

# Adaptive Algorithms for Join Processing in Distributed Database Systems

PETER SCHEUERMANN\*

*Department of Electrical and Computer Engineering, Northwestern University, Evanston, IL 60208*

peters@ece.nwu.edu

EUGENE INSEOK CHONG†

*New England R&D Center, Oracle Corporation, 1 Oracle Drive, Nashua, NH 03062*

echong@us.oracle.com

*Received May 10, 1996; Accepted February 19, 1997*

**Recommended by:** Clement Yu

**Abstract.** Distributed query processing algorithms usually perform data reduction by using a semijoin program, but the problem with these approaches is that they still require an explicit join of the reduced relations in the final phase. We introduce an efficient algorithm for join processing in distributed database systems that makes use of bipartite graphs in order to reduce data communication costs and local processing costs. The bipartite graphs represent the tuples that can be joined in two relations taking also into account the reduction state of the relations. This algorithm fully reduces the relations at each site. We then present an adaptive algorithm for response time optimization that takes into account the system configuration, i.e., the additional resources available and the data characteristics, in order to select the best strategy for response time minimization. We also report on the results of a set of experiments which show that our algorithms outperform a number of the recently proposed methods for total processing time and response time minimization.

**Keywords:** distributed query processing, join algorithms, adaptive algorithms, bipartite graphs

## 1. Introduction

Query processing in distributed database systems often requires the transmission of relations and/or temporary results among different sites via a computer network. Until recently it has been assumed that communication costs were the predominant costs and hence much of the research on query optimization in distributed database systems has concentrated on producing optimal or near-optimal strategies with regard to data transmission costs [1, 3, 5, 6, 8, 12-14].

In order to reduce the communication time, most algorithms for distributed query processing involve three phases [20, 30]: 1) a local processing phase which includes all local operations such as selections and projections, 2) a reduction phase in which an effective sequence of semi-join (and join) operations is further used to reduce the sizes of the relations,

\*The work of this author was partially supported by NSF grant IRI-9303583 and NASA-Ames grant NAG2-846.

†The work of this author was performed while he was with Northwestern University.



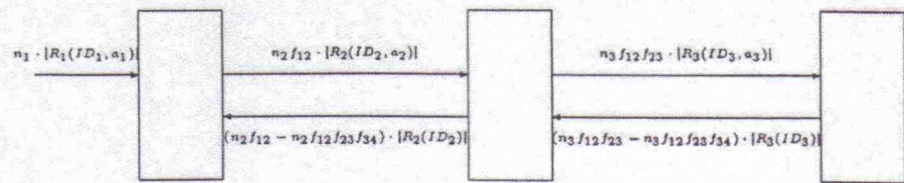


Figure 2. Sizes of transmitted data (PIPE.CHQ).

components. The dominant components of the processing time are the sum  $S_{i,f}$ , the time to construct the bipartite graphs in the forward reduction phase, and  $GT$ , the graph traversal time. Inside the component  $S_{i,f}$ , the major contribution is the I/O cost for the merge-sort step, while inside  $GT$ , the major cost is the I/O time involved in retrieving the target attributes for the tuples in the implicit join. If we denote by  $\bar{P}_i$  the average relation size in pages and by  $K$  the cardinality of the implicit join, then we obtain that the time complexity of the processing time is bounded by  $O(n\bar{P}_i \log_2 \bar{P}_i + nK)$ , with  $n$  being the number of sites.

## 2.2. Disk-based systems

The previous algorithm and cost model were based on the assumption that the bipartite graphs fit in main memory. We shall show now that our algorithm can be easily extended to disk-based systems where only a portion of each graph is memory resident. As each bipartite graph is constructed it is dynamically partitioned into subgraphs that are stored on secondary storage, i.e., whenever the buffer allocated to the construction of the graph  $BG_{R_i, R_{i+1}}$  becomes full its page(s) is(are) flushed to secondary storage. Thus, each page on secondary storage consists of a subgraph. We observe that these subgraphs are not necessarily disjoint, i.e., they may have crossing edges.

In order to minimize the number of I/O operations that are performed on the bipartite graphs during the backward reduction phase, each graph stored at site  $S_i$  is augmented so as to include for every edge  $(id_{i-1}, id_i)$  also the page number in secondary storage associated with the tuple containing  $id_{i-1}$  at  $S_{i-1}$ , to be denoted by  $page(id_{i-1})$ . Thus, we can view now a bipartite graph  $BG_{R_{i-1}, R_i}$  as corresponding to a normalized triary relation  $BG_{R_{i-1}, R_i}(ID_{i-1}, ID_i, PAGE(ID_{i-1}))$ . Similarly, the messages sent in the forward and backward phases are augmented correspondingly to include page numbers. Hence, forward messages sent from  $S_i$ , the site of  $BG_{R_{i-1}, R_i}$ , are now of the form  $(id_i, a_i, \{page(id_i)\})$ , with  $a_i$  being the value of the join attribute  $A_i$  and  $\{page(id_i)\}$  standing for the page number(s) in secondary storage associated with the tuples containing  $id_i$  at  $S_i$ . Backward messages sent from a site  $S_i$  have the format  $(id_{i-1}, page(id_{i-1}))$ . As we shall see below, using these page numbers we can guarantee that during backward reduction each page of a bipartite graph is read and written back to storage only once.

Since the construction of the bipartite graphs during the forward phase does not guarantee that the subgraphs stored on secondary storage are disjoint, at  $S_i$  a particular tuple identifier  $id_i$  may appear in a number of distinct pages, connected in each to disjoint (sets of) tuple

