

# 多数据库查询接口的设计与实现\*

陈微 刘瑞 李昭原

(北京航空航天大学信息工程研究所 北京 100083)

**摘要** MQI(multi-database query interface)是作者研制的、用SQL对多个RDBMS(relational database management system)数据源进行综合查询的客户端工具。该文主要讨论其逻辑结构与算法设计。

**关键词** 多数据库,综合查询。

**中图法分类号** TP311

随着计算机技术的不断进步和应用领域的不断拓展,现代管理信息系统的数据构成呈现出多元和分布的趋势,即一个应用涉及到的数据不仅存在于不同的地点,还可能由不同的数据库管理系统(DBMS)进行管理,这就导致了多数据库系统的产生和发展。与单一的DBMS相比,多数据库系统具有更好的开放性、可扩充性和灵活性。它能够针对不同的领域、不同的部门、不同的需求,使用不同的数据库技术,从而能够提供更好的性能/价格比。然而,多数据库系统也给信息系统的设计和实现带来了更大的复杂性。其中,如何实现对多数据库的查询操作是一个亟待解决的问题。与对单一DBMS的查询相比,对多数据库的查询操作必须隐藏各个DBMS之间的差异和它们在系统中的位置,使用户能够通过一个统一的查询界面,透明地访问所需的数据。

鉴于当前大多数多数据库系统均由关系数据库构成,本文提出和设计了一个适用于由多种RDBMS(relational database management system)组成的多数据库综合查询接口软件MQI(multi-database query interface)。它在不同的RDBMS之上提供了一个新的抽象层面,该抽象层面向下可映射到各种具体的RDBMS,向上(对应用)则隐藏各RDBMS的差异,从而构成了一个虚拟的RDBMS。MQI以SQL作为查询语言,能够支持对分布在多个数据库系统中的数据的查询,具有较好的透明性、可扩充性和容错性。

## 1 MQI的应用环境

目前,市场上支持多数据库系统的商品软件有Ingres/Star和Sybase OmniSQL Gateway等,但它们都是服务器一级的软件,对硬件要求高,安装和维护工作量大,且价格昂贵。鉴于此,MQI的设计原则是小巧和易于使用,它在基于客户/服务器的信息系统中的位置如图1所示。这种结构有两个特点:(1)MQI驻留在客户端,为一个客户机上的应用服务,避免产生新的瓶颈问题,从而降低了对硬件的要求。(2)MQI位于负责网络事务的中间件之上,与具体网络无关。

## 2 MQI的逻辑结构

MQI的模式设计主要参考了Amit P. Sheth和James A. Larson在文献[1]中给出的描述FDBS的5级参考模型。根据实际需要,我们略去了成员模式,并将外模式与联邦模式合二为一,如图2所示。

Component DBS:每种DBMS可看作一个Component DBS。在MQI中,Component DBS的种类可以不一样,但必须是关系的,Local Schema:各成员数据库的概念模式,Export Schema:Local Schema的子集。针对不同的MQI用户,每个成员数据库可有不同的Export Schema。External-Federated Schema:是综合了一个MQI用户所有的Export Schema得到的,它刻画了一个MQI用户所能利用的所有数据的结构。

通过这3层模式,MQI实现了对多RDBMS的综合查询。下面给出加入处理过程的系统结构图如图3所示。

\* 本文研究得到国家航空基金和航空预研项目基金资助。作者陈微,女,1970年生,博士生,主要研究领域为数据库,软件工程。刘瑞,1970年生,讲师,主要研究领域为数据库理论,数据库应用开发工具。李昭原,1938年生,教授,主要研究领域为数据库工程,软件工程,管理工程。

本文通讯联系人:刘瑞,北京100083,北京航空航天大学信息工程研究所

本文1997-05-16收到原稿,1997-06-23收到修改稿

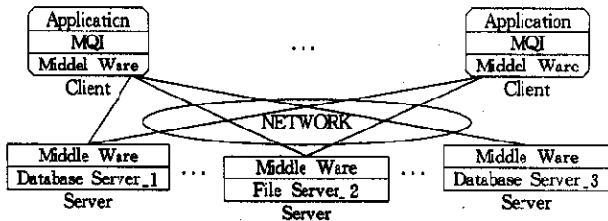


图1

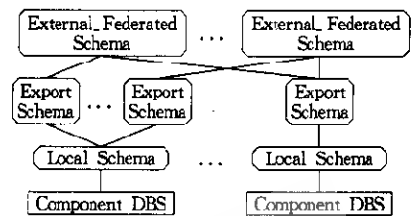


图2

**Filter Processor:** 允许合法的命令和数据通过。它对命令作语法和存取控制检查,对数据作类型转换。

**Constructing Processor:** 通过分割或复制,把一个任务分解为多个任务分别去完成,然后将每个任务产生的数据综合为一个数据集返回给调用者。它通过内部数据字典提供位置和数据源类型的透明性。它主要完成下面的功能:(1) 模式集成;将多个输出模式集成为一个外部-联邦模式。(2) 查询分解与优化;分解和优化建立在外部-联邦模式上的查询。

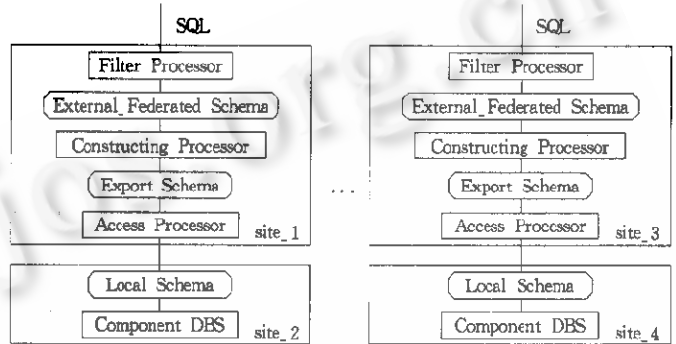


图3

**Access Processor:** 接受命令,按命令要求从数据库中检索出相应的数据。**Construction Processor** 是整个系统实现的关键部分,下一节将详细讨论查询任务分解的算法。

### 3 MQI 的查询分解算法

#### (1) 求出查询的基本数据集合的 SELECT 表达式

关系代数和元组演算是 SQL 查询的数学基础,是证明各种分解算法正确性的手段,但关系代数却不一定能表达 SQL 表示的查询。这是因为 SQL 融进了很多关系代数所没有的特征,如分组、排序、聚集等。为便于算法的设计与证明,我们提出基本数据集合的概念:给定 SELECT 语句的基本数据集合是能关系代数表达,且能产生 SELECT 的最终结果的最小关系。这里,最小是指,去掉任何一行或一列都不足以得到最终结果。每个基本数据集合总能由一个 SELECT 语句表达(因为能用关系代数表达的总能用 SELECT 表达),这个 SELECT 语句称为原 SELECT 语句的基本数据集合的 SELECT 表达。

查询分解的第 1 步就是去掉 SELECT 语句中那些不能由关系代数表达的特征,得到与之对应的基本数据集合的 SELECT 表达。在 MQI 支持的形如

```
SELECT select-list
FROM table-list
WHERE where-list GROUP BY group-list HAVING having-list
ORDER BY order-list
```

的 SELECT 语句中,不能由基本关系代数运算表达的特征有:列别名、表达式、聚集函数、分组、排序。下面逐步除去这些特征,得到基本数据集合的 SELECT 表达:

- ① 去掉 select-list 中的聚集函数,保留作为被去掉的聚集函数的参数的列名(如果和已有列不重复的话)。
- ② 去掉出现在 select-list 中的表达式并将表达式中的列名添到 select-list 中(如果和已有列不重复的话)。
- ③ 去掉出现在 select-list 中的别名。
- ④ 将出现在 where-list 中,却不在 select-list 中的列加到 select-list 中。
- ⑤ 对 group-list, having-list, order-list 重复④中对 where-list 的操作,并去掉相应的子句。

#### (2) 化一般查询为用 UNION 连接的多个合取查询

合取查询是 WHERE 句中的条件只包含 AND 和 NOT 两种连接符,而且每个 NOT 作用域内部都没有 AND 的

查询. J. D. Ullman 在文献[2]中指出,合取查询是一类很自然的查询,它包括了关系代数“选投连”核心的部分. 实际中,大多数查询都属于合取查询. 合取查询的等价性问题是可判定的,所以,大多数系统在设计其查询分解优化的算法时,主要工作都在对合取查询的优化上. MQI 的分解优化策略也主要针对合取查询.

当某个查询不是合取查询时,我们将把它变为多个合取查询来处理. 根据德·摩根定理,每个布尔表达式都可以规范为合取式的析取形式,即具有  $A_1 \text{ OR } \dots \text{ OR } A_n$  (其中  $A_i (1 \leq i \leq n)$  是一个合取式), 则任何查询的基本数据集的 SELECT 表达都可表示为

```
SELECT select-list FROM table-list WHERE  $A_1 \text{ OR } \dots \text{ OR } A_n$ 
```

可以证明,它与下述语句等价.<sup>[3]</sup>

```
SELECT select-list FROM table-list WHERE  $A_1$ 
UNION .....
UNION
SELECT select-list FROM table-list WHERE  $A_n$ 
```

然后对每个已成为合取查询的 SELECT 语句分别处理,最后将其合并为一个关系.

(3) 分割每一个合取查询

分割合取查询的目的是把涉及到多个 DBMS 的数据的查询分解为一组能在单个 DBMS 上完成的子任务和综合这些子任务的结果的子任务.

①合取查询中“选投连”的可重组性

合取查询中“选投连”的可重组性是指 SELECT 语句的这样一种特性: 设有合取查询

```
SELECT select-list FROM from-list WHERE where-list
```

其中,所有在 where-list 中出现的列都出现在 select-list 中. 现将 from-list 分为  $r$  组,相应地,把属于 from-list- $i$  的列分到 select-list- $i$  组,把只与 select-list- $i$  中的列相关的 where-list 中的条件都分在 where-list- $i$  中,将 where-list 中不属于任何 where-list- $i$  中的条件分到 where-list-0 中,忽略属性在元组中出现的先后次序,原 SELECT 语句改写为

```
SELECT select-list-1,select-list-2,...,select-list-r
FROM from-list-1,from-list-2,...,from-list-r
WHERE where-list-1 and ... and where-list-r and where-list-0
```

其中 where-list- $i$  可能为空,则有如下结论.

$$\text{SELECT select-list-}i \text{ FROM from-list-}i \text{ WHERE where-list-}i \text{ 结果存于 } T_i \tag{1}$$

$$\dots$$

$$\text{SELECT select-list-}r \text{ FROM from-list-}r \text{ WHERE where-list-}r \text{ 结果存于 } T_r \tag{r}$$

$$\text{SELECT select-list FROM } T_1, T_2, \dots, T_r \text{ WHERE where-list-0} \tag{r+1}$$

{1}~{r+1}步执行的结果与原 SELECT 语句结果相同.<sup>[3]</sup>

②查询分解中的分割原则

合取查询中“选投连”的可重组性解决了子任务的执行结果是否与原来的查询结果等价的问题. 分组时,应遵循两个原则:①每个分组只能包含存在于同一数据源上的表/视图;②每个分组的查询的连接图必须是连通的.

连接图用于描述查询中涉及到的表/视图之间的连接关系. 连接图的顶点是 table-list 中的每个表/视图,如果两个表/视图有连接关系,则在对应二顶点间画一条边. 语句

```
SELECT C.attrc,A.attra
FROM A,B,C,D,E,F
```

```
WHERE A.ab=B.ba and B.bc=C.cb and C.ce=E.ec and D.de=E.ed and D.df=F.fd and F.fe=E.ef
```

的连接图如图 4(a) 所示.

假设 A,B,D,E,F 在数据源  $\alpha$  中,C 在数据源  $\beta$  中,如果以是否在同一数据源上为分组的唯一原则,无疑,A,B,D,E,F 在同一组,而 C 在另一组. 在查询被分解的同时,原查询图也被分为两个子图(图 4(b),(c)). 图 4(b)是在数据源  $\alpha$  上执行的查询的连接图. 由于这个图不是连通的,所以必然产生笛卡尔积,而笛卡尔积能使结果数据成几何速度增长,这是查询中最糟糕的情况. 在一般情况下,应避免使用包含笛卡尔积(查询连接图不为连通图)的查询. 在例子中,用户给出的查询连接图是连通的,只是由于使用了不合适的分割方法而被分割成不连通的子图,从而使查询效率

大受影响。如果将 A, B, D, E, F 进一步划分为 A, B 和 D, E, F 两个组, 则可避免上述问题。即使在用户给出的查询连接图不连通的情况下, 也要尽量推迟笛卡尔积的操作, 以减少网络通讯量和对其后操作的影响。通过以上分析, 我们给出分组的第 2 条原则。

③分割步骤

我们用一个四元组 (C, T, S, D) 来描述每一个分割出的组, 其中 C 是属于该组的需要输出的列的集合, T 是该组的表/视图的集合, S 是该组选取和连接条件的集合, D 是临时结果库名。对一个合取查询分割的结果, 就是产生若干这样的四元组, 再由四元组组装 SELECT 语句。不失一般性, 设合取查询有如下形式。

SELECT a<sub>1</sub>, ..., a<sub>n</sub> FROM t<sub>1</sub>, ..., t<sub>m</sub> WHERE B<sub>1</sub> and ... and B<sub>r</sub>

分割过程如下:

① 创建 P<sub>0</sub>(C, T, S, D), 其中 C = {a<sub>1</sub>, ..., a<sub>n</sub>}, 创建 P<sub>i</sub>(C<sub>i</sub>, T<sub>i</sub>, S<sub>i</sub>, D<sub>i</sub>) (1 ≤ i ≤ m), 其中 P<sub>i</sub>. T = {T<sub>i</sub>}, P<sub>i</sub>. C = {a<sub>1</sub>, ..., a<sub>n</sub> 中属于 t<sub>i</sub> 的列}, P<sub>i</sub>. S = {B<sub>1</sub>, ..., B<sub>r</sub> 中只包含了 P<sub>i</sub>. C 中元素的布尔表达式}。

② 对于每个不属于任何 P<sub>i</sub> 的 B<sub>k</sub>, 设 A, B 是该 B<sub>k</sub> 涉及到的两个表/视图, 若 A, B 在同一数据源上, 设 A 所在的分割是 P<sub>i</sub>, B 所在的分割是 P<sub>j</sub>, 且 i < j, 则将 P<sub>j</sub> 中的内容添加到 P<sub>i</sub> 中去, 即令 P<sub>i</sub>. C = P<sub>i</sub>. C ∪ P<sub>j</sub>. C; P<sub>i</sub>. T = P<sub>i</sub>. T ∪ P<sub>j</sub>. T; P<sub>i</sub>. S = P<sub>i</sub>. S ∪ P<sub>j</sub>. S ∪ {B<sub>k</sub>}, 然使 P<sub>j</sub> 无效, 否则, 将 B<sub>k</sub> 加到 P<sub>0</sub> 中, 即 P<sub>0</sub>. S = P<sub>0</sub>. S ∪ {B<sub>k</sub>}。

③ 过程②的结果使有的 P<sub>i</sub> 无效了。对每一个有效的 P<sub>i</sub> 分配一个唯一的表名 TN<sub>i</sub>, 即使 P<sub>i</sub>. D = TN<sub>i</sub>, 并使

$$P_0.T = \bigcup_{1 \leq i \leq m, P_i \text{有效}} \{TN_i\}.$$

④ 为每个有效的 P<sub>i</sub> 生成相应的 SQL 语句, 即

SELECT P<sub>i</sub>. C FROM P<sub>i</sub>. T WHERE P<sub>i</sub>. S

⑤ 还原在①中被去掉的特征

第 3 步所产生的 SQL 序列只是查出了基本数据集合, 要得到用户想要的结果, 还要加上被去掉的如列别名、GROUP、ORDER 子句等特征, 具体方法是

SELECT select-list FROM P<sub>0</sub>. D GROUP By group-list  
HAVING having-list ORDER BY order-list

其中 select-list, group-list, having-list, order-list 和用户所给出的相同。

4 小 结

MQI 是国家航空基金和国防科工委基金的资助项目——“开放-分布-异构和基于客户/服务器结构的数据库应用快速开发平台”的子项目。它支持对多数据库的综合查询, 既可嵌入在“开发平台”中使用, 也可独立使用。目前, 我们已经在 Windows 环境下实现了 MQI 综合查询系统。今后, 我们将作更进一步的理论探索和实践, 使 MQI 的算法更加完善, 从而提高效率和支持更多的 SQL 特征。

参考文献

- 1 Sheth Amit P, LARSON James A. Federated database system for managing distributed, heterogeneous, and autonomous databases. ACM Computing Surveys, 1990, 22(3): 183~236
- 2 Jeffrey D Ullman 著, 张作民译. 数据库系统原理. 北京, 国防工业出版社, 1984  
(Jeffery D Ullman. Principles of Database Management System. Rockville, MD: Computer Science Press, 1980)
- 3 陈徽. 多数据库查询接口的设计与实现 [硕士学位论文]. 北京航空航天大学, 1995  
(Chen Wei. Design and implementation of a multi-database query interface [M. S. Thesis]. Beijing University of Aeronautics and Astronautics, 1995)

Design and Implementation of a Multi-database Query Interface

CHEN Wei LIU Rui LI Zhao-yuan

(Information Engineering Institute Beijing University of Aeronautics and Astronautics Beijing 100083)

**Abstract** MQI (multi-database query interface) is a client tool which the authors developed for querying with SQL on several different RDBMS (relational database management system) data sources. Its logical architecture and algorithm design are mainly discussed in this paper.

**Key words** Multi-database, synthetic query.