

# WEB 数据库 XML 数据发布及信息提取系统

赵美艳 王会进 张诗军

(暨南大学信息科学技术学院, 广东 广州 510632)

**摘要** 结合传统数据库和 XML 各自特点,在数据库到 XML 的数据转出、XML 数据发布、信息提取这三大模块的基础上,设计了一个系统原型.同时,对研究过程中提出的关系数据库到 XML 数据的“相关转出”、XML 转出数据的 WEB 服务端处理、XML 数据的提取等问题和关键解决思路、技术方法做出阐述.

**关键词** XML, 相关转出, 数据发布, 信息提取, WEB 服务

**中图分类号** TP 311.132.3

**文献标识码** A

数据库尤其是传统的关系数据库仍然是存储数据的主要手段,是 Internet 信息的一个主要来源.数据库中的数据必须通过适当的数据载体才能在 Internet 传播,而 XML 则可以说是目前最适合在 Internet 上使用的数据载体.因此,我们可以将两者结合起来,综合它们的优势,使 XML 成为数据库数据得以在 Internet 上传播的桥梁<sup>[1]</sup>.目前对关系数据库到 XML 数据的转出,许多数据库应用软件都提供了支持,但提供的只是以 XML 格式表达的简单查询结果集返回支持.从数据库中提取出来的数据只能面向某一特定显示或某一特定应用要求,往往不能被复用.同时,只支持对某个表的数据的转出而不支持多个表关系的同时转出.因此在本文中,除了提出了 SQL 直接查询支持外,还提出了按照应用表集实现基于表集的“相关转出”支持.从而,充分发挥了 XML 的应用优势.另一方面,数据库转出的 XML 格式的数据不能单纯地认为只是将 XML 数据拿来,然后原封不动的发送给请求者.在本文中提出了 XML 数据发布和信息提取的观点.它使得转出的 XML 数据在最大程度上得到复用,同时满足用户根据自己的需求提取相关信息<sup>[2]</sup>.

## 1 系统的运行架构及其实现思想

### 1.1 系统的运行架构图

系统的运行如图 1 所示. WEB 服务器收到 WEB 请求后进行判断.如果请求是 XML 静态文档,则从 WEB 文档发布目录中直接得到相应 XML 文档进行发布.如果需要运行扩展服务程序,则通过扩展程序的接口来运行相应的外部扩展服务程序,得到结果后响应 WEB 请求.扩展服务程序运行时,如果需要直接通过数据库访问得到 XML 数据或查询结果集,则通过数据库访问接口(如 JDBC)访问数据库.如果使用“相关转出”的 XML 数据,则访问 WEB 文档发布目录中的相关 XML 文档并判断该文档是否有效.无效则调用 WEB 文档生成功能来刷新 XML 文档,得到相关 XML 文档后可以经过相应的处理,再将处理结果返回 WEB 服务器.如果是动态的“相关转出”,XML 数据访问则由外部扩展服务程序直接调用 WEB 文档生成功能,得到动态的 XML 数据. WEB 文档生成功能从表关系保存文件目录中得到表关系保存文件,据此生成 XML 文档并放入 WEB 文档发布目录中或实时返回给 XML 数据请求程序.表关系保存文件由表关系保存文件生成功能根据数据库中已经定义的表关系生成.

### 1.2 系统设计的特点

(1) 提出“相关转出”功能的支持.提供了多个表数据及它们之间关系的转出. (2) WEB 服务端扩



```
< ! ELEMENT Parents( # PCDATA)>
< ! ELEMENT Sons( # PCDATA)>
< ! ELEMENT MapInclude( # PCDATA)>
< ! ELEMENT MapIdref( # PCDATA)>
```

其中 MapInclude 对应于被依赖主键所在表只存在当前表外键对它的依赖, 没有其它表的外键依赖它的主键的情况. 多个则需要使用 ID/IDREF 表示的表关系元素 MapIdref, 即被依赖主键所在表除了当前外键对它的依赖外, 还存在其它表的外键依赖于它的主键.

2.4 相关转出 XML 文件生成程序的主要执行步骤

(1) 读取表关系保存文件, 进行初始化. (a) 依据 DataSource 元素中保存的数据库连接描述信息连接数据库. (b) 针对每个表的相关描述信息生成描述对象集 A. 根据需要可以由 A 生成相关的 DTD 或 Schema. (c) 形成 XML 文件的数据框架 S, 主要包括 XML 文件头和根元素名. (2) 遍历表对象集 A, 依次搜索到未被搜索过的没有被任何外键依赖或被多个外键依赖的描述对象 Object(即 Object 的 Parents 属性无值或有多个值). (3) 根据 Object 属性值(主要是 MapInclude 和 MapIdref)找到和 Object 有直接表关系的表, 并根据描述信息生成并执行相关的数据库查询语句. (4) 根据数据库查询结果集, 生成 Object 的 XML 转出数据. (5) 由 Object 的 MapInclude 属性值得到表集 B. 遍历 B, 依次搜索到 B 中未被搜索过的表描述对象 Object, 转向(3)步处理 Object. (6) 由 S 生成 XML 转出文档.

3 结束语

本文从数据库到 XML 的数据转出和 XML 数据发布及信息提取两个方面着手. 分析讨论构建 XML 的 WEB 数据库数据发布及信息提取系统的一些主要问题及相关解决思想. 其中, 对于关系数据库到 XML 的数据转出问题的解决, 除了提供直接查询生成支持外还提出了“相关转出”的概念和实现思路, 在转出时兼顾了“整体性”和“局部性”问题. 对于 XML 数据的发布处理问题的解决, 采用 Java Servlet 技术, WEB 服务器采用 Apache+ Tomcat 技术.

参 考 文 献

1 Seligman L, Rosenthal A. XML's Impact on DataBases and DataSharing[J]. IEEE Computer, 37(6): 59~ 67  
2 Ferandez M, Chiew W. Trading between relations and XML[J]. Computer Networks, 2000, 33: 723~ 745  
3 Jasnowski M 著. Java, XML 和 WEB 服务宝典[M]. 盖江南等译. 北京: 电子工业出版社, 2002. 34~ 63

Data Publishing and Information Extracting System of  
XML-Based WEB Database

Zhao Meiyan    Wang Huijin    Zhang Shijun

( College of Info. Sci. & Tech., Jinan Univ., 510632, Guangzhou, China)

**Abstract** Based on data roll out from database to ZML, XML data publishing and information extracting as three modules, a system prototype is designed. The design is also integrated with the respective characteristic of traditional database and XML. Explanations are made on such problems as “correlated roll out” from relational database to XML data, processing of XML roll out data at WEB server, and extraction of XML data; and also on the thinking of solving and the technology.

**Keywords** XML, correlated roll out, data publishing, information extracting, WEB server