The results are then partitioned according to their majors and sent to the corresponding parent nodes. In the next step, the SQTs represented by the nodes labeled with Geography and Physics are evaluated to find all the students eligible to graduate.

The above process can be improved by using the constants appearing in a query to speed up the computation. For example, if we are considering only the students majoring in Geography, we need to access only part of the graph (marked grey) shown in Figure 9(a) or Figure 9(b).

#### 4 Conclusion and Future work

In this paper, we consider a document type that includes requirements and where a user comprehends these requirements as rules to be followed to achieve a certain designation. As a result, we consider each document a single compound rule that may be assembled from many fragments. When such a document (e.g. 3-Year BSc Geography) is evaluated in a certain context (e.g. for a specific student) there will be a value generated for it. For this type of document, fragmented SQTs succinctly represent document content, evaluation and query requirements; a simple tree traversal is required to evaluate or display a document.

Since there may exist some requirement appearing many times in documents, these documents can exhibit a great deal of redundancy. We have introduced an algorithm to fragment a collection of documents, and described an efficient approach for document evaluation where each fragment/VBSQT is evaluated just once, but many results are made available (i.e. for a set of students).

We have developed a prototype system that assembles and displays requirements documents from fragments and determines on request the graduation status for students on a) an individual basis or b) a set-oriented approach for handling many students at one time. The former is useful by an individual student to measure their own progress, and the latter approach is useful in a university setting at say, the end of term, when students should be graduating. The prototype has been constructed using Java, a SAX parser, and student history data stored in a MySQL relational database. Requirements documents are stored as fragments that are independent XML documents. Various functions such as Logical And, Logical Or, Minimum, and Arithmetic Add required for the general synthesized query tree have been implemented. XML's Xinclude is used for linking (including) one document (fragment) to another. The prototype is being extended as we explore the issues identified here, and as we incorporate additional components of our university's general calendar.

We are examining algorithms for document evaluation, query optimization, and other processing models such as the pipe and filter architecture [11]. An interesting aspect of the queries we see in our sample application is that the queries in the leaves are examining similar sets of records for a given student.

We are examining other situations to apply the query tree approach for incorporating business rules into the documents themselves. We have not used functions that return data in XML format. Such functions can be used to perform an include operation, in the same way that we have used Xinclude to incorporate document fragments. Also, if we allow functions as found in [6, 7, 8] that invoke arbitrary Web Services returning XML then our model allows, as a special case, Active XML documents. We intend to examine other issues related to the processing of these types of query-based documents, such as the implications of non-relational data stores (in particular XML data stores), versioning of the general calendar, applications to workflow, use of event-condition-action (ECA) rules, and the implications of pertinent student questions such as "Which majors, given my current status, would allow me to gradually most quickly?" or "What are the added requirements if I were to do a double major in Geography and Physics, instead of a single major in Geography"?

#### 5 References

- [1] Business Rules Group, Defining business rules: What are they really?, 3rd. edition, July 2000, <http://www.BusinessRulesGroup.org>.
- [2] R. G. Ross, The business rule book, 2nd. edition, Business Rules Solutions, Houston, 1997.
- [3] Ron McFadyen, Yangjun Chen, Fung-Yee Chan, XML-based evaluation of synthesized queries, 1st International Conference on Web Information Systems and Technologies (WEBIST 2005), May 26-28, 2005, Miami, USA.
- [4] Airi Salminen, Frank Wm. Tompa, Requirements for XML document database systems, Document Engineering, November 9-10, 2001, Atlanta, Georgia, USA.
- [5] Vincent Quint, Irene Vatton, Techniques for authoring complex XML documents, ACM 2004
- [6] Angela Bonifati, Stefano Ceri, Stefano Paraboschi, Active rules for XML: A new paradigm for E-services, VLDB Journal 10: 39-47 (2001)
- [7] S. Abiteboul, O. Benjelloun, I. Manolescu, T. Milo, R. Weber, Active XML: peer-to-peer data and web services integration (demo), *Proceedings of VLDB*, 2002.[9] S. Abiteboul, O. Benjelloun, I. Manolescu,

T. Milo, R. Weber, Active XML: peer-to-peer data and web services, Proceedings of VLDB, 2002.

[8] Serge Abiteboul, Angela Bonifati, Gregory Cobena, Ioana Manolescu, Tova Milo, Dynamic XML documents with distribution and replication, SIGMOD 2003, June 9-12, 2003, San Diego, CA, USA.

[9] James Bailey, Alexandra Poulouvasilis, Peter T. Wood, An event-condition-action language for XML, WWW2002, May 7-11, 2002, Honolulu, Hawaii, USA.

[10] <http://www.w3.org/TR/xinclude/>, retrieved June 8, 2005

[11] Stephen T. Albin, The art of software architecture: design methods and techniques, 1<sup>st</sup> edition, ISBN 0471228869, Wiley Publishing, 2003.