

# 偶查询优化的图论方法

OPTIMIZATION OF BIPARTITE JOINS WITH GRAPH THEORY

曾永宁 童 源

Zeng Yongning Tong Fu

(计算机系)

**摘 要** 本文针对一类特殊的多关系查询——偶查询,提出了一种建立在图论、偶图和匹配理论基础上的查询优化方法,这种方法具有多项式复杂性。

**关键词** 数据库; 查询处理; 查询优化; 算法; 图论; 偶图; 配对

**ABSTRACT** An optimization approach for multi-relational joins based on the matching concept in graph theory is introduced. The basic idea lies in constructing a join graph from an expression of N-relation query, and seeking for a maximum matching with minimum total weight. The algorithm for bipartition(X, Y) is proposed in this paper.

**KEY WORDS** database; algorithm; graph theory; bipartite graph; matching; query processing; query optimization

## 一、引 言

关系数据库的查询处理效率极大地影响数据库的系统性能。因此,查询优化问题受到了广泛的重视,不少人对它进行了深入广泛的研究〔2、3、4、5、6、7、8、9、10、11、13〕。其研究结果表明查询优化问题是NP-完全的〔1〕。为了寻找有效的优化方法,通常把查询分为树查询和环查询两类。所谓树查询是指这类查询所对应的查询图是树;而环查询所对应的查询图中包含着回路。对树查询已经找到了具有多项式复杂性的优化方法,但对环查询还没有找到有效的优化方法。

本文从图论的角度出发,把查询分为偶查询和非偶查询两类。所谓偶查询是指其查询所对应的查询图是偶图(Bipartite Graph)。不难证明树查询属于偶查询类,而且一些环查询也属偶查询类。我们应用图论的匹配(Matching)理论,导出了一个有效的偶查询优化算法。本文的第二节介绍图论有关的概念,给出了若干新的定义;第三节讨论估算多关系查询

本文于1988年1月23日收到。

处理开销的方法；第四节给出了偶查询优化算法。

### 二、预 备 知 识

**定义2.1:** 设 $D = (R_1, R_2, \dots, R_n)$ 为数据库模式,  $Q = (q = q_1 \wedge q_2 \wedge \dots \wedge q_n)$ 为查询表达式,  $q$ 为查询约束条件. 查询图 $G_q$ 是一个有序三元组 $(V_q, E_q, \varphi_q)$ . 其中 $V_q = \{R_1, R_2, \dots, R_n\}$ 为顶点集.  $E_q$ 代表边集,  $\varphi_q$ 是映射函数, 它确定了边集 $E_q$ 和顶点集 $V_q$ 间的映射关系:  $\varphi_q(e) = R_i R_j$ 当且仅当 $R_i$ 和 $R_j$ 在 $Q$ 中是可联接的. 例如, 设 $Q = (R_1 \cdot A = R_2 \cdot B) \wedge (R_1 \cdot C = R_3 \cdot D) \wedge (R_2 \cdot E = R_4 \cdot F) \wedge (R_2 \cdot G = R_5 \cdot H)$ 为一查询表达式, 它的查询图如图2.1所示.

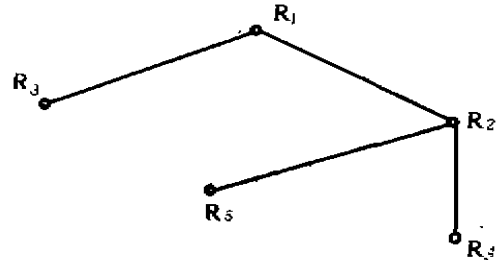


图 1 查 询 图

**定义2.2:** 联接图 $G_j$ 是一个带权查询图. 即 $G_j = (V_q, E_q, \varphi_q, W_q)$ 其中 $W_q$ 是边的权集:

$$W_q = \{w_{ij} | \text{对所有的边 } R_i R_j \in E_q\}$$

其中 $w_{ij}$ 是边 $R_i R_j$ 所带的权.

**定义2.3**[6]: 设 $M$ 为边集 $E_q$ 的一个子集, 如果 $M$ 满足没有任何两条 $M$ 中的边是相邻的. 则称 $M$ 为 $E_q$ 的一个匹配 (Matching), 并称 $M$ 中任一边所关联的两个顶点在 $M$ 下是配对的.

若匹配 $M$ 的某条边与顶点 $R_i$ 关联, 则称 $M$ 饱和顶点 $R_i$ , 并称 $R_i$ 是 $M$ 饱和的, 否则 $R_i$ 是 $M$ 非饱和的. 若 $G_j$ 的每个顶点为 $M$ 饱和的, 则称 $M$ 为 $G_j$ 的完全匹配. 若 $G_j$ 没有另外的匹配 $M'$ , 使 $|M'| > |M|$ 成立, 则称 $M$ 为 $G_j$ 的最大匹配.

设 $M$ 是 $G_j$ 的匹配,  $G_j$ 的 $M$ 交错路径是指其边在 $E_q \setminus M$ 和 $M$ 中交错出现的路径.  $M$ 可扩路径是指其起点和终点都是 $M$ 非饱和的 $M$ 交错路径.

**定理2.1:**  $G_j$ 的匹配 $M$ 是最大匹配, 当且仅当 $G_j$ 不含 $M$ 可扩路径.

**定义2.4**[6]: 若图 $G_j$ 的顶点可划分为两个子集 $X$ 和 $Y$ , 使得 $G_j$ 的任一条边的两个顶点总是分属 $X$ 和 $Y$ , 则称图 $G_j$ 为偶图 (Bipartite Graph), 而 $(X, Y)$ 称为 $G_j$ 的一个二分类.

完全偶图是一个具有二分类 $(X, Y)$ 的简单偶图, 并且 $X$ 的任一个顶点总是和 $Y$ 的每一个顶点相关联.

**定义2.5:** 如果一个查询 $Q$ 所对应的查询图 $G_j$ 是一个具有二分类 $(X, Y)$ 的偶图, 则称 $Q$ 为偶查询.

树查询属于偶查询类, 因为根据图论, 树总可转换为偶图. 有一些环查询, 也可以转换为偶查询. 例如, 对查询图为如图2.2(a)所示的一个环查询, 不难把它转换成其查询图为如图2.2(b)所示的偶查询.

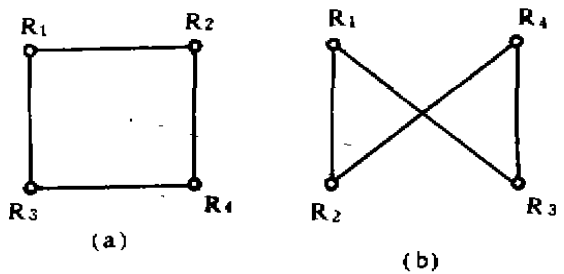


图 2 环 查 询 向 偶 查 询 的 转 换

$$\begin{aligned} X &= \{R_1, R_4\} \\ Y &= \{R_2, R_3\} \end{aligned}$$

**定义2.6:** 设 $S$ 为联接图 $G_j$ 的顶点集 $V_q$ 的一个子集.  $S$ 的邻集 $V_n(s)$ 是由所有 $S$ 相邻的顶点构

成的顶点集。

若 $G_i$ 是一个具有二分类 $(X, Y)$ 的偶图, 我们希望找到一个能够饱和 $X$ 的每一个顶点的匹配, 存在这样的匹配的充要条件由下面的定理给出。

**定理2.2:** 设 $G_i$ 是具有二分类 $(X, Y)$ 的偶图,  $G_i$ 包含饱和 $X$ 的每一个顶点的匹配的充要条件是

$$|V_n(S)| \geq |S| \quad \text{对所有的 } S \subset X$$

### 三、查询处理开销评估

用 $\phi(R_i, A)$ 表示关系 $R_i$ 的属性 $A$ 的值域。 $R_i[A]$ 表示关系 $R_i$ 在属性 $A$ 上的投影。 $R_i[A]$ 是 $\phi(R_i, A)$ 的一个随机子集。所有值域的集合构成了一个值域层次结构。每一个层次由一个最大值域 $\phi_{max}$ 和若干值域 $\phi$ 组成, 满足 $\phi \subset \phi_{max}$ 。如果值域 $\phi_i$ 和 $\phi_j$ 同属一个层次, 则称 $\phi_i$ 和 $\phi_j$ 是可联接的。

下面给出三个假设作为讨论的前提。

(1) 如果 $\phi_i \subset \phi_j$ , 则 $\phi_i$ 是 $\phi_j$ 的随机子集;

(2) 关系 $R_i$ 的元组在 $R_i[A]$ 上是均匀分布的;

(3) 用 $\text{attr}(R_i)$ 表示 $R_i$ 的属性集, 则对关系 $R_i$ 和属性 $\{A, B\} \subset \text{attr}(R_i)$ ,  $R_i[A]$ 和 $R_i[B]$ 是相互独立的。

由以上假设, 可导出任一元素 $x \in \phi_{max}$ 满足 $x \in R_i[A]$ 的概率 $p(R_i, A)$ 为:

$$p(R_i, A) = \eta(R_i[A]) / \eta(\phi_{max}) \quad (1)$$

其中,  $\eta(R_i[A])$ 表示 $R_i[A]$ 的不同值的数目。 $\eta(\phi_{max})$ 表示 $\phi_{max}$ 的不同值的数目。

一般地, 多关系查询的约束条件可表示为 $q = q_1 \wedge \dots \wedge q_m$ 。 $q_k (k = 1, 2, \dots, m)$ 的形式为:  $R_i \cdot A \theta v$ , 或 $R_i \cdot A \theta R_j \cdot B$ ,  $\theta \in \{=, <, \geq, <!, >!\}$ ,  $v$ 是一个算术表达式,  $R_i \cdot A \theta v$ 是单变量约束, 不少文献已对它作了深入的讨论[12]。本文只讨论联接约束 $R_i \cdot A \theta R_j \cdot B$ 。

多关系查询处理的开销主要取决于在查询处理过程中访问存储页面的总数。一旦数据库建立起来后, 各个关系所占用的存储页面总是可计算的。设 $N_i$ ,  $t_i$ 分别为关系的元组数和元组长度。 $P_i$ 为关系 $R_i$ 占用的存储页面。则:

$$P_i = N_i t_i / b \quad (2)$$

$b$ 表示页面大小。

另一方面, 在多关系查询处理过程中, 要形成若干中间临时关系, 这些临时关系所占用的总页数随着选择执行二元联接的次序的不同而变化。每一个约束子条件为 $R_i \cdot A \theta R_j \cdot B$ 的子查询 $Q'$ 都要形成一个临时关系 $R_{i,j}$ 。处理 $Q'$ 的联接开销 $C_{i,j}$ 取决于关系 $R_i$ ,  $R_j$ 的大小和联接算法。然而, 当我们从一个多关系联接表达式确定二元联接的次序时, 不仅要考虑二元联接的开销 $C_{i,j}$ , 而且也要考虑形成的中间临时关系 $R_{i,j}$ 的大小。因为后者对多关系查询处理的总开销有很大的影响。因此, 我们把联接图 $G$ 各边的权 $W_{i,j}$ 定义为:

$$W_{i,j} = C_{i,j} + P_{i,j} \quad (3)$$

而

$$P_{i,j} = N_{i,j} t_{i,j} / b \quad (4)$$

其中 $N_{i,j}$ 是 $R_{i,j}$ 的元组数, 它是未知的。下面我们来讨论 $N_{i,j}$ 的计算方法。

首先定义 $R_{i,j}$ 的选择因子。

**定义3.1:** 临时关系 $R_{i,j}$ 的选择因子 $p_{i,j}$ 是由关系 $R_i$ 和 $R_j$ 的元组构成的元组对满足约束子条件 $q_{i,j} = R_i \cdot A \theta R_j \cdot B$ 的概率。

因此只要 $p_{i,j}$ 已知,  $N_{i,j}$ 容易由下述公式求得:

$$N_{i,j} = p_{i,j} N_i N_j \quad (5)$$

在相等联接 (equijoin) 的情况下,  $p_{i,j} = p(R_i \cdot A) \cdot p(R_j \cdot B)$ 。这时(5)式变为:

$$N_{i,j} = p(R_i \cdot A) p(R_j \cdot B) N_i N_j \quad (6)$$

为了在多关系查询处理过程中, 重新计算压缩后的联接图各边的数, 需要估算联接属性 $R_{i,j} \cdot AB$ 的不同值的数目 $\eta(R_{i,j}[AB])$ , 以及 $R_{i,j}$ 的其他属性 $R_{i,j} \cdot A'$ 的不同值的数目 $\eta(R_{i,j}[A'])$ 。

$$\eta(R_{i,j}[AB]) = p(R_i \cdot A) p(R_j \cdot B) \eta(\phi_{max}) \quad (7)$$

然而, 估算 $\eta(R_{i,j}[A'])$ 较困难。这里我们利用文献[5]提出的一个近似公式:

$$\eta(R_{i,j}[A']) = \begin{cases} N_{i,j} & \text{当 } N_{i,j} < \frac{1}{2} \eta'(R_{i,j}[A']) \\ 1/3 (N_{i,j} \cdot \eta'(R_{i,j}[A'])) & \text{当 } \frac{1}{2} \eta'(R_{i,j}[A']) < N_{i,j} < 2\eta'(R_{i,j}[A']) \\ \eta'(R_{i,j}[A']) & \text{当 } 2\eta'(R_{i,j}[A']) < N_{i,j} \end{cases} \quad (8)$$

其中

$$\eta'(R_{i,j}[A']) = \begin{cases} \eta(R_i \cdot A') & \text{当 } R_{i,j} \cdot A' = R_i \cdot A' \\ \eta(R_j \cdot A') & \text{当 } R_{i,j} \cdot A' = R_j \cdot A' \end{cases} \quad (9)$$

#### 四、偶查询优化算法

Kuhn-Munkres 曾提出了一个有效的算法[6], 在一个具有二分类 $(X, Y)$ ,  $X = \{x_1, x_2, \dots, x_n\}$ ,  $Y = \{y_1, y_2, \dots, y_n\}$ 的带权完全偶图上寻找具有最大权的完美匹配。我们汲取了Kuhn-Munkres 算法的基本思想, 提出了在具有二分类 $(X^*, Y^*)$ ,  $X^* = \{Rx_1, Rx_2, \dots, Rx_n\}$ ,  $Y^* = \{Ry_1, Ry_2, \dots, Ry_r\}$ 的带权偶图 $G_i$ 上寻找具有最小权的最大匹配的算法, 从而导出了偶查询优化算法。该算法从以下三点对Kuhn-Munkres 算法作了修改:

- (1) 取消了带权偶图必须是完全图的限制;
- (2) 我们的算法是求具有最小权而不是最大权的最大匹配;
- (3) 放宽了对二分类的限制。 $X$ 的元素数目可以不等于 $Y$ 的元素数目。即 $k \leq r$ 。

##### 1. 偶查询优化算法的基本思想

首先从查询表达式构造联接图 $G_i$ 。并根据式(1)~(9)计算各边的权。然后在图上寻找一个具有最小权的最大匹配 $M$ 。 $M$ 中的边代表了一组相互独立的二元联接操作。然后, 把 $M$ 中的边压缩为一个顶点,  $G_i$ 变为一个新的联接图 $G'_i$ ,  $G'_i$ 的边数和顶点数都要比 $G_i$ 少。重新计算各边的权。重复上述过程, 直到 $G_i$ 变成一个平凡图(只有一个顶点的图)。至此我们得到了一个二元联接的最优次序, 算法终止。

不难看出, 算法的关键在于寻求具有最小权的最大匹配(以下简称最佳匹配)。下面介绍寻找最佳匹配的算法的思想。

设 $G_1$ 为具有二分类 $(X^*, Y^*)$ 的带权偶图。 $X^* = \{Rx_1, Rx_2, \dots, Rx_r\}$ ,  $Y^* = \{Ry_1, Ry_2, \dots, Ry_r\}$ ,  $k \leq r$ 。 $G_1$ 的每条边的权 $W_{ij}$ 为:

$$W_{ij} = C_{ij} + p_{ij} \quad (10)$$

然后人为地在 $G_1$ 上加若干虚边, 使之成为一个完全偶图。这些虚边的权定义为:

$$W'_{ij} = W_{\max} + 1 \quad (11)$$

$W_{\max} = \max \{W_{ij} \text{ 对所有的 } e(Rx_i, Ry_j) \in E_0\}$ 。有了以上假想, 在以后的讨论中, 都把 $G_1$ 看成一个完全偶图。

**定义4.1:** 若在顶点集 $X^* \cup Y^*$ 上的实值函数 $L$ 满足下述条件: 对所有的 $R_x \in X^*$ ,  $R_y \in Y^*$ , 均有

$$L(R_x) + L(R_y) \leq w(R_x, R_y) \quad (12)$$

则把这个函数定义为图 $G_1$ 的一个可行顶点标号。实数 $L(v)$ 称为顶点 $v$ 的标号。可行顶点标号是这样的标号, 它使每条边的两个顶点的标号之和不会超过这条边的权。不管边的权是什么, 总存在一个可行顶点标号:

$$\begin{aligned} L(R_x) &= \min w(R_x, R_y) \text{ 若 } R_x \in X^* \\ L(R_y) &= 0 \text{ 若 } R_y \in Y^* \end{aligned} \quad (13)$$

若 $L$ 是可行顶点标号, 则用 $E_l$ 表示使(12)式等式成立的那些边的集。即:

$$E_l = \{R_x, R_y \in E_0 \mid L(R_x) + L(R_y) = w(R_x, R_y)\} \quad (14)$$

具有边集 $E_l$ 的生成子图不妨称为对应于可行顶点标号的相等子图, 用 $G_l$ 表示。相等子图 $G_l$ 与最佳匹配 $M$ 间的联系由下述定理给出。

**定理4.1:** 设 $L$ 是 $G_1$ 的可行顶点标号, 若 $G_l$ 包含最大匹配 $M^*$ , 则 $M^*$ 是 $G_1$ 的最佳匹配。

证: 假设 $G_1$ 包含最佳匹配。由于 $G_l$ 是 $G_1$ 的生成子图, 所以 $M^*$ 也就是 $G_l$ 最佳匹配。于是

$$w(M^*) = \sum_{e \in M^*} w(e) = \sum_{v \in V_0} l(v) \quad (15)$$

这是因为每个边 $e \in M^*$ 都属于这个相等子图, 并且 $M^*$ 的边的端点覆盖 $V_0$ 的每个顶点恰好一次。另一方面, 若 $M$ 是 $G_1$ 的任一最佳匹配, 则有:

$$w(M) = \sum_{e \in M} w(e) \leq \sum_{v \in V_0} V_0 l(v) \quad (16)$$

从(15)式和(16)式推出 $w(M^*) \geq w(M)$ 。于是 $M^*$ 是最佳匹配, 证毕。

下面叙述最佳匹配算法的思想。

首先给出任一可行顶点标号 $l$ , 然后决定 $G_l$ 。在 $G_l$ 中选择任一匹配 $M$ , 并且利用匈牙利方法[6], 若在 $G_l$ 中已经找到一个最大匹配, 则由定理4.1, 该匹配就是最佳的。否则匈牙利方法将终止于一个非完全的匹配 $M'$ 和一棵既不包含 $M'$ 可扩路径, 又不能使 $G_l$ 中进一步生长 $M'$ 的交错树 $H$ 。随后, 把 $l$ 修改为具有下述性质的另一个可行顶点标号 $l'$ :  $M'$ 和 $H$ 都包

含在 $G'$ 中, 并且 $H$ 能够在 $G'$ 中伸展。每当必要时, 连续不断地进行这种可行顶点标号的修改, 直到一个最大匹配在某个相等子图中找到为止。

## 2. 寻找最佳匹配的算法OPTIM

从任一可行顶点标号开始, 然后决定 $G_1$ , 并且在 $G_1$ 中选取任一匹配 $M$ 。

1). 若 $X^*$ 是 $M$ 饱和的, 则 $M$ 是最大匹配, 并由定义4.1可知 $M$ 是最佳匹, 转到第4步。否则, 令 $u \in X^*$ 是一个 $M$ 非饱和顶点, 置 $s = \{u\}$ ,  $T = \phi$ 。

2). 若 $V_u(s) \supset T$ , 则转到第3步。否则,  $V_u(s) = T$ , 计算:

$$dl = \min_{\substack{R_x \in S \\ R_y \in T}} \{c(R_x R_y) - (L(R_x) + L(R_y))\} \quad (17)$$

且由

$$L'(v) = \begin{cases} L(v) + dl & \text{若 } v \in S \\ L(v) - dl & \text{若 } v \in T \\ L(v) & \text{其它} \end{cases} \quad (18)$$

给出可行顶点标号 $L'$  (注意 $dl > 0$ , 且 $V_u(S) \supset T$ )。以 $L'$ 代替 $L$ ,  $G'_1$ 代替 $G_1$ 。

3). 在 $V_u(S) \setminus T$ 中选择一个顶点 $v'$ , 考察 $v'$ 是否 $M$ 饱和。若 $v'$ 是 $M$ 饱和的, 则一定存在某条边 $e \in M$ , 其关联的两个顶点为 $v'$ 和 $Z \in X^*$ , 用 $S \cup \{Z\}$ 代替 $S$ , 用 $T \cup \{v\}$ 代替 $T$ , 再转到第2步。否则, 设 $P$ 是 $G_1$ 中的 $M$ 可扩 $(u, v')$ 路, 用 $M' = M \Delta E(P)$ 代替 $M$ , 并转到第1步。

4). 从 $M$ 中删除那些具有权为 $W_{ij} = W_{max} + 1$ 的边。算法终止。

下面来分析算法OPTIM的复杂性。

**定理4.2:** 算法OPTIM的时间复杂性为 $O(r^4)$ 。

证: 在算法的第2步, 计算 $G_1$ 所需的计算次数显然是 $|V_0|^2 = (r+k)^2$ 级的。由于在找出 $M$ 可扩路之前, 算法的循环最多可能有 $|X|$ 次。所以算法总的计算次数为

$$|V_0|^2 * |X|^2 = (r+k)^2 k^2 = k^4 + 2rk^3 + r^2 k^2 \quad (19)$$

所以, 算法的复杂性为 $O(r^4)$ 。

## 3. 偶查询优化算法GEN

在第1节已经介绍了偶查询优化算法的基本思想, 下面给出算法描述。

ALGORITHM GEN

INPUT:  $G_j$

OUTPUT: An Optimal Join Order.

1. WHILE  $G_j$  NOT a isolated vertex  
BEGIN
2. OPTIM( $G_j$ );
3. WHILE M NOT NIL  
BEGIN
4. FOR all  $e(u, v) \in M$ , REDUCE( $e(u, v)$ )
5.  $M = M - \{e(u, v)\}$ ;
- END

6. COMPUT-WC( $G_j$ );

END

7. END of ALGORITHM.

过程REDUCE把 $M$ 中的每条边压缩为一个顶点, 并去掉新形成的联接图 $G'_j$ 中的重复边; 过程COMPUT-W根据式(1)~(9)重新计算边的权。计算时, 只有那些受到影响的边的权才需重新计算。

## 五、结 语

本文针对一类特殊的多关系查询——偶查询, 提出了一种有效的优化方法。这种方法是建立在图论的偶图和匹配理论基础上的。把确定最优次序问题转化为在一个带权的偶图上寻求一个具有最小权的最大匹配问题。算法具有多项式复杂性。

## 参 考 文 献

- [1] Aho A V, Sasiv Y, Ulman J D. Efficient Optimization of A Class of Relational Expression. ACM trans, Database Syst, 1979; 4(4): 435~454
- [2] Apers P M G, Hevner A R, Yao S B. Optimization Algorithms for Distributed Queries. IEEE trans. on Soft. Eng., 1983; 9(1): 57~68
- [3] Bernstein P A, Chiu D W. Using Semi-joins to Solve Relational Queries. J. ACM, 1981; 28(1): 25~40
- [4] Bernstein P A etc. Query Processing in a System for Distributed Databases (RDD-1). ACM trans. Databases Syst., 1981; 6(4): 602~625
- [5] Bondy J. A, Murty U R. Graph Theory with Applications. The MACMILLAN Press ltd, 1976;
- [6] Chandra A K, Merlin P M. Optimal Implementation of Conjunctive Queries in Relatioal Data Bases. Proc. Ninth Ann. ACM Symp. on Theory of Uompt., 1976; 77~90
- [7] Epstein R, Stonebraker M. Distributed Query Processing in A Relational Data Base System. Proc. of ACM-SIGMOD, 1978; 169~180
- [8] Goodman N, Shmueli O. Tree Queries: A Simple Class of Relational Queries. ACM trans. Database Syst., 1982; 7(4): 653~677
- [9] Ibaraki T, Kameda T. On the Optimal Nesting Order for Computing N-Relational Joins. ACM trans. Database Syst., 1984; 9(3): 428~502
- [10] Selinger P G, Adiba M. Access Path Selection in Distributed Data Management System. Proc. Conf. Databases, 1980; 204~215
- [11] Stonebraker M, Wang E, Kroeps P. The Design and Implementation of INGRES ACM trans. Database Syst., 1976; 1(3): 189~222
- [12] Wong E, Youssefi K. Decomposition-A Strategy for Query Processing. ACM trans. Databases Syst., 1976; 1(3): 223~241