

OrientXA:一种有效的 XQuery 查询代数*

孟小峰⁺, 罗道锋, 蒋瑜, 王宇

(中国人民大学 信息学院, 北京 100872)

OrientXA: An Effective XQuery Algebra

MENG Xiao-Feng⁺, LUO Dao-Feng, JIANG Yu, WANG Yu

(Information School, Renmin University of China, Beijing 100872, China)

+ Corresponding author: Phn: +86-10-62519453, E-mail: xfmeng@public.bta.net.cn, <http://www.ccf-dbs.org.cn/idke/xfmeng/>

Received 2004-05-11; Accepted 2004-09-07

Meng XF, Luo DF, Jiang Y, Wang Y. OrientXA: An effective XQuery algebra. *Journal of Software*, 2004, 15(11):1648-1660.

<http://www.jos.org.cn/1000-9825/15/1648.htm>

Abstract: XQuery is the recommended standard for XML Query. XQuery processing strategies can be classified into two categories: core syntax based strategy (node-oriented) and algebra based strategy (set-oriented). Neither of them can handle XQuery well. The syntax based strategy is inefficient and hard to optimize, while the current algebra based strategies can not satisfy the flexible programming characteristics of XQuery. After summarizing the current stage and unsolved problems of former algebra based works, an effective XQuery algebra system, OrientXA, is proposed in this paper, ideas from both strategies are embodied in it. OrientXA introduces the notion of Construct Pattern Tree for the first time. The Construct operator in it materializes the flexible characteristics of XQuery. Corresponding to its expressive operators, it is able to express all the queries in W3C use cases and XMark benchmark.

Key words: XML; query processing; Xquery; algebra; core syntax

摘要: XQuery是XML查询语言的推荐标准.现有的XQuery处理策略有两种方法,一种是基于核心语法一次一结点的方法,另一种是基于查询代数的一次一集合的方法.单独使用这两种方法都不能很好地处理XQuery查询.一方面,基于核心语法树的方法效率低,而且很难优化,而现有基于代数的方法又不能很好地满足XQuery的过程化查询语言的特性.描述了一个有效的XQuery代数系统——OrientXA.它是在总结现有的基于代数的XQuery查询工作、分析未解决的问题,并吸收了前述的两种策略的思想的基础上提出来的.它第一次提出了构造模式树的概念,其中的构造操作符能体现XQuery的灵活特征.由于操作符具有很强的表达能力,OrientXA能

* Supported by the National Natural Science Foundation of China under Grant No.60273018 (国家自然科学基金); the National High-Tech Research and Development Plan of China under Grant No.2002AA116030 (国家高技术研究发展计划(863)); the Key Project of the Chinese Ministry of Education under Grant No.03044 (国家教育部科学技术研究重点项目)

作者简介: 孟小峰(1964—),男,博士,教授,主要研究领域为XML数据库,Web数据集成,移动数据库技术;罗道锋(1978—),男,硕士,主要研究领域为XML存储和查询;蒋瑜(1981—),男,硕士,主要研究领域为XML查询,索引技术;王宇(1973—),女,博士,主要研究领域为XML查询优化技术.

够表达 W3C 的用例和 XMark 测试集的所有查询.

关键词: XML;查询处理;Xquery;代数;核心语法

中图法分类号: TP311 文献标识码: A

随着 XML 数据被越来越广泛地使用,人们对 XML 数据的查询提出了越来越高的要求.前期的研究主要集中在 XPath 的查询处理上.XPath 相对比较简单,表达能力有限,比如,不能表示连接操作等.作为 W3C 的推荐标准,XQuery 的表达能力比 XPath 强大很多,同时,XQuery 也比 XPath 复杂很多.

XQuery 兼有结构化查询语言和过程化编程语言的特点.一方面,XQuery 的 FLWR 子句在一定程度上类似于 SQL 的 Select-From-Where 子句,是 XQuery 的最重要的表达式;另一方面,XQuery 支持表达式的任意嵌套,支持诸如条件表达式(IF-THEN-ELSE)、循环表达式(FOR)、返回值(RETURN)等,还有变量和谓词的作用域问题,这些都是一种编程语言的重要特征.

因为 XQuery 兼有这两种特点,因此,XQuery 的处理方法分为两大类:基于核心语法的处理和基于 XML 代数的处理.

- 基于核心语法的处理(一次一结点)

现在流行的 XQuery 查询引擎,比如 Galax^[1],IPSI^[2]等,使用的都是这种方法.这种方法把 XQuery 看作一门编程语言来处理,通常先把 XQuery 转换成核心语法(core syntax),然后以 XML 树的根结点作为输入,依次执行该核心语法树,最后得到结果.

比如,对于以下查询:

Q1:

```
FOR $b in (document("bib.xml")/bib/book
WHERE $b/price<50
RETURN
<cheap-books>{$b/title}</cheap-books>
```

转换成的核心语法是(由于核心语法表示很复杂,这里用简化的伪码来表示):

1) for \$bib in root/bib do	8) endfor
2) for \$book in \$bib/book do	9) break;
3) for \$price in \$book/price do	10) else ()
4) if \$price<50 then	11) endfor
5) construct a <cheap-books>	12) endfor
6) for \$title in \$book/title do	13) endfor
7) append the content of \$title as child of the element <cheap-books>	

可以看出,这种方法的特点是一次一结点(node-oriented).这种方法比较简单,其优点是:它天生符合 XQuery 的程序化的特点,能正确地处理各种 XQuery 语句的嵌套,并能自动地保证结点的正确的顺序.它最大的劣势在于很难进行优化.关于一次一结点和一次一集合的方式,人们的共识是,一次一集合方式的效率要优于一次一结点,而且有利于优化.因此,XQuery 的另一条处理路线是一次一集合的方式.

- 基于 XML 代数的处理(一次一集合)

这种方法提出了一套类似于关系代数的代数体系.每一个代数操作符的输入都是 1 个或者多个 XML 树集合,输出也是一个 XML 树集合.通常,为了生成 XML 树集合,从 XQuery 语句中生成 1 个或多个模式树(pattern tree),模式树表示了该查询感兴趣的(包括谓词结点和目标结点)一组有祖先后代关系的变量绑定,然后,用模式树从输入的 XML 树中抽取实例树.比如,对于查询 1,会先抽取模式树 bib(book(price,title)),然后用这个模式树对 XML 树进行匹配,得到模式树的实例树集.随后在这个实例树集上进行谓词判断(price<5),最后输出希望的结果.

可以看出,这种方式是一次一集合的处理方式.它最大的优点是能够对逻辑操作树进行优化,从而获得较高的效率.然而,一次一集合的方式与 XQuery 的编程语言特点和 XML 数据本身的特点存在许多不相容的地方,

从而带来一些新的问题.比如,如何保证结点顺序、如何处理谓词的作用域等.

无论是哪种方法,都不能很好地解决 XQuery 查询处理的问题.一方面,一次一结点的操作效率很低,而且很难优化.在关系数据库中,一次一集合的操作被证明其效率远高于一次一结点的操作.另一方面,一次一集合的方法又和 XQuery 过程化编程语言的特点不相符.XQuery 是任意嵌套的语言,特别是在结果构造符中,一次一集合的操作有时不可能表达复杂的结果嵌套情况,而必须借用编程语言的线性执行的特点.

因此,本文提出的 OrientXA 试图结合这两种方法,为 XQuery 提供一个有效的代数系统.这种结合主要体现在结果构造符中.

本文的主要贡献有:

(a) 总结并分析了现有的查询代数的研究现状和问题所在,根据 XQuery 语言的特点,明确提出了包括强绑定/弱绑定、结点绑定/序列绑定等概念;

(b) 提出了一套有效的 XML 查询代数系统,引入 Sequence 操作符,有效结合了两种 XQuery 处理策略的特点,能有效地表达 XQuery 查询,这是 OrientXA 区别于现有 XML 代数系统的最大不同之处;

(c) 提出了大量 XQuery 中有别于关系代数的逻辑优化的问题.

本文第 1 节分析 XML 代数设计的现状.第 2 节提出 XML 代数设计中新的问题.第 3 节介绍 OrientXA 查询,并给出 3 个示例.第 4 节讨论 XQuery 代数处理中优化的问题.第 5 节是结论和未来工作.

1 相关工作

XML 代数的目的是使用一次一集合的方式来处理查询.逻辑操作符的输入是记录的集合.问题是:什么是记录;XML 代数与关系代数相比有什么不同.

目前,人们提出了很多 XML 代数^[3~11],其中著名的包括 SAL^[3],Xtasy^[4],TTX^[5],XAL^[6],TAX^[7]等.

1.1 XML代数处理的对象

XML 代数处理的对象是记录的集合.在 TAX 中,记录是 XML 树.TAX 引入了 pattern tree 和 witness tree 的概念.pattern tree 就是一组相关的变量绑定,witness tree 是 pattern tree 的实例,如图 1 所示.

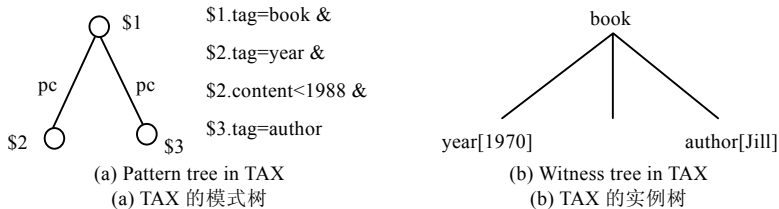


Fig.1 Pattern tree and Witness tree in TAX

图 1 TAX 的模式树和实例树

在 Xtasy 和 TTX 中,记录是一种叫做 Env 的结构.它们首先使用 input filter(类似于 pattern tree)对输入文档用 path 操作符进行过滤,过滤的结果就是 Env 的集合.Input filter 的定义是:

$$F:=(op;var;binder)label[F].$$

它的含义是:从当前上下文结点找到满足 op 关系(/,//,或者@)的结点名为 lable 的结点,并用 binder 方式(in 绑定或者=绑定)绑定到变量 var,然后以当前变量绑定为上下文结点,处理嵌套的 F.

比如,对于 XQuery 查询:

```
FOR $b in book,
  $a in $b//author,
```

input filter 是 (_;\$b;in) book [(//;\$a;in)author[\emptyset]],表示以当前上下文结点找到 book 结点并“in”绑定到\$b,然后以 \$b 为上下文结点,找到满足//关系的 author 结点并“in”绑定到\$a.

```
FOR $b in book,
  LET $a in $b//author,
```

input filter 是 (`_;$b;in`)`book[(//;$a;=)author[\emptyset]]`,唯一的区别是对 `$a` 的绑定是“=”绑定.

绑定的结果是 Env 结构的集合.一个 Env 结构就是一个记录.如图 2 所示.其中, o_i 表示 Object ID.在 W3C 的 XML data model 中,每一个结点都有一个全局唯一的 ID.

TTX 中的处理与 Xtasy 很相似,只不过 TTX 把 input filter 变成树状表示,把 Env 结构也表示成树状.

<code>\$b:o₁</code>	<code>\$a:o₃</code>
<code>\$b:o₁</code>	<code>\$a:o₄</code>
<code>\$b:o₁</code>	<code>\$a:o₅</code>
...	...

<code>\$b:o₁</code>	<code>\$a: {o₁,o₂,o₃}</code>
...	...

(a) Env Set in Xtasy (in binding) (b) Env Set in Xtasy (= binding)

Fig.2 Flat Env in Xtasy
图 2 XTasy 的扁平 Env 结构

1.2 XML代数的操作符

XML 代数与关系代数一样,都需要表达 `select,project,product,rename,join` 等操作.我们把这类操作符叫做类关系操作符.那么,XML 代数还需要哪些额外的操作符呢?由于 XQuery 语言和 XML 数据本身的特点,XML 代数还应该有两类操作符:Extraction 操作和 Construction 操作.

• Extraction 操作符

在关系数据查询处理中,操作符处理的对象是结构化的扁平的元组.但是在 XML 查询处理中,XML 查询只是对 XML 文档上的某些特定结点有兴趣(谓词结点和目标结点),而这些结点往往只是整个文档的一小部分.所以,必须先感兴趣结点从输入文档中抽取出来,才能做下一步的处理.

TAX 把数据抽取操作合并在选择操作或者投影操作中,选择或者投影操作符的输入是 `pattern tree` 和谓词列表 PL.在 Xtasy 和 TTX 中,使用 `path` 操作符来进行数据抽取,`path` 操作符的输入就是 `input filter`.

实际上,`pattern tree` 或者 `input filter` 可以视为 XPath 的树状表示.所以,Extraction 操作实际上是 XPath 处理,而 XPath 是 XQuery 的核心部分.所以,Extraction 操作是 XML 代数中最影响查询效率的操作符之一.目前,处理 Extraction 操作的方法主要有两种:

(1) Navigation:对文档树进行遍历,找到满足 `pattern tree` 的实例树.这种方法比较适合绝对路径的处理.

(2) Structure join:利用对 XML 数据的编码和 Tag Index,快速地找到满足祖先后代关系的结点.这种方法比较适合相对路径(`//`)的处理.

如何混合使用这两种方法,是目前的研究热点之一.

• Construction 操作符

在 SQL 中,除了重命名以外,并没有结果构造的需求.但是,在 XML 查询处理中,往往要求查询结果以一定格式输出.所以,XML 代数中还需要结果构造符来满足这个需要.

在 Xtasy 和 TTX 中,结果构造符是 `return` 操作符.`return` 操作符的输入是 Env 结构和一个 `output filter,output filter` 就是想要输出的结果的框架.`output filter` 的定义是:

$$\begin{aligned}
 OF &::= OF_1; \dots ; OF_n \\
 &| \text{label}[OF] \\
 &| @\text{label}[val] \\
 &| val \\
 val &::= v_B | var | vvar
 \end{aligned}$$

`label[OF]`表示元素构造,元素结点是 `label`,内容由 `[OF]`指定;`@label[val]`是属性构造符,属性名是 `label`,值是 `val.val`可以是自定义的原子值 v_b ,或者拷贝(`vvar`)或者引用(`var`).`vvar`是对结点的拷贝,不会保留原来的 `oid`,从而与源数据不再有联系;`var`是对结点的引用,会保留原来的 `oid`,从而与原来的数据仍有联系.所以 `var`可以用来定义视图.

比如,对于以下查询:

```
FOR $b in document("bib.xml")//book
WHERE $b/year<2001
```

```
RETURN <books>{ $b }</books>
```

转换成操作树是:

```
returnbooks[$b]( $\sigma_{Sy < 2001}$ (path(_;$b;in)book
[(/;$y;in)year[ $\phi$ ]](db1)))
```

在 TAX 中,没有一个明确的结构构造符,它用 grout-by,copy-and-paste,rename 来完成结果构造的工作.

2 问题的提出

OrientXA 借鉴了上面的代数的思想.在 OrientXA 中,记录是 XML 树,操作符处理的对象是 XML 树的集合.OrientXA 使用 pattern tree 来表示感兴趣的结点,使用选择操作符来做抽取操作,使用构造操作来进行结果构造.在设计并实现 OrientXA 的过程中,我们遇到了一些新的问题,需要新的方法来处理,同时这些问题对查询优化提出了新的要求.

2.1 序列化操作符

XQuery 是可以任意嵌套的查询语言.这给代数式处理带来了额外的困难.文献[12]对 OO 的嵌套化进行了讨论,但是 XQuery 的嵌套比 OO 的嵌套更复杂.XQuery 更像一个第三代编程语言的代码,而不像 SQL 那样的结构化查询语言.有时候,由于表达式的任意嵌套,一些操作符必须序列化执行.比如:

<pre>Q2: <bib> { FOR \$b in doc("http://bstore1.example.com/bib.xml")//book WHERE count(\$b/author)>0 RETURN <book> { \$b/title } { FOR \$a in \$b/author[position()<=2]</pre>	<pre>RETURN \$a } { if (count(\$b/author)>2) then <et-al/> else () } } } </bib></pre>
--	--

从 Q2 产生的 Pattern Tree 是 book(author,title).在第 1 个 RETURN 子句中,嵌套了 3 个表达式:一个 XPath 表达式、一个 XPath 表达式、一个 FLWR 表达式和一个 IF-THEN-ELSE 表达式.给定一个实例树,这 3 个表达式应该被序列化地执行.需要注意的是,序列化执行和流水线执行是不同的.给定操作 op_1, op_2, \dots, op_n ,所谓流水线执行是指, op_{i+1} 的输入是 op_i 的输出;而序列化操作是指,对于一个记录 r ,顺序地执行 op_1, op_2, \dots, op_n .

序列化操作实际上引入了一次一结点的处理思想.现有的 XML 代数和传统的关系代数都没有这样的操作符,因此不能很好地处理任意嵌套化的 XML 查询.OrientXA 引入了序列化操作符,很好地结合了一次一集合和一次一结点的处理策略.

2.2 强绑定和弱绑定

不是所有的变量绑定都要求一定绑定到结点,有些绑定可以绑定到空结点(或空结点集).我们把前者叫做强绑定,后者叫做弱绑定.pattern tree 是相关变量绑定的树状表示,一个实例树应该满足 pattern tree 的所有的强绑定,而不要求必须满足弱绑定.在 OrientXA 中,对 pattern tree 的变量绑定,都标明了是强绑定还是弱绑定.文献[13]注意到了这个问题,但是并没有讨论哪些是强绑定.

我们发现,除了 FOR 语句的绑定是强绑定之外,其他语句中的绑定(如 LET 语句,WHERE 语句和 RETURN 语句等)都是弱绑定.

比如对于以下查询:

```
FOR $b in (document(bib.xml)//bib/book
LET $a := $b/author
WHERE not ($b/year<1998)
```

```

RETURN
  <authorss isbn="{ $b/@isbn }">
    { $a }
  </authors>

```

FOR 子句把 book 显式地绑定到 \$b,同时还包含了对 bib 的隐式绑定;LET 子句把 author 显式地绑定到 \$a,而 WHERE 子句中 \$b/year 对 year 的绑定和 return 子句中 \$b/@isbn 对 isbn 的绑定则是隐式绑定.

在这些绑定中,\$b 是强绑定.需要注意的是,如果 pattern tree 上的一个结点是强绑定,那么对 pattern tree 上它的祖先结点也是强绑定.所以,对 bib 的隐式绑定是强绑定,而其他隐式绑定都是弱绑定.如果一个 book 结点没有 author,year 等子结点或 isbn 属性,该 book 结点仍应该被视为符合该 pattern tree 的实例树.

强绑定可以过滤更多的结点,使结果集更小,从而有利于提高查询效率.有些情况下,一些弱绑定结点可以转化为强绑定结点,比如,如果 WHERE 子句变为

```
WHERE $b/year<1998
```

那么,我们可以认为对 year 的绑定是强绑定,因为如果 book 结点没有 year 子结点,那么它肯定不满足这个谓词,因此该 book 结点就会被过滤掉.

尽可能地把弱绑定转化为强绑定,是 XML 查询处理中一个很大的优化点.

2.3 结点绑定和序列绑定

在 XQuery 中,一个变量可以绑定到一个结点,如 FOR 绑定,也可以绑定到一个结点序列(sequence),如 LET 绑定、隐式绑定等.我们把前者叫做结点绑定,后者叫做序列绑定.在 XQuery 中,只有 FOR 绑定是结点绑定.结点绑定和序列绑定对结果的构造是有影响的.比如,下列两个相似的查询:

<pre> Q3: <books> FOR \$b in document("bib.xml")//book LET \$a:=\$b/author WHERE \$b/year<2001 RETURN <authors> <aNum>{count(\$a)} </aNum> { \$a } </authors> </books> </pre>	<pre> Example Results: <books> <authors> <aNum>2</aNum> <author>Smith</author> <author>Tom</author> </authors> ... </books> </pre>
<pre> Q4: <books> FOR \$b in document("bib.xml")//book \$a:=\$b/author WHERE \$b/year<2001 RETURN <authors> <acouont>{count(\$a)} </acouont> { \$a } </authors> </books> </pre>	<pre> Example Results: <books> <authors> <aNum>1</aNum> <author>Smith</author> </authors> <authors> <aNum>1</aNum> <author>Smith</author> </authors> ... </books> </pre>

Q3 对 author 的绑定是序列绑定,而 Q4 对 author 的绑定是结点绑定,它们的结果就很不一样.

在 Xtasy 和 TTX 中,用“in”绑定来表示结点绑定,用“=”来表示序列绑定.

TAX 中没有显式地标明.在 TAX 中,如果一个 pattern tree 被用在选择操作符上,则应用结点绑定;如果用在

投影操作符上,则应用序列绑定.但是,如果一个 pattern tree 上同时有结点绑定和序列绑定,就不好处理了.

文献[11]提出了一个 Tree Logical Class 的概念,用来表示序列绑定的概念.它对 Pattern Tree 上的结点做标记,来表明它是哪种绑定类型.

OrientXA 借鉴 Xtasy 和文献[11]的思路,显式地标明 pattern tree 的结点哪些是结点绑定,哪些是序列绑定.

在 Xtasy 中,进行 path 操作的时候,对于结点绑定,每一个结点应该有一个单独的 instance tree(Env 结构),如图 2(a)所示;对于序列绑定,一个序列里的所有结点应该在一个 instance tree 中,如图 2(b)所示.

但是,问题似乎不这么简单.如果对一个结点的绑定是序列绑定,但是对它的子结点的绑定是结点绑定,那么 instance 怎么来表示?比如,对于以下变量绑定:

```
for $b in document("bib.xml")/bib
let $book := $b/book
for $a in $book/author
```

假设某一个 bib 下面有 4 个 book 子结点,由于是序列绑定,按照结点绑定和序列绑定的概念,instance tree 应该被表示为 bib1(book1,book2,book3,book4).接下来,对于 author 是结点绑定,所以,每一个不同的 author 结点都应该有一个 instance tree.于是,就出现了这种奇怪的变量绑定结果:

```
bib1(book1(author1), book2, book3, book4), bib1(book1(author2),book2,book3,book4)
```

这里,(book1, book2, book3, book4)在每一个实例中都要重复一遍,是很大的冗余,这对查询处理的效率影响很大,因为它会大大增大中间结果的大小.

OrientXA 正在研究如何在保证语义正确性的情况下减少这种冗余.这是我们未来优化工作的研究点之一.

2.4 选择操作和投影操作

在 SQL 中,选择操作和投影操作的区别是很明显的.选择操作选择表中的某些行,投影操作选择表中的某些列.在 XML 查询处理中,记录不是扁平的,它是嵌套的树结构.那么,在 XML 查询处理中,选择操作和投影操作的区别是什么呢?

Xtasy 和 TTX 认为不需要投影操作.path 操作用来抽取感兴趣的结点,然后用选择操作在上面作谓词判断.

在 TAX 中,选择操作和投影操作被用来区别结点绑定和序列绑定.

OrientXA 综合了 Xtasy 的思想:不引入投影操作符,也不引入专门的 Extract 操作符,而是直接用选择操作符来完成这个功能.

3 OrientXA 设计

OrientXA 的代数设计充分借鉴了已有的 XML 代数的优点,并且充分考虑了上述新的问题,并针对这些问题作了一些特殊的设计.本节介绍 OrientXA 的具体设计.

3.1 定义

Source Pattern Tree:一个 Source Pattern Tree 是一棵树(N, E, R),其中, N 是结点集合, E 是边的集合, R 是树的根.

每个结点有结点名,结点名可以为空,表示可以取任意结点名.结点名在 Pattern Tree 上不唯一.为了唯一地标识一个结点,每个结点分配一个 PID.每一个结点还有若干个修饰符,主要包括:(1) 是否序列绑定:用双圆圈表示序列绑定的结点,单圆圈表示结点绑定的结点;(2) 是否强绑定:用实线圆圈表示强绑定,虚线圆圈表示弱绑定;(3) 是否连带绑定所有子孙结点:用 p 表示需要连带绑定所有子孙结点.

边用来表达结点间的关系.Pattern Tree 上有 3 种边:(1) 父子边,用单实线表示;(2) 祖先后代边,用双实线表示;(3) 元素属性边,用单虚线表示.

Constructor Pattern Tree:为了表示结果构造操作,引入 Constructor Pattern Tree.注意,Constructor Pattern Tree 上的边都是父子边或者元素属性边.

结点定义如下:

$$N := n(\text{TagName}) [= val] \\ | c(\text{PID}) \\ val := \text{atom} | \text{PID}$$

其中, $n(\text{TagName})=val$ 表示新建一个结点,结点名为 TagName ,结点值为 val ; $c(\text{PID})$ 表示拷贝 input pattern tree 上特定 PID 结点,包括结点名、结点值,以及所有子孙结点; val 可以取原子值 atom,也可以取某个 PID 结点的文本值.

Constructor Pattern Tree 的结点可以有以下修饰符:是否序列构造:用双圆圈表示序列构造.序列构造只用于拷贝构造结点,表示输入集所有拷贝的结点挂载在同一个父结点下面.

为了方便描述,给出以下定义:

- 如果 X 是一个 XML 树集合, x 是该集中一棵树,记做 $x \in X$.
- 如果 s 是 x 的一棵子树,记做 $s \in x, s \prec X$.
- 如果 P 是一个 Pattern Tree, x 是 P 的实例树,记做 $P(x)$.
- 如果 PE 是一个值谓词表达式,实例树 x 上的相应结点满足该谓词,记做 $PE(x)$.

3.2 OrientXA操作符

每个操作符实际上都有一个 P_i 表示输入数据的格式,有一个 P_o 表示输出数据的格式.我们把 P_i 叫做输入 Pattern,把 P_o 叫做输出 Pattern.这里,有几个规则:

- (1) 选择操作符的 P_i 可以为空,其他操作的 P_i 均不能为空;
- (2) 结构构造操作符的 P_o 是 Constructor Pattern Tree,其他均为 Source Pattern Tree;
- (3) 每一个操作的 P_o 是下一个操作的 P_i .

选择操作符: $\sigma_{P_i, P_o, PE}(X) = \{x | x \prec X, P_o(x), PE(x)\}$.

输入: P_i :输入 pattern tree,是 Source Pattern Tree,可以为空; P_o :输出 pattern tree,是 Source Pattern Tree,不能为空; PE :谓词链表,可以为空; X :XML 树集合.

输出: P_o 的实例树集合.

输出 Pattern: P_o .

结果构造操作符: $\chi_{P_i, P_o}(X)$.

输入: P_i :输入 pattern tree,是 Source Pattern Tree,不能为空; P_o :输出 pattern tree,是 Constructor Pattern Tree,不能为空; X :XML 树集合.

输出:按照 P_o 构造的实例树集合.

输出 Pattern: P_o .

序列操作符: $\xi_{P_i, (op_1, op_2, \dots, op_n)}(X) = \{s(op_1(x), op_2(x), \dots, op_n(x)) | x \in X\}$.

输入: P_i :输入 Pattern,是 Source Pattern Tree,不能为空. op_i :一元操作符; X :XML 树集合.

输出: $s(op_1(x), op_2(x), \dots, op_n(x))$,其中, s 是虚根结点.

输出 Pattern: $s(P_{o1}, P_{o2}, \dots, P_{on})$,其中, P_{oi} 是 op_i 的输出 Pattern.

连接操作符:

$$\triangleright_{P_i, P_2, C}(X_1, X_2) = \{j(x_1, x_2) | x_1 \in X_1, x_2 \in X_2, C(x_1, x_2)\}, \\ \otimes_{P_i, P_2, C}(X_1, X_2) = \{j(x_1, (x_{21}, x_{22}, \dots, x_{2n})) | x_1 \in X_1, x_{2i} \in X_2, C(x_1, x_{2i})\}.$$

输入: P_i :输入 pattern tree,是 Source Pattern Tree,不能为空; P_o :输出 pattern tree,是 Source Pattern Tree,不能为空; C :连接谓词; X :XML 树集合.

输出:满足连接谓词的 $j(x_1, x_2)$ 集合,其中, j 是虚根结点.

输出 Pattern: $j(P1, P2)$.

还可以在此基础上定义左连接.注意,没有右连接,因为由于 XQuery 的特性,不会出现右连接的情况:

$$L \triangleright_{P_i, P_2, C}(X_1, X_2), L \otimes_{P_i, P_2, C}(X_1, X_2).$$

分组操作符: $\gamma_{p,f}(X) = \{g(v, r(x_1, x_2, \dots, x_n)) \mid x_i \in X, v = f(x_i)\}$.

输入: P : 输入 pattern tree, 是 Source Pattern Tree, 不能为空; f : 分组函数; X : XML 树集合.

输出: $g(v, r(x_1, x_2, \dots, x_n))$, 其中 x_1, x_2, \dots, x_n 具有相同的分组函数值 v ; g 是虚根.

输出 Pattern: $g(v, r(P))$.

聚集操作符: $A_{p,f}(X) = \{a(v, x) \mid x \in X, v = f(X)\}$.

输入: P : 输入 pattern tree, 是 Source Pattern Tree, 不能为空; f : 聚集函数, 如 sum, count, average 等; X : XML 树集合.

输出: $a(v, x)$, v 是聚集函数值.

输出 Pattern: $a(v, P)$.

消重操作符: $\delta_{p,f}(X) = \{x \mid x \in X, \forall x_i \in \delta(X), \text{if } x_i \neq x \Rightarrow f(x_i) \neq f(x)\}$.

输入: P : 输入 pattern tree, 是 Source Pattern Tree, 不能为空; f : 消重函数, 对于每一个输入的 XML 树, 返回一个可比较的值; X : XML 树集合.

输出: 消除重复值之后的 XML 树集合.

输出 Pattern: P .

排序操作符: $\tau_{p,f}(X)$.

输入: P : 输入 pattern tree, 是 Source Pattern Tree, 不能为空; f : 排序函数, 对于每一个输入的 XML 树, 返回一个可排序的值; X : XML 树集合.

输出: 排序后的 XML 树集合.

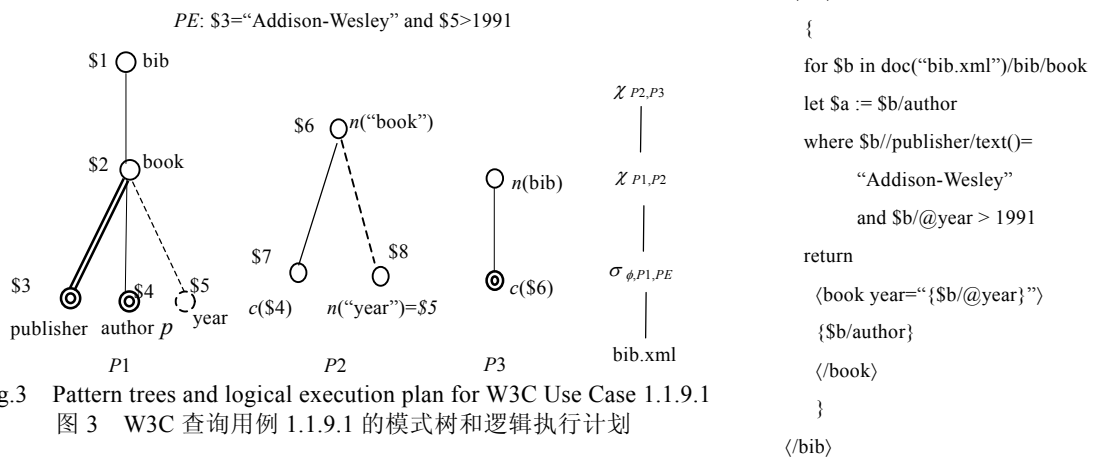
输出 Pattern: P .

3.3 OrientXA 示例

图例:

- 父子边
- ==== 祖先后代边
- 元素属性边
- p : 绑定所有后代结点
- n : 新建结点
- c : 拷贝结点
- 强结点绑定
- 弱结点绑定
- ◎ 序列绑定或序列构造

示例 1: W3C Use Case 1.1.9.1 (如图 3 所示).



这个例子很简单,先用 $P1$ 和 PE 抽取出感兴趣的数据,然后进行结果构造.

示例 2: W3C Use Case 1.1.9.4 (如图 4 所示).

```

<results>
{
  LET $a:=doc("http://bstore1.example.com/bib/bib.xml")//
  author
  FOR $last in distinct-values($a/last),
    $first in distinct-values($a[last=$last]/first)
  order by $last, $first
  RETURN
  <result>
  <author>
  <last>{ $last }</last>
  <first>{ $first }</first>
}
</results>

</author>
{
  FOR $b in doc("http://bstore1.example.com/bib.xml")/bib
  /book
  WHERE some $ba in $b/author
  satisfies ($ba/last = $last and $ba/first=$first)
  RETURN $b/title
}
</result>
</results>
    
```

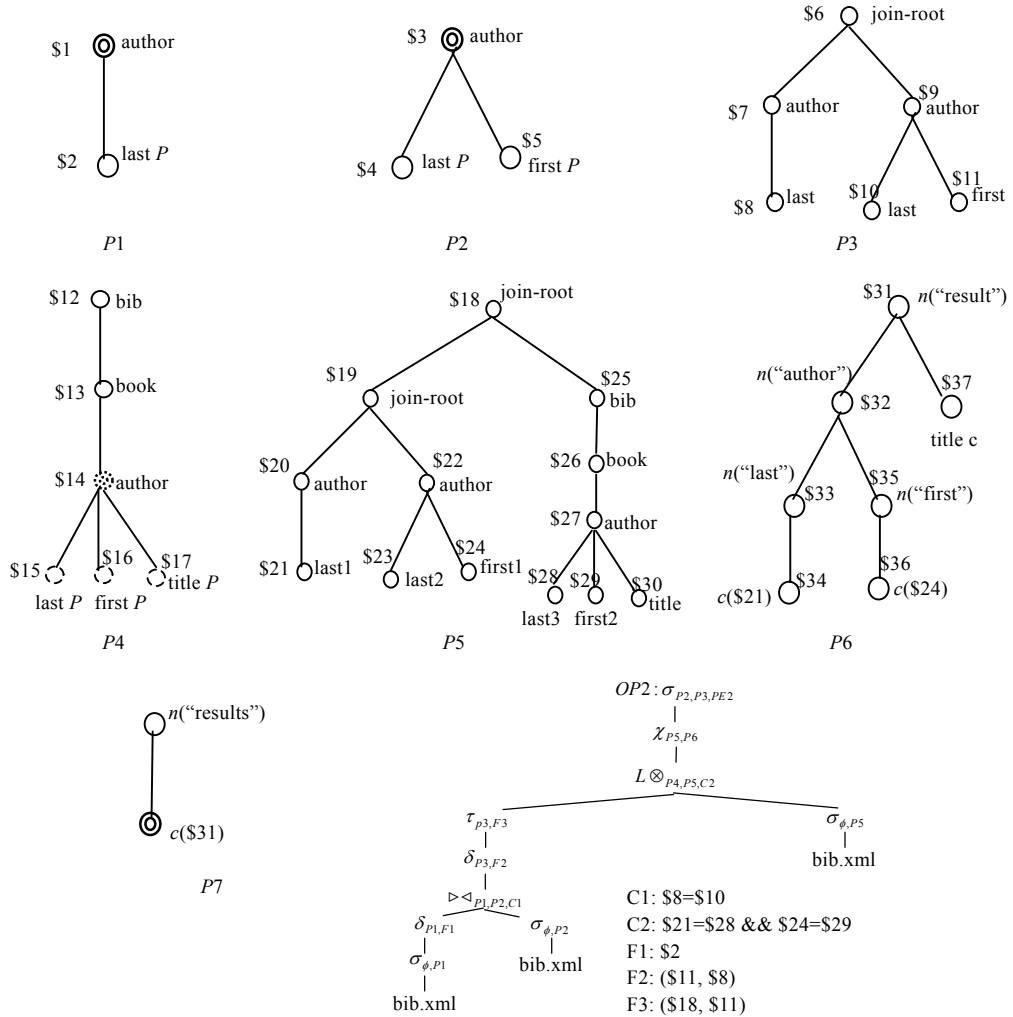


Fig.4 Pattern trees and logical execution plan for W3C Use Case 1.1.9.4
 图4 W3C 查询用例 1.1.9.4 的模式树和逻辑执行计划

示例 2 比较复杂.先用 $P1$ 抽取出数据,做一个消重操作;然后用 $P2$ 抽取出数据,做一个连接操作,并对连接的结果消重并排序;然后用 $P4$ 抽取出数据,跟上次的连接结果做一个左连接,最后进行结果构造.

示例 3: W3C Use Case 1.1.9.6(如图 5 所示).

```

<bib>
  {
    FOR $b in
doc("http://bstore1.example.com/bib.xml")//book
  WHERE count($b/author)>0
  RETURN
    (book)
      { $b/title }
      {
        FOR $a in $b/author[position()<=2]
      }
  }
</bib>
  
```

```

RETURN $a
}
{
  if (count($b/author)>2)
  then (et-al/)
  else ()
}
}
</book>
}
  
```

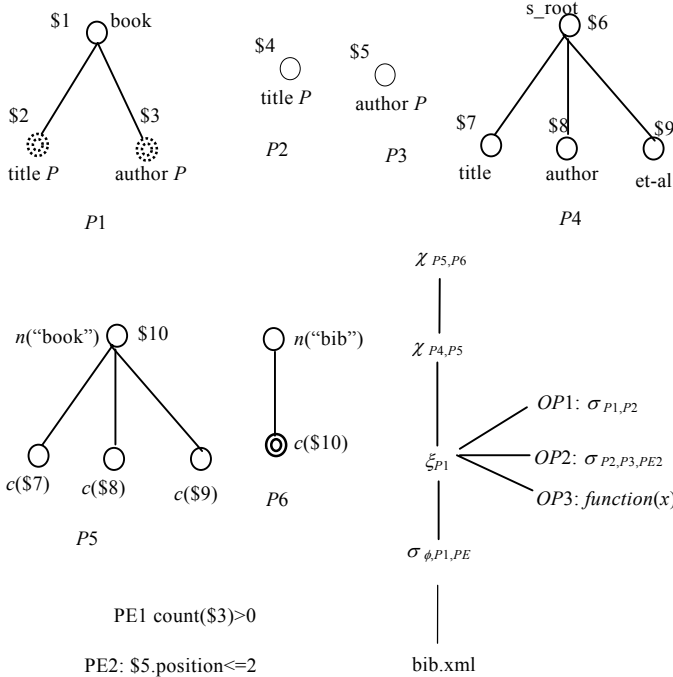


Fig.5 pattern trees and logical execution plan for W3C Use Case 1.1.9.6
 图5 W3C 查询用例 1.1.9.6 的模式树和逻辑执行计划

在这个例子中,有一个序列化操作符,注意,序列化操作符中的 3 个序列操作都是一元操作.

4 XML 代数新的优化问题

由于 XML 查询语言和 XML 数据的特点,XML 查询处理中出现了一些新的需要研究的优化问题.主要包括以下几方面.

4.1 谓词的作用域

在 XQuery 查询中,由于表达式可以任意嵌套,谓词可以出现在任意地方.谓词是有作用域的.同一个谓词在不同的地方,会产生不同的查询结果.这同时意味着原来关系代数优化中普遍使用的选择下推规则,在 XML 查询处理中可能会导致不正确的结果.比如以下查询:

```

Q5:
<bib>
{
for $b in doc("http://bstore1.example.com/bib.xml")//book,
where $b/publisher="Adison"
return
  <book-reviews>
  {
    { $b/title }
    for $v in doc("http://bstore2.example.com/reviews.xml")//review
      where $b/title=$v/title
      return
        <reviews>
        { $v/content }
        { $b/author[first="Smith"] }
      }
    }
  }
}
</bib>

```

在这个查询中,对于 bib.xml 上的结点,一共出现了两个选择谓词:\$b/publisher="Adison"和\$b/author/first="Smith"和一个连接谓词\$b/title=\$v/title.这里,两个选择谓词的作用域是不同的,谓词\$b/author/first="Smith"不能在连接之前做,否则结果会不正确.

在另一些情况下,谓词是可以被提到外层子句的.比如下面这个查询:

```

<bib>
{
for $b in doc("http://bstore1.example.com/bib.xml")//book,
where $b/publisher="Adison"
return
  <book-reviews>
  {
    { $b/title }
    for $v in doc("http://bstore2.example.com/reviews.xml")//review
      where $b/title=$v/title and $b/author/first="Smith"
      return
        <reviews> { $v/content } </reviews>
      }
    }
  }
}
</bib>

```

这时候,谓词\$b/author/first="Smith"就可以先于连接谓词进行判断,以使参与连接的集合变小.一般来说,如果一个内层谓词提到外层之后,该谓词所约束的子句是一样的,那么该谓词就可以提到外层,以提高查询效率.

如何快速地判断谓词的作用域并且判断谓词什么情况下可以提到外层(即选择下推),是 OrientXA 进行优化的一个方面.

4.2 更多的分组操作

XQuery 查询中,除了显式的分组操作,还有很多隐式的分组操作.这是由于 XQuery 的表达式可以任意嵌套.当存在相关嵌套子查询的时候,一般要转换为连接操作.这时候,为了构造期望的结果,就产生了分组的需要.比如以下查询:

```

<bib>
{
for $b in doc("bib.xml")//book,
where $b/publisher="Adison"
return
  <book-reviews>
  {
    { $b/title }
    for $v in doc("reviews.xml")//review
      where $b/title=$v/title
      return
        <reviews> { $v/content } </reviews>
      }
    }
  }
}
</bib>

```

一般来说,连接操作的形式是

$$\triangleright \triangleleft_C (X_1, X_2) = \{r(x_1, x_2) \mid x_1 \in X_1, x_2 \in X_2, C(x_1, x_2)\}.$$

其中, X_1 和 X_2 是参与连接的两个 XML 树的集合, C 是连接谓词. r 是虚的根结点,为了把连接结果以树的方式输出.

产生连接结果 $r(X_1, X_2)$ 之后, XQuery 的输出要求 X_2 的项嵌套在 X_1 的项中输出.这时候,就需要根据 X_1 .OID 进行一个分组操作,以构造期望的结果.

一个可能的避免这种分组操作的方法是引入一种新的连接操作:

$$\otimes_C (X_1, X_2) = \{r(x_1, (x_{21}, x_{22}, \dots, x_{2n})) \mid x_1 \in X_1, x_{2i} \in X_2, C(x_1, x_{2i})\}.$$

这种连接的意思是说:给定 x_1 ,把满足连接条件的所有 x_2 挂载在同一个虚根结点下.

在处理两层嵌套的查询中,这种连接确实可以避免上面的分组操作.但是,在3层或者以上的嵌套中,这种连接操作就不起作用了.

事实上,TAX 没有专门的结果构造符,它用 `group-by,copy-and-paste` 等操作符来完成结果构造的任务,这就反映出,在结果构造过程中,往往需要隐式的分组操作.

分组操作往往需要一次排序,是比较费时费力的操作.使用上述新的连接操作,可以避免两层嵌套的查询.如何避免或者减少由于结果构造带来的隐含的分组操作,也是 OrientXA 进行查询优化研究的一个方面.

4.3 更多的排序操作

XML 数据的一个特点是有序性.除非查询语句特别指定,否则应该保持结点在源文档中的结点顺序.这就带来一个重大的问题:原来一些处理连接的算法,比如排序连接、Hash 连接,会打乱结点顺序,在连接完成后,需要对连接结果做一个额外的根据结点顺序的排序操作.而如果用 `nest-loop` 连接算法,则可以省去这趟额外的排序.在进行代价估计的过程中,这个额外的排序操作要被考虑在内.这就给进行查询优化的过程带来新的考虑因素.

一个可能的解决方法是,在所有操作中忽略结点的有序性,在最后构造结果的时候,再对结点按照文档顺序排序.

究竟是使用 `nest-loop`,还是最后增加额外的排序的方法,这是 OrientXA 的查询优化的一个研究点.

5 结论和未来工作

本文研究 XQuery 的处理策略,提出了一套 OrientXA 查询代数系统.它主要借鉴了关系代数和现有的 XML 查询代数的思想,同时,又在结果构造符上体现了 XQuery 程序化语言的特点,因此能够有效地表达各种 XQuery 查询语言.本文还探讨了 OrientXA 代数优化的一些研究点.这些都是我们未来的工作之一.

References:

- [1] Galax XQuery engine. <http://db.bell-labs.com/galax>
- [2] IPSI-XQ XQuery engine. http://www.ipsi.fraunhofer.de/oasys/projects/ipsi-xq/index_e.html
- [3] Beeri C, Tzaban Y. SAL: An algebra for semistructured data and XML. In: Proc. of the Workshop on the Web and Databases. 1999. 37~42. <http://www-rocq.inria.fr/~cluet/WEBDB/procwebdb99.html>
- [4] Sartiani C, Albano A. Yet another query algebra for XML data. In: Proc. of the Workshop Iberoamericano de Ingenieria de Requisitosy Ambientes Software 2002. 2002. 106~115. <http://www.csd.uch.gr/~hy561/Papers/t2.2-r4b.pdf>
- [5] Colazzo D, Sartiani C, Albano A, Manghi P, Ghelli G, Lini L, Paoli M. A typed text retrieval query language for XML documents. Journal of the American Society for Information Science and Technology, 2002,53:467~488.
- [6] Frasinca F, Houben G, Pau C. XAL: An algebra for XML query optimization. In: Proc. of the Australasian Database Conf. 2002. 49~56.
- [7] Jagadish HV, Lakshmanan LVS, Srivastava D, Thompson K. Tax: A tree algebra for XML. In: Proc. of the 8th Biennial Workshop on Data Bases and Programming Languages. 2001. 149~164. <http://www.csd.uch.gr/~hy561/Papers/tax-full-dbp101.pdf>
- [8] Fernández M, Siméon J, Suciú D, Wadler P. A data model and algebra for XML query. Technical Report, 2000. <http://homepages.inf.ed.ac.uk/wadler/papers/xquery-algebra/xquery-algebra.html>
- [9] Fernandez M, Simeon J, Wadler P. An algebra for XML query. In: Proc. of the 20th Conf. on the Foundations of Software Technology and Theoretical Computer Science. 2000. 11~46. <http://homepages.inf.ed.ac.uk/wadler/papers/xalgebra-india/xalgebra-india.pdf>
- [10] Christophides V, Cluet S, Siméon J. On wrapping query languages and efficient XML integration. In: Proc. of the ACM SIGMOD Int'l Conf. on Management of Data. 2000. 141~152.
- [11] Pappas S, Wu Y, Lakshmanan LVS, Jagadish HV. Tree logical classes for efficient evaluation of XQuery. In: Proc. of the ACM SIGMOD Int'l Conf. on Management of Data. 2004. 71~82. <http://www.informatik.uni-trier.de/~ley/db/indices/a-tree/p/Pappas:Stelios.html>
- [12] Cluet S, Moerkotte G. Classification and optimization of nested queries in object bases. Technical Report, 1994. <http://citeseer.ist.psu.edu/cluet94classification.html>
- [13] Chen ZM, Jagadish HV, Lakshmanan LVS. From tree patterns to generalized tree patterns: On efficient evaluation of XQuery. In: Proc. of the 29th Int'l Conf. on Very Large Data Bases. 2003. 237~248. <http://www.informatik.uni-trier.de/~ley/db/indices/a-tree/c/Chen:Zhimin.html>