

Query Evaluation Techniques for Large Databases

Goetz Graefe
University of Colorado at Boulder
Computer Science Department

Draft of Friday, June 5, 1992, 16:32

Revised and expanded version of
Technical Report CU-CS-579-92, January 1992

Abstract

Database management systems will continue to manage large data volumes. Thus, efficient algorithms for accessing and manipulating large sets and sequences will be required to provide acceptable performance. The advent of object-oriented and extensible database systems will not solve this problem; on the contrary, modern data models exacerbate the problem. In order to manipulate large sets of complex objects as efficiently as today's database systems manipulate simple records, query processing algorithms and software will become more complex, and a solid understanding of algorithm and architectural issues is essential for the designer of database management software.

This survey provides a foundation for the design and implementation of query execution facilities in new database management systems. It describes a wide array of practical query evaluation techniques for both relational and post-relational database systems, including iterative execution of complex query evaluation plans, the duality of sort- and hash-based set matching algorithms, types of parallel query execution and their implementation, and special operators for emerging database application domains.

Index Terms

Relational, Extensible, and Object-Oriented Database Systems; Query Execution Architecture; Iterators; Complex Query Evaluation Plans; Set Matching Algorithms; Sort-Hash Duality; Dynamic Query Evaluation Plans; Operator Model of Parallelization; Parallel Algorithms; Emerging Database Application Domains.

Introduction

Effective and efficient management of large data volumes is necessary in virtually all computer applications, from business data processing to library information retrieval systems, multimedia applications with images and sound, computer aided design and manufacturing, real-time process control, and scientific computation. While database management systems are standard tools in business data processing, they are only slowly being introduced to all the other emerging database application areas.

In most of these new application domains, database management systems have traditionally not been used for two reasons. First, restrictive data definition and manipulation languages can make application development and maintenance unbearably cumbersome. Research into semantic and object-oriented data models and into persistent database programming languages has been addressing this problem and will eventually lead to acceptable solutions. Second, data volumes might be so large or complex that the real or perceived performance advantage of file systems is considered more important than all other criteria, e.g., the higher level of abstraction and programmer productivity typically achieved with database management systems. Thus, object-oriented database management systems that are designed for non-traditional database application domains and extensible database management systems toolkits that support a variety of data models must provide excellent performance to meet the

challenges of very large data volumes, and techniques for manipulating large data sets will find renewed and increased interest in the database community.

The purpose of this paper is to survey the software architecture of database query execution engines and efficient algorithms for executing complex queries over large databases. A "complex" query is one that requires a number of query processing algorithms to work together, and a "large" database uses files with sizes from several megabytes to many terabytes, which are typical for database applications in the present and the near future [252]. This survey discusses a large variety of query execution techniques that must be considered when designing and implementing the query execution module of a new database management system: algorithms and their execution costs, sorting vs. hashing, parallelism, resource allocation and scheduling issues in complex queries, special operations for emerging database application domains such as statistical and scientific databases, and general performance-enhancing techniques such as precomputation and compression. While many, although not all, techniques discussed in this paper have been developed in the context of relational database systems, most of them are applicable to and useful in the query processing facility for any database management system and any data model, provided the data model permits non-procedural queries over "bulk" data types such as sets and lists.

It is assumed that the reader possesses basic textbook knowledge of database query languages, in particular of relational algebra, and of file systems, including some basic knowledge about indices. As shown in Figure 1, query processing fills the gap between database query languages and file systems, It can be divided into optimization and execution. A query optimizer translates a query expressed in a high-level query language into a sequence of operations that are provided by the query execution engine or the file system. The goal of query optimization is to find a query evaluation plan that minimizes the most relevant performance measure, which can be elapsed time, the database user's waiting time for the first or last result item, the number of I/O operations, the total resource usage, the amount of memory required for query execution, all of the above, or some other performance measure. Query optimization is a special form of planning, employing techniques from artificial intelligence such as plan representation, search including directed search and pruning, dynamic programming, branch-and-bound algorithms, etc. The query execution engine is a collection of query execution operators and mechanisms for operator communication and synchronization — it uses concepts from algorithm design, operating systems, networks, and parallel and distributed computation. The facilities of the query execution engine define the space of possible plans that can be chosen by the query optimizer.

A general outline of the steps required for processing a database query are shown in Figure 2. Of course, this sequence is only a general guideline, and different database systems may use different steps or merge multiple steps into one, in particular extensible and object-oriented systems. After a query or request has been entered into the database system, be it interactively or by an application program, the query is parsed into an internal form.

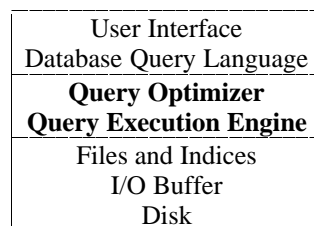


Figure 1. Query Processing in a Database System.

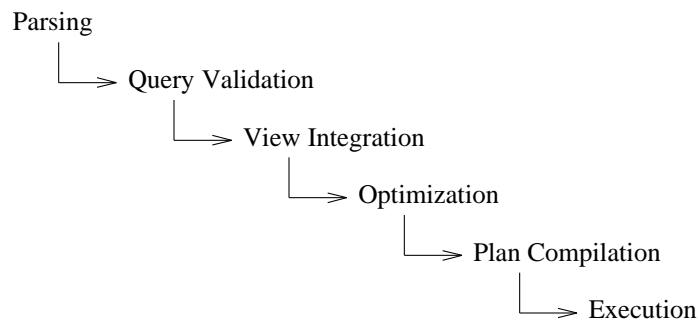


Figure 2. Query Processing Steps.

Next, the query is validated against the meta-data (data about data, also called schema or catalogs) to ensure that the query references only existing database objects. If the database system provides a macro facility such as relational views, referenced macros and views are expanded into the query [262]. Integrity constraints might be expressed as views (externally or internally) and would also be integrated into the query at this point in most systems [193]. The optimizer then maps the expanded query expression into an optimized plan that operates directly on the stored database objects. This mapping process can be very complex and might require substantial search and cost estimation effort. Optimization is not discussed in this paper; a relatively recent survey can be found in [150]. The optimizer's output is called a query execution plan, query evaluation plan, QEP, or simply plan. This plan is translated using a simple tree traversal algorithm for execution by the database's query execution engine; the result of this translation can be compiled machine code or a semi-compiled or interpreted language.

This survey explicitly discusses only read-only queries; however, most of the techniques are also applicable to update requests. In most database management systems, update requests may include a search predicate to specify which database objects are to be modified. Standard query optimization and execution techniques apply to this search; the actual update procedure can be either applied in a second phase, a method called *deferred updates*, or merged into the search phase if there is no danger of creating ambiguous update semantics.¹ The problem of ensuring *ACID* semantics for updates, — making updates *Atomic* (all-or-nothing semantics), *Consistent* (translating any consistent database state into another consistent database state), *Isolated* (from other queries and requests), and *Durable* (persistent across all failures) — is beyond the scope of this paper; suitable techniques have been described by many other authors [20, 22, 117, 127].

Embedded queries, i.e., database queries that are contained in an application program written in a standard programming language such as Cobol, PL/1, C, or Fortran, are also not addressed specifically in this paper because all techniques discussed here can be used for interactive as well as embedded queries. Embedded queries usually are optimized when the program is compiled in order to avoid the optimization overhead when the program runs. This method was pioneered in System R, including mechanisms for storing optimized plans and invalidating stored

¹ A standard example for this danger is the "Halloween" problem: Consider the request to "give all employees with salaries greater than \$30,000 a 3% raise." If (i) these employees are found using an index on salaries, (ii) index entries are scanned in increasing salary order, and (iii) the index is updated immediately as index entries are found, each such employee will get an infinite number of raises.

plans when they become infeasible, e.g., when an index is dropped from the database [49].

Recursive queries are omitted from this survey because the entire field of recursive queries — optimization rules and heuristics, selectivity and cost estimation, algorithms and their parallelization — is still developing rapidly; suffice it to point to a survey by Bancilhon and Ramakrishnan [9].

Section 1 discusses the architecture of query execution engines. Sorting and hashing, the two general approaches to managing and matching elements of large sets, are described in Section 2. Section 3 focuses on accessing large data sets on disk, including a discussion of indexing methods and disk arrays. Section 4 begins the discussion of actual data manipulation methods with algorithms for aggregation and duplicate removal, continued in Section 5 with binary matching operations such as join and intersection and in Section 6 with operations for universal quantification. Section 7 reviews the many dualities between sorting and hashing and points out their differences that have an impact on the performance of algorithms based on either one of these approaches. Execution of very complex query plans with many operators and with non-trivial plan shapes is discussed in Section 8. Section 9 is devoted to mechanisms for parallel execution, including architectural issues, load balancing, and tuning, and Section 10 discusses specific parallel algorithms. Section 11 outlines some non-standard operators for emerging database applications such as statistical and scientific database management systems. Section 12 is a potpourri of additional techniques that enhance the performance of many algorithms, e.g., compression, precomputation, and specialized hardware. The final section contains a brief summary and an outlook on query processing research and its future.

1. Architecture of Query Execution Engines

This survey focuses on useful mechanisms for processing sets of items. These items can be records, tuples, entities, or objects. Furthermore, most of the techniques discussed in this survey apply to sequences, not only sets, of items, although most query processing research has assumed relations and sets. All query processing algorithm implementations iterate over the members of their input sets; thus, sets are always represented by sequences. Sequences can be used to represent not only sets but also other one-dimensional "bulk" types such as lists, arrays, and time series, and many database query processing algorithms and techniques can be used to manipulate these other bulk types as well as sets. The important point is to think of these algorithms as algebra operators consuming zero or more inputs (sets or sequences) and producing one (or sometimes more) outputs. A complete query execution engine consists of a collection of operators and mechanisms to execute complex expressions using multiple operators, including multiple occurrences of the same operator. Taken as a whole, the query processing algorithms form an algebra which we call the *physical algebra* of a database system.

The physical algebra is equivalent to, but quite different from, the *logical algebra* of the data model or the database system. The logical algebra is more closely related to the data model and defines what queries can be expressed in the data model; for example, the relational algebra is a logical algebra. A physical algebra, on the other hand, is system-specific. Different systems may implement the same data model and the same logical algebra but may use very different physical algebras. For example, while one relational system may use only nested loops joins, another system may provide both nested loops join and merge-join, while a third one may rely entirely on hash join algorithms. (Join algorithms are discussed in detail later in the section on binary matching operators and algorithms.)

Another significant difference between logical and physical algebras is the fact that specific algorithms and therefore cost functions are associated only with physical operators, not with logical algebra operators. Because of

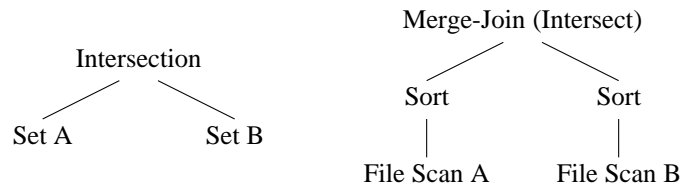


Figure 3. Logical and Physical Algebra Expressions.

the lack of an algorithm specification, a logical algebra expression is not directly executable and must be mapped into a physical algebra expression. For example, it is impossible to determine the execution time for the left expression in Figure 3, i.e., a logical algebra expression, without mapping it first into a physical algebra expression such as the query evaluation plan on the right of Figure 3. This mapping process can be trivial in some database systems but usually is fairly complex in real database systems because it involves algorithms choices and because logical and physical operators frequently do not map directly into one another, as shown in these three examples. First, some operators in the physical algebra may implement multiple logical operators. For example, all serious implementations of relational join algorithms include a facility to output fewer than all attributes, i.e., a relational delta-project (a projection without duplicate removal) is included in the physical join operator. Second, some physical operators implement only part of a logical operator. For example, a duplicate removal algorithm implements only the "second half" of a relational projection operator. Third, some physical operators do not exist in the logical algebra. For example, a sort operator has no place in pure relational algebra because it is an algebra of sets and sets are, by their definition, unordered.

Finally, some properties that hold for logical operators do not hold, or only with some qualifications, for the counterparts in physical algebra. For example, while intersection and union are entirely symmetric and commutative, algorithms implementing them (e.g., nested loops or hybrid hash join) do not treat their two inputs equally.

The purpose of the query execution engine is to execute physical algebra expressions produced by the query optimizer from a logical algebra expression. Thus, the premier design goal for the query execution engine is to implement a variety of efficient query execution mechanisms. The policies for using these mechanisms are built into an optimizer.

Synchronization and data transfer between operators is the main issue to be addressed in the architecture of the query execution engine. Imagine a query with two joins and consider how the result of the first join is passed to the second one. The simplest method is to create (write) and read a temporary file. The need for temporary files, whether they are kept in the buffer or not, is a direct result of executing an operator's input subplans completely before starting the operator. Alternatively, it is possible to create one process for each operator and then to use interprocess communication mechanisms (e.g., pipes) to transfer data between operators, leaving it to the operating system to schedule and suspend operator processes as pipes are full or empty. While such data-driven execution removes the need for temporary disk files, it introduces another cost, that of operating system scheduling and interprocess communication. In order to avoid both temporary files and operating system scheduling, Freytag proposed writing rule-based translation programs that transform a plan represented as a tree structure into a single iterative program with nested loops and other control structures [89]. However, the required rule set is not simple, in particular for algorithms with complex control logic such as sorting, merge-join, or even hybrid hash join (to be

discussed later in the section on matching).

The most practical alternative is to implement all operators in such a way that they *schedule each other within a single operating system process*. The basic idea is to define a granule, typically a single record, and to iterate over all granules comprising an intermediate query result². Each time an operator needs another granule, it calls its input (operator) to produce one. This call is a simple procedure call, much cheaper than inter-process communication since it does not involve the operating system at all. The calling operator waits (just as any calling routine waits) until the input operator has produced an item. That input operator, in a complex query plan, might require an item from its own input to produce an item; in that case, it calls its own input (operator) to produce one. Two important features of operators implemented in this way are that they can be combined into arbitrarily complex query evaluation plans and that any number of operators can execute and schedule each other in a single process without assistance from or interaction with the underlying operating system. This model of operator implementation and scheduling resembles very closely those used in relational systems, e.g., System R (and later SQL/DS and DB2), Ingres, Informix, and Oracle; as well as in extensible systems, e.g., Exodus' E programming language [214], Genesis [11, 12], and Starburst [125, 126]. Operators implemented in this model are called *iterators*, streams, synchronous pipelines, row-sources, or similar names in the "lingo" of commercial systems.

To make the implementation of operators a little easier, it makes sense to separate the functions (a) to

Iterator	<i>Open</i>	<i>Next</i>	<i>Close</i>
Print	<i>open</i> input	call <i>next</i> on input; format the item on screen	<i>close</i> input
Scan	open file	read next item	close file
Select	<i>open</i> input	call <i>next</i> on input until an item qualifies	<i>close</i> input
Hash join (without overflow resolution)	allocate hash directory; <i>open</i> left "build" input; build hash table calling <i>next</i> on build input; <i>close</i> build input; <i>open</i> right "probe" input	call <i>next</i> on probe input until a match is found	<i>close</i> probe input; deallocate hash directory
Merge-Join	<i>open</i> both inputs	get <i>next</i> item from input with smaller key until a match is found	<i>close</i> both inputs
Sort	<i>open</i> input; build all initial run files calling <i>next</i> on input and quicksort or replacement selection; <i>close</i> input; merge run files until only one merge step is left; open the remaining run files	determine next output item; read new item from the correct run file	destroy remaining run files

Table 1. Examples of Iterator Functions.

² It is possible to use multiple granule sizes within a single query processing system, and to provide special iterators with the sole purpose of translating from one granule size to another. An example is a query processing system that uses records as iteration granule except for the inputs of merge-join (see later in the section on binary matching), for which it uses "value packets," i.e., groups of records with equal join attribute values.

prepare an operator for producing data, (b) to produce an item, and (c) to perform final house-keeping. In a file scan, these functions are called *open*, *next*, and *close* procedures; we adopt these names for all operators. Table 1 gives a rough idea what the *open*, *next*, and *close* procedures for some operators do. The first three examples are trivial, but the *hash join* operator shows how an operator can schedule its inputs in a non-trivial manner. The interesting observations are that (i) the entire query plan is executed within a single process, (ii) operators produce one item at a time on request, (iii) this model effectively implements, within a single process, (special-purpose) *co-routines* and *demand-driven dataflow*, (iv) items never wait in a temporary file or buffer between operators because they are never produced before they are needed, (v) therefore this model is very efficient in its time-space-product memory costs, (vi) iterators can schedule any tree, including bushy trees (see below), (vii) no operator is affected by the complexity of the whole plan, i.e., this model of operator implementation and synchronization works for simple as well as very complex query plans. As a final remark, there are effective ways to combine the iterator model with parallel query processing, as will be discussed later.

Since query plans are algebra expressions, they can be represented as trees. Query plans can be divided into prototypical shapes, and query execution engines can be divided into groups according to which shapes of plans they can evaluate. Figure 4 shows prototypical left-deep, right-deep, and bushy plans for a four-way join. Left-deep and right-deep plans are different because join algorithms use their two inputs in different ways; for example, in the nested loops join algorithm, the outer loop iterates over one input (usually drawn as left input) while the inner loop iterates over the other input. The set of bushy plans is the most general as it includes the sets of both left-deep and right-deep plans. These names are taken from [100]; left-deep plans are also called "linear processing trees" [173] or "plans with no composite inner" [204].

For queries with common subexpressions, the query evaluation plan is not a tree but an acyclic directed graph (DAG). Most systems, if they identify and exploit common subexpressions, execute the plan equivalent to a common subexpression separately, saving the intermediate result in a temporary file to be scanned repeatedly and destroyed after the last scan. Each plan fragment that is executed as a unit is indeed a tree. The alternative is a "split" iterator that can deliver data to multiple consumers, i.e., that can be invoked as iterator by multiple consumer iterators. The split iterator paces its input subtree as fast as the fastest consumer requires it and holds items until the slowest consumer has consumed them. If the consumers request data at about the same rate, the split operator does not require a temporary spool file; such a file and its associated I/O cost is required only if the data rate required by the consumers diverges above some predefined threshold.

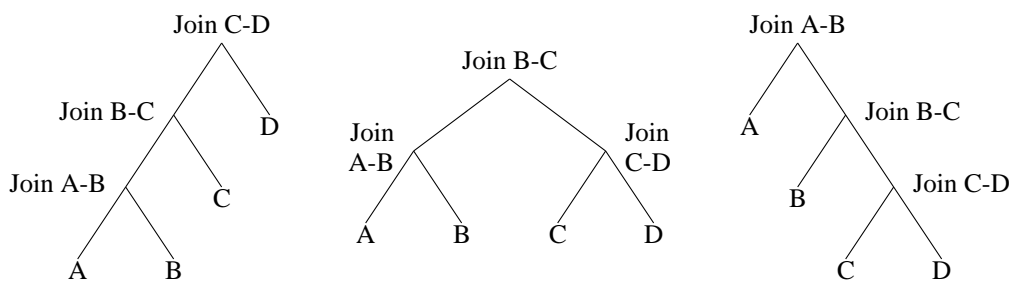


Figure 4. Left-Deep, Bushy, and Right-Deep Plans.

Among the implementations of iterators for query processing, one group can be called "stored-set-oriented" and the other "algebra-oriented." In System R, an example for the first group, complex join plans are constructed using binary join iterators that "attach" one more set (stored relation) to an existing intermediate result [4, 182], a design that supports only left-deep plans. This design led to a significant simplification of the System R optimizer which could be based on dynamic programming techniques but ignores the optimal plan for some queries³ [239]. A similar design was used, although not strictly required by the design of the execution engine, in the Gamma database machine [72, 74, 92]. On the other hand, some systems use binary operators for which both inputs can be intermediate results, i.e., the output of arbitrarily complex subplans. This design is more general as it also permits bushy plans. Examples for this approach are the second query processing engine of Ingres based on Kooi's thesis [169, 170], the Starburst execution engine [125], and the Volcano query execution engine [111]. The tradeoff between left-deep and bushy query evaluation plans is reduction of the search space in the query optimizer against generality of the execution engine and efficiency for some queries. Right-deep plans have only recently received more interest and may actually turn out to be very efficient, in particular in systems with ample memory [233, 234].

The remainder of this section provides more details of how iterators are implemented in the Volcano extensible query processing system. We use the Volcano system as an example repeatedly in this survey because it provides a large variety of mechanisms for database query processing, but mostly because its model of operator implementation and scheduling resembles very closely those used in relational and extensible systems. The purpose of this section is to provide implementation concepts from which a new query processing engine could be derived.

Figure 5 shows how iterators are represented in Volcano. A box stands for a record structure in Volcano's implementation language (C [160]), and an arrow represents a pointer. Each operator in a query evaluation plan consists of two record structures, a small structure of four pointers and a *state record*. The small structure is the same for all algorithms. It represents the stream or iterator abstraction and can be invoked with the *open*, *next*, and *close* procedures. The function of state records is similar to activation records allocated by compiled-generated code upon entry into a procedure. Both hold values local to the procedure or the iterator. Their main difference is that activation records reside on the stack and vanish upon procedure exit while state records must persist from one invocation of the iterator to the next, e.g., from the invocation of *open* to the each invocation of *next* and the invocation of *close*. Thus, state records do not reside on the stack but in heap space.

The type of state records is different for each iterator as it contains iterator-specific arguments and local variables (state) while the iterator is suspended, e.g., currently not active between invocation of the operators *next* procedure. Query plan nodes are linked together by means of *input* pointers, which are also kept in the state records. Since pointers to functions are used extensively in this design, all operator code (i.e., the *open*, *next*, and *close* procedures) can be written in such a way that the names of input operators and their iterator procedures are not "hard-wired" into the code, and the operator modules do not need to be recompiled for each query. Furthermore, all operations on individual items, e.g., printing, are imported into Volcano operators as functions, making the operators independent of the semantics and representation of items in the data streams they are processing. This organization using function pointers for input operators is fairly standard in commercial database management systems.

³ Since each operator in such a query execution system will access a permanent relation, the name "access path selection" used for System R optimization, although including and actually focusing on join optimization, was entirely correct and more descriptive than "query optimization."

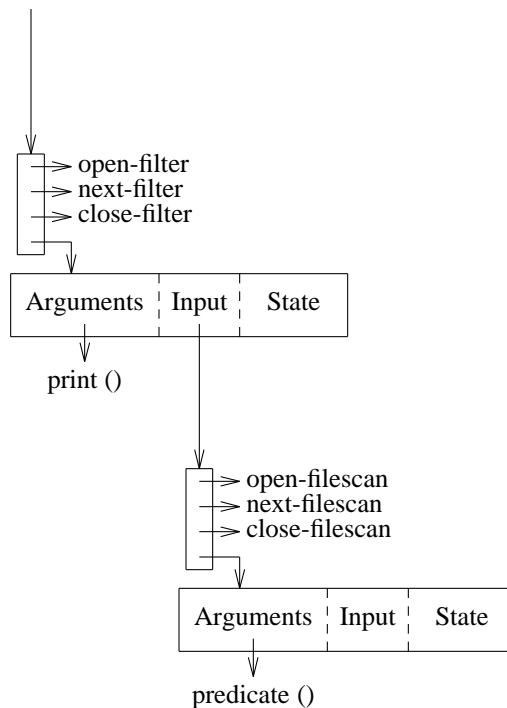


Figure 5. Two Operators in a Volcano Query Plan.

In order to make this discussion more concrete, Figure 5 shows two operators in a query evaluation plan that prints selected records from a file. The purpose and capabilities of the *filter* operator in Volcano includes printing items of a stream using a *print* function passed to the filter operator as one of its arguments. The small structure at the top gives access to the filter operator's iterator functions (the *open*, *next*, and *close* procedures) as well as to its state record. Using a pointer to this structure, the *open*, *next*, and *close* procedures of the filter operator can be invoked and their local state can be passed to them as a procedure argument. Filter's iterator functions themselves, e.g., *open-filter*, can use the input pointer contained in the state record to invoke the input operator's functions, e.g., *open-file-scan*. Thus, the filter functions can invoke the file scan functions as needed, and can pace the file scan according to the needs of the filter.

In this section, we have discussed general physical algebra issues and synchronization and data transfer between operators. Iterators are relatively straightforward to implement and are suitable building blocks for efficient, extensible query processing engines. In the following sections, we consider individual operators and algorithms including a comparison of sorting and hashing, detailed treatment of parallelism, special operators for emerging database applications such as scientific databases, and auxiliary techniques such as precomputation and compression.

2. Sorting and Hashing

Before discussing specific algorithms, two general approaches to managing sets of data are introduced. The purpose of many query processing algorithms is to perform some kind of matching, i.e., bringing items that are "alike" together and performing some operation on them. There are two basic approaches used for this purpose, sorting and hashing. This pair permeates many aspects of query processing, from indexing and clustering over aggregation and join algorithms to methods for parallelizing database operations. Therefore, we discuss these approaches first in general terms, without regard to specific algorithms. After a survey of specific algorithms for unary (aggregation, duplicate removal) and binary (join, semi-join, intersection, division, etc.) matching problems in the following sections, the duality of sort- and hash-based algorithms is discussed in detail.

2.1. Sorting

Sorting is used very frequently in database systems, both for presentation to the user in sorted reports or listings and for query processing in sort-based algorithms such as merge-join. Therefore, the performance effects of the many algorithmic tricks and variants of external sorting deserve detailed discussion in this survey. All sorting algorithms used in database systems use merging, i.e., the input data are written into initial sorted runs and then merged into larger and larger runs until only one run is left, the sorted output. Only in the unusual case that a data set is smaller than the available memory can in-memory techniques such as quicksort be used. An excellent reference for many issues discussed here is Knuth [167].

In order to ensure that the sort module interfaces well with the other operators, e.g., file scan or merge-join, sorting should be implemented as an iterator, i.e., with *open*, *next*, and *close* procedures as all other operators of the physical algebra. In the Volcano query processing system (which is based on iterators) most of the sort work is done during *open-sort* [105, 111]. This procedure consumes the entire input and leaves appropriate data structures for *next-sort* to produce the final, sorted output. If the entire input fits into the sort space in main memory, *open-sort* leaves a sorted array of pointers to records in I/O buffer memory which is used by *next-sort* to produce the records in sorted order. If the input is larger than main memory, the *open-sort* procedure creates sorted runs and merges them until only one final merge phase is left. The last merge step is performed in the *next-sort* procedure, i.e., when demanded by the consumer of the sorted stream, e.g., a merge-join. The input to the sort module must be an iterator, and sort uses *open*, *next*, and *close* procedures to request its input; therefore, sort input can come from a scan or a complex query plan, and the sort operator can be inserted into a query plan at any place or at several places.

There are two alternative methods for creating initial runs, also called "level-0 runs" here. First, an in-memory sort algorithm can be used, typically quicksort. Using this method, each run will have the size of memory and the number of initial runs W will be $W = \lceil R / M \rceil$ for input size R and memory size M . Second, runs can be produced using replacement selection. Replacement selection starts by filling memory with items which are organized into a priority heap, i.e., a data structure that efficiently supports the operations *insert* and *remove-smallest*. Next, the item with the smallest key is removed from the priority heap and written to a run file, and then immediately replaced in the priority heap with another item from the input. With high probability, this new item has a key larger than the item just written, and therefore will be included in the same run file. Notice that if this is the case, the first run file will be larger than memory. Now the second item (the currently smallest item in the priority heap) is written to the run file, and also replaced immediately in memory by another item from the input. This process repeats, always keeping the memory and the priority heap entirely filled. If a new item has a key smaller than the

last key written, the new item cannot be included in the current run file and is marked for the next run file. In comparisons among items in the heap, items marked for the current run file are always considered "smaller" than items marked for the next run file. Eventually, all items in memory are marked for the next run file, at which point the current run file is closed and a new one is created.

Using replacement selection, run files are typically larger than memory. If the input is already sorted or almost sorted, there will be only one run file. This situation could arise, for example, if a file is sorted on field A but should be sorted on A as major and B as minor sort key. If the input is sorted in reverse order, which is the worst case, each run file will be exactly as large as memory. If the input is random, the average run file will be twice the size of memory, except the first few runs (which get the process started) and the last run. On the average, the expected number of runs is about $W = \lceil R / (2 \times M) \rceil + 1$, i.e., about half as many runs as created with quicksort. A more detailed discussion and an analysis of replacement selection was provided by Knuth [167].

An additional difference between quicksort and replacement selection is the resulting I/O pattern during initial run creation. Quicksort results in bursts of reads and writes for entire memory loads from the input file and to initial run files, while replacement selection alternates between individual read and write operations. If only a single device is used, quicksort may result in faster I/O because fewer disk arm movements are required. However, if different devices are used for input and temporary files, or if the input comes as a stream from another operator, the alternating behavior of replacement selection may permit more overlap of I/O and processing and therefore result in faster sorting.

The problem with replacement selection is memory management. If input items are kept in their original pages in the buffer (in order to save copying data, a real concern for large data volumes) each page must be kept in the buffer until its last record has been written to a run file. On the average, half a page's records will be in the priority heap. Thus, the priority heap must be reduced to half the size (the number of items in the heap is one half the number of records that fit into memory), cancelling the advantage of longer and fewer run files. The solution to this problem is to copy records into a holding space and to keep them there while they are in the priority heap and until they are written to a run file. If the input items are of varying sizes, memory management is more complex than for quicksort because a new item may not fit into the space vacated in the holding space by the last item written into a run file. Solutions to this problem will introduce memory management overhead and some amount of fragmentation, i.e., the size of runs will be less than twice the size of memory. Thus, the advantage of having fewer runs must be balanced with the different I/O pattern and the disadvantage of more complex memory management.

The level-0 runs are merged into level-1 runs, which are merged into level-2 runs, etc., to produce the sorted output. During merging, a certain amount of buffer memory must be dedicated to each input run and the merge output. We call this amount of memory the *cluster size* C in this survey. The number of I/O clusters that fit in memory is the quotient of memory size and cluster size. The maximal merge fan-in F , i.e., the number of runs that can be merged at one time, is this quotient minus one cluster for the output. Thus, $F = \lfloor M / C - 1 \rfloor$. Since the sizes of runs grow by a factor F from level to level, the number of merge levels L , i.e., the number of times each item is written to a run file, is logarithmic with the input size, namely $L = \left\lceil \log_F (W) \right\rceil$.

There are four considerations that can improve the merge efficiency. The first two issues pertain to scheduling of I/O operations. First, scans are faster if read-ahead and write-behind are used; therefore, double buffering using two pages of memory per input run and two for the merge output might speed the merge process [228, 229]. The obvious disadvantage is that the fan-in is cut in half. However, instead of reserving $2 \times F + 2$ clusters, a

predictive method called *forecasting* can be employed in which the largest key in each input buffer is used to determine from which input run the next cluster will be read. Thus, the fan-in can be set to any number in the range $\lfloor M / (2 \times C) - 2 \rfloor \leq F \leq \lfloor M / C - 1 \rfloor$. One or two read-ahead buffers per input disk are sufficient, and $F = \lfloor M / C \rfloor - 3$ will be reasonable in most cases because it uses maximal fan-in with one forecasting input buffer and double-buffering for the merge output.

Second, if the operating system and the I/O hardware support them, using large cluster sizes for the run files is very beneficial. Larger cluster sizes will reduce the fan-in and therefore may increase the number of merge levels. However, each merging level is performed much faster because fewer I/O operations and disk seeks and latency delays are required. Furthermore, if the unit of I/O is equal to a disk track, rotational latencies can be avoided entirely with a sufficiently smart disk controller. Usually, relatively small fan-ins with large cluster sizes are the optimal choice, even if the sort requires multiple merge levels [105]. The precise tradeoff depends on disk seek, latency, and transfer times. It is interesting to note that the optimal cluster size and fan-in basically do not depend on the input size.

As a concrete example, consider sorting a file of $R = 50 \text{ MB} = 51,200 \text{ KB}$ using $M = 160 \text{ KB}$ of memory. The number of runs created by quicksort will be $W = \lceil 51200 / 160 \rceil = 320$. Depending on the disk access and transfer times (e.g., 25 ms disk seek and latency, 2 ms transfer time for a page of 4 KB), $C = 16 \text{ KB}$ will typically be a good cluster size for fast merging. If one cluster is used for read-ahead and two for the merge output, the fan-in will be $F = \lfloor 160 / 16 \rfloor - 3 = 7$. The number of merge levels will be $L = \lceil \log_7(320) \rceil = 3$. If a 16 KB I/O operation takes $T = 33 \text{ ms}$, the total I/O time, including a factor of two for writing and reading at each merge level, for the entire sort will be $2 \times L \times \lceil R / C \rceil \times T = 10.56 \text{ min}$.

An entirely different approach to determining optimal cluster sizes and the amount of memory allocated to forecasting and read-ahead is based on processing and I/O bandwidths and latencies. The cluster sizes should be set such that the I/O bandwidth matches the processing bandwidth of the CPU. Bandwidths for both I/O and CPU are measured here in record or bytes per unit time; instructions per unit time (MIPS) are irrelevant. It is interesting to note that the CPU's processing bandwidth is largely determined by how fast the CPU can assemble new pages; in other words, how fast the CPU can copy records within memory. This performance measure is usually ignored in modern CPU and cache designs geared towards high MIPS or MFLOPS numbers [205].

Tuning the sort based on bandwidth and latency proceeds in three steps. First, the cluster size is set such that the processing and I/O bandwidths are equal or very close to equal. If the sort is I/O-bound, the cluster size is increased for less disk access overhead per record and therefore faster I/O; if the sort is CPU-bound, the cluster size is decreased to slow the I/O in favor of a larger merge fan-in. Next, in order to ensure that the two processing components (I/O and CPU) never (or almost never) have to wait for one another, the amount of space dedicated to read-ahead is determined as the I/O time for one cluster multiplied by the processing bandwidth. Typically, this will result in one cluster of read-ahead space per disk used to store and read inputs run into a merge. Of course, in order to make read-ahead effective, forecasting must be used. Finally, the same amount of buffer space is allocated for the merge output (access latency times bandwidth) to ensure that merge processing never has to wait for the completion of output I/O. — It is an open issue whether these two alternative approaches to cluster size and read-ahead space result in different allocations and sorting speeds and whether one of them is more effective than the other.

The third and fourth merging issues focus on using (and exploiting) the maximal fan-in as effectively and often as possible. Both issues require adjusting the fan-in of the first merge step using the formula given below,

either the first merge step of all merge steps or, in semi-eager merging [105], the first merge step after the end of the input has been reached. This adjustment is used for only one merge step, called the *initial merge* here, not for an entire merge level.

The third issue to be considered is that the number of runs W is typically not a power of F ; therefore, some merges proceed with fewer than F inputs, which creates the opportunity for some optimization. Instead of always merging runs of only one level together, the optimal strategy is to merge as many runs as possible using the smallest run files available. The only exception is the fan-in of the first merge, which is determined to ensure that all subsequent merges will use the full fan-in F .

Let us explain this idea with the example shown in Figure 6a. Consider a sort with a maximal fan-in $F = 10$ and an input file that requires $W = 12$ initial runs. Instead of merging only runs of the same level as shown in Figure 6a, merging is delayed until the end of the input has been reached. In the first merge step, only 3 of the 12 runs are combined and the result is then merged with the other 9 runs, as shown in Figure 6b. The I/O cost (measured by the number of memory loads that must be written to disk to any of the runs created) for the first strategy is $12 + 10 + 2 = 24$, while for the second strategy it is $12 + 3 = 15$, meaning that the first strategy requires 60% more I/O to temporary files than the second one. The general rule is to merge just the right number of runs after the end of the input file has been reached, and to always merge the smallest runs available for merging. More detailed examples are given in [105]. One consequence of this optimization is that the merge depth L , i.e., the number of runs files a record is written to during the sort or the number of times is written and read from disk, is not uniform for all records. Therefore, it makes sense to calculate an average merge depth (as required in cost estimation during query optimization), which may be a fraction.

Fourth, since some operations require multiple sorted inputs, for example merge-join (to be discussed in the section on matching) and sort output can be passed directly from the final merge into the next operation (as is natural when using iterators), memory must be divided among multiple final merges. Thus, the final fan-in f and

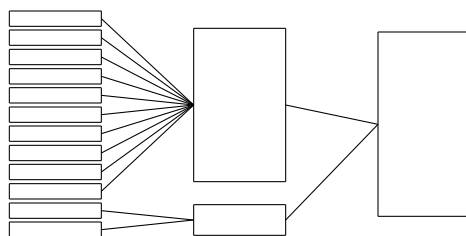


Figure 6. Naive Merging.

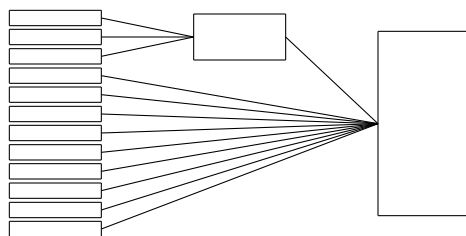


Figure 6. Optimized Merging.

the "normal" fan-in F should be specified separately in an actual sort implementation. Using a final fan-in of 1 also allows the sort operator to produce output into a very slow operator, e.g., a display operator that allows scrolling by a human user, without occupying a lot of buffer memory for merging input runs over an extended period of time.

Considering the last two optimization options for merging, the following formula determines the fan-in of the first merge. Each merge with normal fan-in F will reduce the number of run files by $F - 1$ (removing F runs, creating one new one), and the goal is to reduce the number of runs from W to f and then to 1 (the final output). Thus, the first merge should reduce the number of runs to $f + k(F - 1)$ for some integer k . In other words, the first merge should use a fan-in of $F_0 = ((W - f - 1) \bmod (F - 1)) + 2$. In the example of Figure 6a and Figure 6b, $(12 - 10 - 1) \bmod (10 - 1) + 2$ results in a fan-in for the initial merge of $F_0 = 3$. If the sort of Figure 6b were the input into a merge-join and a final fan-in of 5 were desired, the initial merge should proceed with a fan-in of $F_0 = (12 - 5 - 1) \bmod (10 - 1) + 2 = 8$.

If multiple sort operations produce input data for a common consumer operator, e.g., a merge-join, the two final fan-ins should be set proportionally to the size of the two inputs. For example, if two merge-join inputs are 1 MB and 9 MB, and 20 clusters are available for inputs into the two final merges, 2 clusters should be allocated for the first and 18 clusters for the second input ($1 / 9 = 2 / 18$).

Sorting is sometimes criticized because it requires, unlike hybrid hashing (discussed in the next subsection), that the entire input be written to run files and then retrieved for merging. This difference has a particularly large effect for files only slightly larger than memory, e.g., $1\frac{1}{4}$ times the size of memory. In hybrid hashing, only slightly more than $\frac{1}{4}$ of the memory size must be written to temporary files on disk while the remainder of the file remains in memory. In sorting, the entire file ($1\frac{1}{4}$ memory sizes) is written to one or two run files and then read for merging. Thus, sorting seems to require five times more I/O for temporary files than hybrid hashing. However, this is not necessarily true. The simple trick is to write initial runs in decreasing (reverse) order. When the input is exhausted and merging in increasing order commences, buffer memory is still full of useful pages with small sort keys that can be merged immediately without I/O and that never have to be written to disk.

To demonstrate the effect of cluster size optimizations (the second of the four merging issues discussed above), we sorted 100,000 100-byte records, about 10 MB, with the Volcano query processing system, which includes all merge optimizations described above with the exception of read-ahead and forecasting.⁴ We used a sort space of forty pages (160 KB) within a fifty-page (200 KB) I/O buffer, varying the cluster size from one page (4 KB) to fifteen pages (60 KB). The initial run size was 1,600 records, for a total of 63 initial runs. We counted the number of I/O operations and the transferred pages for all run files, and calculated the total I/O cost by charging 25 ms per I/O operation (for seek and rotational latency) and 2 ms for each transferred page (assuming 2 MB/sec transfer rate). As can be seen in Table 2 and Figure 7, there is an optimal cluster size with minimal I/O cost. It is clearly suboptimal to always choose the smallest cluster size (one page) to obtain the largest fan-in and fewest merge levels. Furthermore, it seems that the range of cluster sizes that result in near-optimal total I/O costs is fairly large; thus, it is not as important to determine the exact value as it is to use a cluster size "in the right ball park." The optimal fan-in is typically fairly small; however, it is not e or 3 as derived by Bratbergsengen under the (unrealistic) assumption that the cost of an I/O operation is independent of the amount of data being transferred [41].

⁴ This experiment has been reported before in [105].

Cluster Size [× 4 KB]	Fan-in	Average Depth	Disk Operations	Pages Transferred [× 4 KB]	Total I/O Cost [sec]
1	40	1.376	6874	6874	185.598
2	20	1.728	4298	8596	124.642
3	13	1.872	3176	9528	98.456
4	10	1.936	2406	9624	79.398
5	8	2.000	1984	9920	69.440
6	6	2.520	2132	12792	78.884
7	5	2.760	1980	13860	77.220
8	5	2.760	1718	13744	70.438
9	4	3.000	1732	15588	74.476
10	4	3.000	1490	14900	67.050
11	3	3.856	1798	19778	84.506
12	3	3.856	1686	20232	82.614
13	3	3.856	1628	21164	83.028
14	2	5.984	2182	30548	115.646
15	2	5.984	2070	31050	113.850

Table 2. Effect of Cluster Size Optimizations.

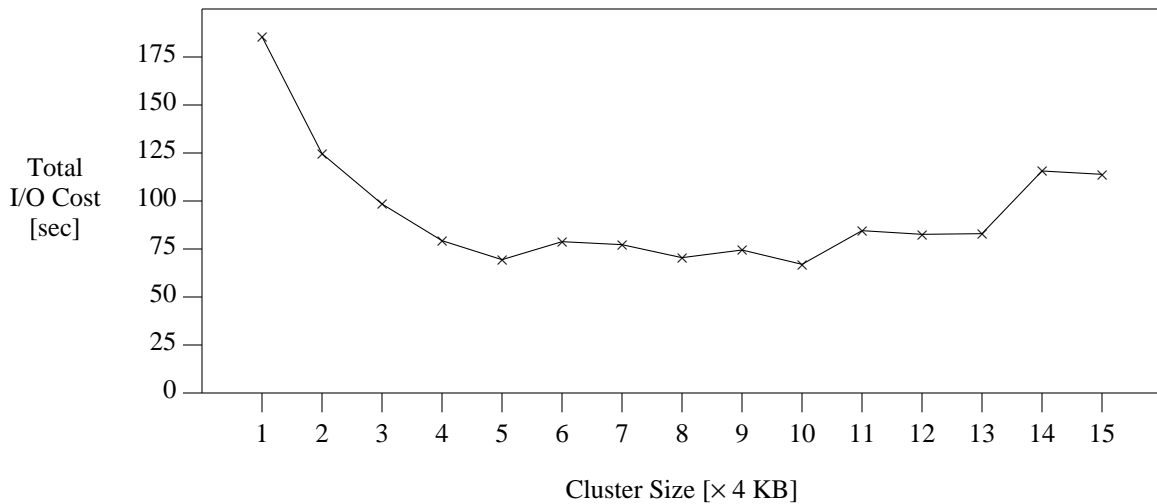


Figure 7. Effect of Cluster Size Optimizations.

2.2. Hashing

For many matching tasks, hashing is an alternative to sorting. In general, when equality matching is required, hashing should be considered because the expected complexity of set algorithms based on hashing is $O(N)$ rather than $O(N \log N)$ as for sorting. Of course, this makes intuitive sense if hashing is viewed as radix sorting on a virtual key [167].

Hash-based query processing algorithms use an in-memory hash table of database objects to perform their matching task. If the entire hash table (including all records or items) fits into memory, hash-based query processing algorithms are very easy to design, understand, and implement, and outperform sort-based alternatives. Note that for binary matching operations, such as join or intersection, only one of the two inputs must fit into memory. However, if the required hash table (i.e., the input into the hash-based query processing algorithm) is larger than

memory, *hash table overflow* occurs and must be dealt with.

There are basically two methods for managing hash table overflow, namely *avoidance* and *resolution*. In either case, the input is divided into multiple partition files such that partitions can be processed independently from one another and the concatenation of the results of all partitions is the result of the entire operation. Partitioning should ensure that the partitioning files are of roughly even size, and can be done using either hash-partitioning or range-partitioning, i.e., based on keys estimated to be quantiles. Usually, partition files can be processed using the original hash-based algorithm. The maximal partitioning *fan-out* F , i.e., number of partition files created, is determined by the memory size M divided over the cluster size C minus one cluster for the partitioning input, i.e., $F = \lfloor M / C - 1 \rfloor$, just like the fan-in for sorting.

In hash table overflow avoidance, the input set is partitioned into F partition files before any in-memory hash table is built. If it turns out that fewer partitions than have been created would have been sufficient to obtain partition files that will fit into memory, bucket tuning (collapsing multiple small buckets into larger ones) and dynamic destaging (determining which buckets should stay in memory) can improve the performance of hash-based operations [164, 194].

Algorithms based on hash table overflow resolution start with the assumption that overflow will not occur, but resort to basically the same set of mechanisms as hash table overflow avoidance once it does occur. No real system uses this naive hash table overflow resolution because so-called hybrid hashing is as efficient but more flexible. Hybrid hashing combines in-memory hashing and overflow resolution [70, 244]⁵. Hybrid hash algorithms start out with the (optimistic) premise that no overflow will occur; if it does, however, they partition the

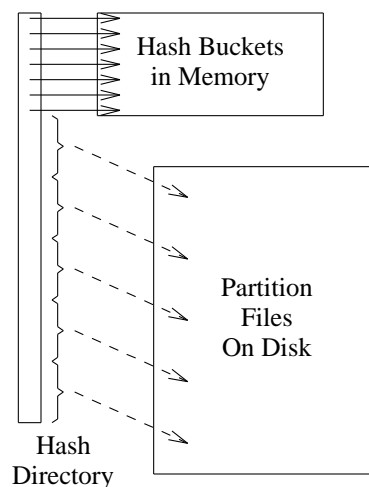


Figure 8. Hybrid Hashing.

⁵ Although invented for relational join and known as hybrid hash join, hybrid hashing is equally applicable to all hash-based query processing algorithms.

input into multiple partitions of which only one is written immediately to temporary files on disk. The other $F - 1$ partitions remain in memory. If another overflow occurs, another partition is written to disk. If necessary, all F partitions are written to disk. Thus, hybrid hash algorithms use all available memory for in-memory processing, but at the same time are able to process large input files by overflow resolution. Figure 8 shows the idea of hybrid hash algorithms. As many hash buckets as possible are kept in memory, e.g., as linked lists as indicated by solid arrows. The other hash buckets are spooled to temporary disk files, called the overflow or partition files, and are processed in later stages of the algorithm. Hybrid hashing is useful if the input size R is larger than the memory size M but smaller than the memory size multiplied by the fan-out F , i.e., $M < R \leq F \times M$.

In order to predict the number of I/O operations (which actually is not necessary for execution because the algorithm adapts to its input size but may be desirable for cost estimation during query optimization), the number of required partition files on disk must be determined. Call this number K , which must satisfy $0 \leq K \leq F$. Presuming that the assignment of buckets to partitions is optimal and each partition file is equal to the memory size M , the amount of data that may be written to K partition files is equal to $K \times M$. The number of required I/O buffers is 1 for the input and K for the output partitions, leaving $M - (K + 1) \times C$ memory for the hash table. The optimal K for a given input size R is the minimal K for which $K \times M + (M - (K + 1) \times C) \geq R$. Solving this inequality and taking the smallest such K results in $K = \lceil (R - M + C) / (M - C) \rceil$. The minimal possible I/O cost, including a factor of 2 for writing and reading the partition files and measured in the amount of data that must be written or read, is $2 \times (R - (M - (K + 1) \times C))$. To determine the I/O time, this amount must be divided by the cluster size and multiplied with the I/O time for one cluster.

For example, consider an input of $R = 240$ pages, a memory of $M = 80$ pages, and a cluster size of $C = 8$ pages. The maximal fan-out is $F = \lfloor 80 / 8 - 1 \rfloor = 9$. The number of partition files that need to be created on disk is $K = \lceil (240 - 80 + 8) / (80 - 8) \rceil = 3$. In other words, in the best case, $K \times C = 3 \times 8 = 24$ pages will be used as output buffers to write $K = 3$ partition files of no more than $M = 80$ pages, and $M - (K + 1) \times C = 80 - 4 \times 8 = 48$ pages of memory will be used as hash table. The total amount of data written to and read from disk is $2 \times (240 - (80 - 4 \times 8)) = 384$ pages. If writing or reading a cluster of $C = 8$ pages takes 40 msec, the total I/O time is $384 / 8 \times 40 = 1.92$ sec.

In the calculation of K , we assumed an optimal assignment of hash buckets to partition files. If buckets were assigned in the most straightforward way, e.g., by dividing the hash directory into F equal-size regions and assigning the buckets of one region to a partition as indicated in Figure 8, all partitions were of nearly the same size and either all or none of them will fit into their output cluster and therefore into memory. In other words, once hash table overflow occurred, all input were written to partition files. Thus, we presumed in the earlier calculations that hash buckets were assigned more intelligently to output partitions.

There are three ways to assign hash buckets to partitions. First, each time the hash table overflow occurs, a fixed number of hash buckets is assigned to a new output partition. In the Gamma database machine, the number of disk partitions is chosen "such that each bucket⁶ can be reasonably be expected to fit in memory" [71], e.g., 10% of the hash buckets in the hash directory for a fan-out of 10 [233]. In other words, the fan-out is set a priori by the query optimizer based on the expected (estimated) input size. Since the page size in Gamma is relatively small, only a fraction of memory is needed for output buffers, and an in-memory hash table can be used even while output partitions are being written to disk. Second, in bucket tuning and dynamic destaging [164, 194], a large

⁶ Bucket in [71] means what is called an output partition in this survey.

number of small partition files is created and then collapsed into fewer partition files no larger than memory. In order to obtain a large number of partition files and, at the same time, retain some memory for a hash table, the cluster size is set quite small, e.g. $C = 1$ page, and the fan-out is very large though not maximal, e.g., $F = M / C / 2$. In the example above, $F = 40$ output partitions with an average size of $R / F = 6$ pages could be created, even though only $K = 3$ output partitions are required. The smallest partitions are assigned to fill an in-memory hash table of size $M - K \times C = 80 - 3 \times 1 = 77$ pages. Hopefully, the dynamic destaging rule — when an overflow occurs, assign the largest partition still in memory to disk — ensures that indeed the smallest partitions are retained in memory. The partitions assigned to disk are collapsed into $K = 3$ partitions of no more than $M = 80$ pages, to be processed in $K = 3$ subsequent phases. In binary operations such as intersection and relational join, bucket tuning is quite effective for *skew* in the first input, i.e., if the hash value distribution is non-uniform and the partition files are of uneven sizes. It avoids spooling parts of the second (typically larger) input to temporary partition files because the partitions in memory can be matched immediately using a hash table in the memory not required as output buffer and because a number of small partitions have been collapsed into fewer, larger partitions, increasing the memory available for the hash table. For skew in the second input, bucket tuning and dynamic destaging has no advantage. Another disadvantage of bucket tuning and dynamic destaging is that the cluster size has to be relatively small, thus requiring a large number of I/O operations with disk seeks and rotational latencies to write data to the overflow files. Third, statistics gathered before hybrid hashing commences can be used to assign hash buckets to partitions [110].

Unfortunately, it is possible that one or more partition files are larger than memory. In that case, partitioning is used recursively until the file sizes have shrunk to memory size. Figure 9 shows how a hash-based algorithm for a unary operation such as aggregation or duplicate removal partitions its input file over multiple recursion levels. The recursion terminates when the files fit into memory. In the deepest recursion level, hybrid hashing may be employed.

If the partitioning (hash) function is good and creates a uniform hash value distribution, the file size in each recursion level shrinks by a factor equal to the fan-out, and therefore the number of recursion levels L is logarithmic with the size of the input being partitioned. After L partitioning levels, each partition file is of size

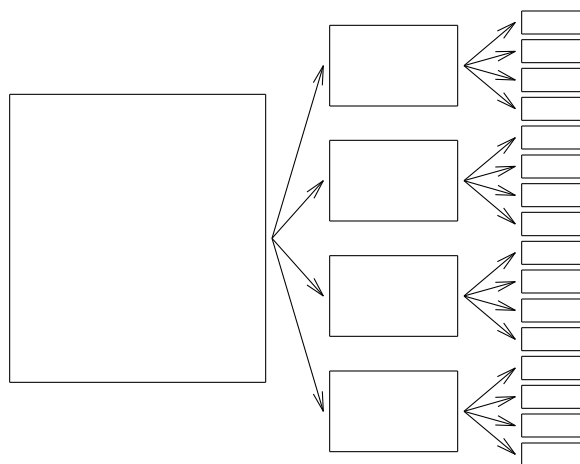


Figure 9. Recursive Partitioning.

$R' = R / F^L$. In order to obtain partition files suitable for hybrid hashing (with $M < R' \leq F \times M$), the number of full recursion levels L , i.e., levels at which hybrid hashing is not applied, is $L = \left\lceil \log_F (R / M) \right\rceil$. The I/O cost of the remaining step using hybrid hashing can be estimated using the hybrid hash formula above with R replaced by R' and multiplying the cost with F^L because hybrid hashing is used for this number of partition files. Thus, the total I/O cost for partitioning an input and using hybrid hashing in the deepest recursion level is

$$\begin{aligned} & 2 \times R \times L + 2 \times F^L \times \left[R' - (M - K \times C) \right] \\ &= 2 \times \left[R \times (L + 1) - F^L \times (M - K \times C) \right] \\ &= 2 \times \left[R \times (L + 1) - F^L \times \left[M - \lceil (R' - M) / (M - C) \rceil \times C \right] \right]. \end{aligned}$$

A major problem with hash-based algorithms is that their performance depends on the quality of the hash function. In many situations, fairly simple hash functions will perform reasonably well. Remember that the purpose of using hash-based algorithms usually is to find database items with a specific key or to bring like items together; thus, methods as simple as using the value of a join key as hash value will frequently perform satisfactorily. For string values, good hash values can be determined by using binary exclusive "or" operations or by determining cyclic redundancy check (CRC) values as used for reliable data storage and transmission. If the quality of the hash function is a potential problem, universal hash functions should be considered [46].

If the partitioning is skewed, the recursion depth may be unexpectedly high, making the algorithm rather slow. This is analogous to the worst-case performance of quicksort, $O(N^2)$ comparisons for an array of N items, if the partitioning pivots are chosen extremely poorly and do not divide arrays into nearly equal subarrays.

Skew is the major danger for inferior performance of hash-based query processing algorithms. There are several ways to deal with skew. For hash-based algorithms using overflow avoidance, bucket tuning and dynamic destaging are quite effective. Another method is to obtain statistical information about hash values and to use it to carefully assign hash buckets to partitions. Such statistical information can be kept in the form of histograms, and can either come from permanent system catalogs (meta-data), from sampling the input, or from previous recursion levels. For example, for an intermediate query processing result for which no statistical parameters are known a priori, the first partitioning level might have to proceed naively pretending that the partitioning hash function is perfect, but the second recursion and further levels should be able to use statistics gathered in earlier levels to ensure that each partitioning step creates even partitions, i.e., that the data is partitioned with maximal effectiveness [42]. As a final resort, if skew cannot be managed otherwise or if not distribution skew but duplicates are the problem, some systems resort to algorithms that are not affected by data or hash value skew. For example, Tandem's hash join algorithm resorts to nested loops join (to be discussed later) [298].

As for sorting, larger cluster sizes result in faster I/O at the expense of smaller fan-outs, with the optimal fan-out being fairly small [112]. Thus, multiple recursion levels are not uncommon for large files, and statistics gathered on one level to limit skew effects on the next level are a realistic method for large files to control the performance penalties of uneven partitioning.

3. Disk Access

All query evaluation systems have to access base data stored in the database. For databases in the megabyte to terabyte range, base data are typically stored on secondary storage in form of rotating random-access disks. However, deeper storage hierarchies including optical storage, (maybe robot-operated) tape archives, and remote storage servers will also have to be considered in future high-functionality high-volume database management systems, e.g., as outlined by Stonebraker [272]. Research into database systems supporting and exploiting a deep storage hierarchy is still in its infancy.

3.1. File Scans

The first operator to access base data is the file scan, typically combined with a built-in selection facility. There is not much to be said about file scan except that it can be made very fast using read-ahead, particularly large-chunk (e.g., "track-at-a-crack") read-ahead. Efficient read-ahead requires contiguous file allocation, which is supported by many operating systems. Such contiguous disk regions are frequently called extents. The UNIX operating system does not provide contiguous files, and many database systems running on UNIX use "raw" devices instead, even though this means that the database management system must provide operating system functionality such as file structures, disk space allocation, and buffering.

The disadvantages of large units of I/O are buffer fragmentation and the waste of I/O and bus bandwidth if only individual records are required. Permitting different page sizes may seem to be a good idea, even at the added complexity in the buffer manager [43, 251], but this does not solve the problem of mixed sequential scans and random record accesses within one file. The common solution is to choose a middle-of-the-road page size, e.g., 8 KB, and to support multi-page read-ahead.

3.2. Associative Access using Indices

In order to reduce the number of accesses to secondary storage (which is relatively slow compared to main memory), most database systems employ associative search techniques in the form of indices that map key or attribute values to locator information with which database objects can be retrieved. The best-known and most-often used database index structure is the B-tree [14, 64]. A large number of extensions to the basic structure and its algorithms have been proposed, e.g., B⁺-trees for faster scans, for fast loading from a sorted file, increased fan-out and reduced depth by prefix and suffix truncation, B^{*}-trees for better space utilization in random insertions, and top-down B-trees for better locking behavior through preventive maintenance [119]. Interestingly, B-trees seem to be having a renaissance as a research subject, in particular with respect to improved space utilization [8], concurrency control [258], recovery [176], parallelism [236], and on-line creation of B-trees for very large databases [259, 260]. On-line reorganization and modification of storage structures, though not a new idea [200], is likely to become an important research topic within database research over the next few years as databases become larger and larger and are spread over many disks and nodes in parallel and distributed systems.

While most current database system implementations only use some form of B-trees, there is an amazing variety of index structures described in the literature, e.g., [16, 18, 118, 121-123, 133, 143-145, 149, 156, 168, 171, 172, 181, 197, 215, 230, 253]. One of the few multi-dimensional index structures actually implemented in a complete database management system are R-trees in Postgres [123, 270].

The large variety of index types can be described by the following six characteristics. First, does the index support range retrievals and ordered scans, or only exact-match equality lookups? This issue is the main difference

between sort-based indices such as B-trees and hash-based indices. Indices that support ordered key domains tend to have logarithmic insertion, deletion, and search costs, while index and storage structures based on hashing typically have constant average maintenance complexity.

Second, is the index structure static (e.g., ISAM) or dynamic (e.g., B-tree)? In other words, either the index structure allocates a fixed number of "buckets" when it is first created and resorts to overflow pages if buckets cannot hold all data items that logically belong in them, or it reorganizes itself incrementally as items are inserted and deleted.

Third, an index structure can be programmed to permit only single-attribute search or it can support multiple attributes in a hierarchical fashion. Such index implementations supporting composite keys, e.g., last name and first name, are still single-dimensional indices because the components of the composite key are ordered hierarchically into a major and a minor key. Note that this is an implementation detail; logically, multiple attributes are concatenated into one search key.

Fourth, does the index support only single-dimensional data or also data representing multiple dimensions? True multi-dimensional indices support all dimensions as equals, for example the x- and y-axes in a geometric application. Figure 10 shows a node in a quadtree [17, 230, 281], the simplest multi-dimensional index structure. The region represented by this node is divided along both dimensions, and each of the four subregions is represented by its own node.

The use of multi-dimensional indices for record-keeping applications has been largely ignored, despite their effectiveness for conjunctive queries. Consider a search for employees within both a certain age and salary range, e.g., $20 \leq \text{employee.age} \leq 30$ and $60 \leq \text{employee.salary} \leq 70$. In a system with single-dimensional indices only, either only one index can be utilized or two pointer lists must be intersected. Using only one index results in more data accesses since all employee records satisfying one clause must be inspected to evaluate the other clause. Intersection of two lists (or sets) can be quite expensive as will be seen in the subsequent section on binary matching, which includes relational join and set intersection. A two-dimensional index, on the other hand, permits more direct access to only the required employees because it supports both restriction clauses simultaneously. Presuming that the two dimensions represent age and salary, the relevant search region is shaded in Figure 10. It is clear

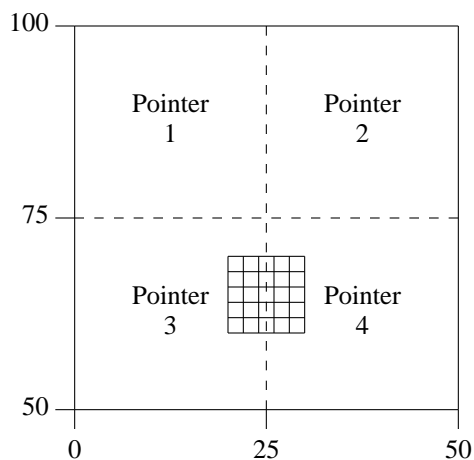


Figure 10. A Node in a Quadtree.

from the figure that more than one pointer from one node may need to be followed for some range queries, possibly at each index level. However, a multi-dimensional index promises to be still faster for large data sets than any method using single-dimensional indices. Of course, multi-dimensional indices can also be used for disjunctive queries, although their performance advantage is not as obvious for disjunctive as for conjunctive queries.

Fifth, do the indices support point data or range data? Range data have two data points in each dimension; the standard example is the case of two-dimensional rectangles. One method to support N -dimensional range data is to use an index structure for point data in $2 \times N$ dimensions. For example, the region shaded in Figure 10 could be represented in a four-dimensional index for point data as $x_1 = 20$, $x_2 = 30$, $x_3 = 60$, and $x_4 = 70$. One of the problems with this solution is that pairs of dimensions (the start and end value in the original dimensions) will likely be correlated, and the data structure may or may not include space- and search-efficient balancing mechanisms.

Sixth, most index implementations can be switched to accept or reject duplicate keys. Thus, indices are used to enforce the uniqueness constraint for identifying keys in database systems.

Finally, storage structures typically thought of as index structures may be used as primary structures to store actual data or as redundant structure ("access paths") that do not contain actual data but pointers to the actual data items in a separate data file. For example, Tandem's NonStop SQL system uses B-trees for actual data as well as for redundant index structures. In this case, a redundant index structure contains not absolute locations of the data items but keys used to search the primary B-tree. If indices are redundant structures, they can still be used to cluster the actual data items, i.e., the order or organization of index entries determines the order of items in the data file. Such indices are called *clustering* indices; other indices are called *non-clustering* indices. Clustering indices do not necessarily contain an entry for each data item in the primary file, but only one entry for each page of the primary file; in this case, the index is called *sparse*. Non-clustering indices must always be *dense*, i.e., there are the same number of entries in the index as there are items in the primary file.

The common theme for all index structures is that they associatively map some attribute of a data object to some locator information that can then be used to retrieve the actual data object. Typically, in relational systems, an attribute value is mapped to a tuple or record identifier (TID or RID). Different systems use different approaches, but it seems that most new designs do not firmly attach the record lookup to the index scan.

There are several advantages to separating index scan and record lookup. First, it is possible to scan the index only. For example, if only salary values are needed (e.g., to determine the count or sum of all salaries), it is sufficient to access the salary index only without actually retrieving the data records. The advantages are that (i) fewer I/O's are required (consider the number of I/O's for retrieving N successive index entries and those to retrieve N index entries plus N full records, in particular if the index is non-clustering [185], and (ii) the remaining I/O operations are basically sequential along the leaves of the index (at least for a B^+ -tree; other index types behave differently). The optimizers of several commercial relational products have recently been revised to recognize situations in which an index-only scan is sufficient. Second, if two or more indices can be used for a query, it may be more effective to union or intersect two RID lists obtained from two index scans than using only one index (algorithms for union and intersection are discussed below as join algorithms in the section on binary matching). Third, for non-clustering indices, sets of RID's can be sorted by physical location and the records retrieved very efficiently, reducing substantially the number of disk seeks and their seek distances. Obviously, the second and third advantages can be combined.

Record access performance for non-clustering indices can also be addressed without performing the entire index scan first (as required if all RID's are to be sorted) by using a "window" of RID's. Instead of obtaining one RID from the index scan, retrieving the record, getting the next RID from the index scan, etc., the lookup operator (sometimes called "functional join") could load N RID's, sort them into a priority heap, retrieve the most conveniently located record, get another RID, insert it into the heap, retrieve a record, etc. Thus, a functional join operator using a window always has N open references to items that must be retrieved, giving the functional join operator significant freedom to fetch items from disk efficiently. Of course, this technique works most effectively if no other transactions or operators use the same disk drive at the same time.

This idea has been generalized to assemble complex objects. In object-oriented systems, objects can contain pointers to (identifiers of) other objects or components, which in turn may contain further pointers, etc. If multiple objects and all their unresolved references can be considered concurrently when scheduling disk accesses, significant savings in disk seek times can be achieved [154].

3.3. Faster Storage Devices

Both for file scans and for index retrieval, the raw performance of the underlying storage (disk) system is crucial. Numerous ideas for faster disk access have been proposed, including RAM disks, i.e., simulation of disk drives using semi-conductor memory, and improvements in disk hardware speed, including physically smaller disks for faster seeks, faster disk rotation for shorter rotational latencies, and improved channel transfer rates. Other commonly used techniques are dual- or even quadruple-ported memory for concurrent transfer to or from multiple disks and processors, multiple paths from disks arms to I/O channels, and the insertion of RAM caches at various points in the data staging hierarchy, e.g., the disk controller or the drive. Write-only disk caches are an interesting proposal for cache use because most disk writes can be delayed in safe RAM and later piggy-backed transparently onto disk reads to the right cylinder or track, thus giving the illusion of instantaneous writes [257].

Recently, the idea of using multiple disk devices as a single, more powerful device has received considerable attention, and is now commonly known as Redundant Array of Inexpensive Disks (RAID) [209], although essentially the same idea had been proposed earlier as disk striping [226]. Put simply, by distributing the data and blocks of a file over multiple disk devices, higher transfer rates can be achieved. Furthermore, by using controlled redundancy, the mean time to failure as well as the mean time to repair can be improved substantially [30, 53, 65, 91, 97]. The idea can be further generalized from multiple disks to multiple nodes in a distributed system [269].

One problem with all disk array designs is ensuring proper placement of data on the disks to obtain the benefits of parallel I/O without incurring the additional overhead of controlling multiple devices for relatively simple and small requests [266, 267, 286]. Furthermore, RAID disk controller prohibit accessing the individual disks in a disk array separately, although it could be very useful for database query processing, in particular while merging multiple sorted runs or while partitioning large files. In sort- and hash-based query processing algorithms, access to individual disk drives could guarantee sequential I/O for temporary files including the merge input and partitioning output. Without solving these problems, it will be difficult to maximize disk array benefits; nonetheless, several vendors are making disk arrays commercially available, including the implementation of RAID control logic within device controllers to permit replacing a disk drive by a disk array.

While RAID-style disk arrays provide increased capacity, reliability, availability, and bandwidth for permanent data, they might not be the optimal I/O architecture for temporary data. The first reason is that writing parity blocks will be unnecessary (if a disk fails, most probably the entire transaction will abort) and reduce

processing performance (two writes instead of one). Second, important operations using multiple temporary files might not make optimal use of multiple disk drives if they are controlled as a single unit in a disk array. Consider sorting⁷ a very large file using a disk array, with each run file striped across all disk drives. Merging is faster using a disk array than using a single disk because many disk arms perform the required number of disk seeks. If each merge input file were stored contiguously on its own disk, however, merging could proceed without any disk seeks at all. Recall that our earlier discussion on cluster size and fan-in justifies and even recommends fairly small fan-ins; therefore, a modest disk array will permit one input file per disk. Thus, for run files in a sort operation, disk striping as built into disk array controllers is not optimal. It would be desirable if disk array controllers provided disk striping and parity blocks for permanent files but at the same time permitted sophisticated application software such as database management systems to set their own disk space allocation strategies. Current disk array controllers are examples of useful mechanisms packaged with fixed strategies deployed in all papers on the interface between operating systems and database systems, e.g. [263]. Space allocation for temporary files on multiple disk drives for sorting and partitioning in database query processing is an interesting issue that is not completely resolved at this time.

3.4. Buffer Management

I/O cost can be further reduced by caching data in an I/O buffer. A large number of buffer management techniques have been devised; we point out only a few references. Effelsberg surveys many of the buffer management issues, including those pertaining to issues of recovery, e.g., write-ahead logging [80]. In his survey paper on the interactions of operating systems and database management systems, Stonebraker pointed out that the "standard" buffer replacement policy, LRU (least recently used), is wrong for many database situations [263]. For example, a file scan reads a large set of pages but uses them only once, "sweeping" the buffer clean of all other pages, even if they might be useful in the future and should be kept in memory. Sacco and Schkolnick focused on the non-linear performance effects of buffer allocation to many relational algorithms, e.g., nested loops join [223, 224]. Chou and DeWitt combined these two ideas in their DBMIN algorithm which allocates a fixed number of buffer pages to each scan, depending on its needs, and uses a local replacement policy for each scan appropriate to its reference pattern [59, 60]. A recent study into buffer allocation is the study by Faloutsos et al. on using marginal gain for buffer allocation [85, 198]. A very promising research direction for buffer management in object-oriented database systems is the work by Palmer and Zdonik on saving reference patterns and using them to predict future object faults and to prevent them by prefetching the required pages [208].

The interactions of index retrieval and buffer management were studied by Sacco as well as Mackert and Lohman [185, 225], while several authors studied database buffer management and virtual memory provided by the operating system, e.g., [249, 263, 278].

⁷ For hash-based query processing algorithms, the same concern about disk arrays applies because merging and partitioning are quite similar, as discussed later in the section on the duality of sorting and hashing.

3.5. Physical Database Design

In order to minimize the I/O costs in a database system, it is important that (a) the data structures on disk permit efficient retrieval of only relevant data through effective access paths, and (b) data be arranged and placed on disk such that the I/O cost for relevant data is minimized. Both of these concerns are addressed in physical database design. Because physical database design is a wide area in which there are many studies on individual techniques but no comprehensive set of rules or guidelines on how to consider all of them in combination, we only list a number of choices to be made in physical database design and a few selected references:

- (1) index selection, i.e., indexed attributes and index structure, e.g. [118, 129, 156],
- (2) clustering, i.e., assignment of data items such as records to disk locations, e.g. [50, 56, 95, 132, 202, 255],
- (3) declustering (striping) over multiple disks or nodes [116, 209, 226, 269, 286],
- (4) replication for reliability and performance, e.g. [30, 65, 139, 209],
- (5) physical representation types for abstract data types,
- (6) management of derived information, e.g. [34, 142, 151, 211, 220, 277],
- (7) physical pointers to represent relationships, e.g. [44, 52, 216, 248],
- (8) data compression, e.g. [66, 106, 148, 179, 184, 199],
- (9) assignment of data to deeper archival storage levels, e.g. [272], and
- (10) automatic staging (rules) of data between storage levels, e.g. [96, 208].

It is important to recognize that most of these choices exist independently of the data model. For example, replication has beneficial availability and performance effects in network, relational, semantic, and object-oriented databases alike. On the other hand, clustering might have more effect in systems with logical or physical references, i.e., network and object-oriented databases, but master-detail clustering is used in relational system as well and is particularly effective in conjunction with index and pointer joins (to be discussed later in the section on binary matching). Thus, extensible database systems designed to build high-performance database systems must allow for a wide array of physical database design options.

While the number of choices for physical database design is confusing, the most significant source of complexity in physical database design is that many decisions are interdependent. For example, the optimal clustering of data items depends on whether or not replication of data items is supported. For example, when the clustering module cannot decide among two advantageous locations for a data item, it might choose to place a copy in each location. Replication can increase system performance by giving the optimizer or retrieval algorithm the choice of which copy to use; however, it must not be used too freely since it can also decrease overall performance if the cost of updating and maintaining multiple copies dominates their benefits [32, 142]. Moreover, the performance effects of replication are different depending on whether or not replicas of collections are declustered over multiple storage media as wholes or by means of partitioning, and whether or not the replicas are clustered equally. There is only limited research into making physical database design easier, e.g., [87, 152]. Considering the complexity of physical database design, automating physical database design in a comprehensive and extensible way seems to be an extremely fruitful area for database research, in particular in light of the added choices and complexity faced by the database implementor and administrator in extensible and object-oriented database management systems.

4. Aggregation and Duplicate Removal

Aggregation is a very important statistical concept to summarize information about large amounts of data. The idea is to represent a set of items by a single value or to classify items into groups and determine one value per group. Most database systems support aggregate functions for minimum, maximum, sum, count, and average (arithmetic mean). Other aggregates, e.g., geometric mean or standard deviation, are typically not provided, but may be constructed in some systems with extensibility features. Aggregation has been added to relational calculus and algebra and adds the same expressive power to each of them [166].

Aggregation is typically supported in two forms, called *scalar aggregates* and *aggregate functions* [82]. Scalar aggregates calculate a single scalar value from a unary input relation, e.g., the sum of the salaries of all employees. Scalar aggregates can easily be determined using a single pass over a data set. Some systems exploit indices, in particular for minimum, maximum, and count.

Aggregate functions, on the other hand, determine a set of values from a binary input relation, e.g., the sum of salaries for each department. Aggregate functions are relational operators, i.e., they consume and produce relations. Figure 11 shows the output of the query "count of employees by department." The "by-list" or grouping attributes are the key of the new relation, the Department attribute in this example.

Algorithms for aggregate functions require grouping, e.g., employee items may be grouped by department, and then one output item is calculated per group. This grouping process is very similar to duplicate removal in which equal data items must be brought together, compared, and removed. Thus, aggregate functions and duplicate removal are always implemented by the same module. There are only two differences between aggregate functions and duplicate removal. First, in duplicate removal, items are compared on all their attributes, but only on the attributes in the by-list of aggregate functions. Second, an identical item is immediately dropped from further consideration in duplicate removal whereas in aggregate functions some computation is performed before the second item of the same group is dropped. Both differences can easily be dealt with using a switch in an actual algorithm implementation. Because of their similarity, duplicate removal and aggregation are described and used interchangeably here.

In most existing commercial relational systems, aggregation and duplicate removal algorithms are based on sorting, following Epstein's work [82]. Since aggregation requires that all data be consumed before any output can be produced, and since main memories were significantly smaller 15 years ago when the prototypes of these systems were designed, these implementations used temporary files for output, not streams and iterator algorithms. However, there is no reason why aggregation and duplicate removal cannot be implemented using iterators.

Department	Count
Toy	3
Shoe	9
Hardware	7

Figure 11. Count of Employees by Department.

4.1. Aggregation Algorithms Based on Sorting

There are basically two types of algorithms for duplicate removal, one based on sorting and one based on hashing. Sorting will bring equal items together, and duplicate removal will then be easy. The cost of duplicate removal is dominated by the sort cost, and the cost of this naive duplicate removal algorithm based on sorting can be assumed to be that of the sort operation. For aggregation, items are sorted on their grouping attributes.

This simple method can be improved by detecting duplicate removal as possible, easily implemented in the routines that write run files during sorting. With early duplicate removal, a run file can never contain more items than the final output (because otherwise it would contain duplicates!), which may speed up the final merges significantly [29].

Since aggregation implemented by sorting using replacement selection can perform aggregation of matching items in the priority heap, it is only necessary that the output, not the input, fit into memory to make the entire sort and aggregation operation proceed without I/O to temporary run files. A corresponding optimization for sort-based aggregation using quicksort can be designed, but would be fairly cumbersome to implement.

As for any external sort operation, the optimizations discussed in the section on sorting, namely read-ahead using forecasting, merge optimizations, large cluster sizes, and reduced final fan-in for binary consumer operations, are fully applicable when sorting is used for aggregation and duplicate removal. However, to limit the complexity of the formulas, we derive I/O cost formulas without the effects of these optimizations.

The amount of I/O in sort-based aggregation is determined by the number of merge levels and the effect of early aggregation on each merge step. The total number of merge levels is unaffected by aggregation; in sorting with quicksort and without optimized merging, the number of merge levels is $L = \left\lceil \log_F (R / M) \right\rceil$ for input size R , memory size M , and fan-in F . In the first merge levels, the likelihood is negligible that items of the same group end up in the same run file, and we therefore assume that the sizes of run files are unaffected until their sizes would exceed the size of the final output. Runs on the first few merge levels are of size $M \times F^i$ for level i , and runs of the last levels have the same size as the final output. Assuming the output cardinality (number of items) is G -times less than the input cardinality ($G = R / O$), where G is called the average group size or the reduction factor, only the last $\left\lceil \log_F (G) \right\rceil$ merge levels, including the final merge, are affected by early aggregation because in earlier levels, more than G runs exist and items from each group are distributed over all those runs, giving a negligible chance of early aggregation.

In the first merge levels, all input items participate, and the cost for these levels can be determined without explicitly calculating the size and number of run files on these levels. In the affected levels, the size of the output runs is constant, equal to the size of the final output $O = R / G$, while the number of run files decreases by a factor equal to the fan-in F in each level. The number of affected levels that create run files is $L_2 = \left\lceil \log_F (G) \right\rceil - 1$; the subtraction of 1 is necessary because the final merge does not create a run file but the output stream. The number of unaffected levels is $L_1 = L - L_2$. The number of input runs is W / F^l on level l (recall the number of initial runs $W = R / M$ from the discussion of sorting). The total cost, including a factor 2 for writing and reading, is⁸

$$2 \times R \times L_1 + 2 \times O \times \sum_{l=L_1}^{L-1} W / F^l$$

$$= 2 \times R \times L_1 + 2 \times O \times W \times \left[1/F^{L_1} - 1/F^L \right] / \left[1 - 1/F \right].$$

For example, consider aggregating $R = 100$ MB input into $O = 1$ MB output (i.e., reduction factor $G = 100$) using a system with $M = 100$ KB memory and fan-in $F = 10$. Since the input is $W = 1,000$ times the size of memory, $L = 3$ merge levels will be needed. The last $L_2 = \log_F(G) - 1 = 1$ merge level into temporary run files will permit early aggregation. Thus, the total I/O will be

$$\begin{aligned} & 2 \times 100 \times 2 + 2 \times 1 \times 1000 \times \left[1/10^2 - 1/10^3 \right] / \left[1 - 1/10 \right] \\ & = 400 + 2 \times 1000 \times 0.009 / 0.9 = 420 \text{ MB} \end{aligned}$$

which has to be divided by the cluster size used and multiplied by the time to read or write a cluster to estimate the I/O time for aggregation based on sorting. Naive separation of sorting and subsequent aggregation would have required reading and writing the entire input file three times, for a total of 600 MB I/O. Thus, early aggregation realizes a 30% savings in this case.

Aggregate queries may require that duplicates be removed from the input set to the aggregate functions, e.g., using the SQL *distinct* keyword. If such an aggregate function is to be executed using sorting, early aggregation can be used only for the duplicate removal part. However, the sort order used for duplicate removal can be suitable to permit the subsequent aggregation as a simple filter operation on the duplicate removal's output stream.

4.2. Aggregation Algorithms Based on Hashing

Hashing can also be used for aggregation by hashing on the grouping attributes. Items of the same group (or duplicate items in duplicate removal) can be found and aggregated when inserting them into the hash table. Since only output items, not input items, are kept in memory, hash table overflow occurs only if the output does not fit into memory. However, if overflow does occur, the partition files (all partitioning files in any one recursion level) will basically be as large as the entire input because once a partition is being written to disk, no further aggregation can occur until the partition files are read back into memory.

The amount of I/O for hash-based aggregation depends on the number of partitioning (recursion) levels required before the output (not the input) of one partition fits into memory. This will be the case when partition files have been reduced to the size $G \times M$. Since the partitioning files shrink by a factor of F at each level (presuming hash value skew is absent or effectively counter-acted), the number of partitioning (recursion) levels is $\left\lceil \log_F(R/G/M) \right\rceil = \left\lceil \log_F(O/M) \right\rceil$ for input size R , output size O , reduction factor G , and fan-out F . The costs at each level are equal to the input file size R . The total I/O cost for hashing with overflow avoidance, including a factor of 2 for writing and reading, is

$$2 \times R \times \left\lceil \log_F(O/M) \right\rceil.$$

The last partitioning level may use hybrid hashing, i.e., it may not involve I/O for the entire input file. In that case, $L = \left\lceil \log_F(O/M) \right\rceil$ complete recursion levels involving all inputs records are required, partitioning the input into files of size $R' = R/F^L$. In each remaining hybrid hash aggregation, the size limit for overflow files is

⁸ Using $\sum_{i=0}^N a^i = \left[1 - a^{N+1} \right] / \left[1 - a \right]$ and $\sum_{i=K}^N a^i = \sum_{i=0}^N a^i - \sum_{i=0}^{K-1} a^i = \left[a^K - a^{N+1} \right] / \left[1 - a \right]$.

$M \times G$ because such an overflow file can be aggregated in memory. The number of partition files K must satisfy $K \times M \times G + (M - K \times C) \times G \geq R'$, meaning $K = \lceil (R'/G - M)/(M - C) \rceil$ partition files will be created. The total I/O cost for hybrid hash aggregation is

$$\begin{aligned}
 & 2 \times R \times L + 2 \times F^L \times \left[R' - (M - K \times C) \times G \right] \\
 & = 2 \times \left[R \times (L + 1) - F^L \times (M - K \times C) \times G \right] \\
 & = 2 \times \left[R \times (L + 1) - F^L \times \left[M - \lceil (R'/G - M)/(M - C) \rceil \times C \right] \times G \right].
 \end{aligned}$$

As for sorting, if an aggregate query requires that duplicates be removed from the input set to the aggregate function, the group size or reduction factor of the duplicate removal step determines the performance of hybrid hash duplicate removal. The subsequent aggregation can be performed as a simple filter operation on the duplicate removal output stream.

4.3. A Rough Performance Comparison

It is interesting to note that the performance of both sort- and hash-based aggregation is logarithmic and improves with increasing reduction factors. Figure 12 compares the performance of sort- and hash-based aggregation using the formulas developed above for 100 MB input data, 100 KB memory, clusters of 8 KB, fan-in or fan-out of 10, and varying group sizes or reduction factors. The output size is the input size divided by the group size.

It is immediately obvious in Figure 12 that sorting without early aggregation is not competitive because it does not limit the sizes of run files, confirming the results of Bitton and DeWitt [29]. The other algorithms all exhibit similar, though far from equal, performance improvements for larger reduction factors. Sorting with early aggregation improves once the reduction factor is large enough to affect not only the final but also previous merge

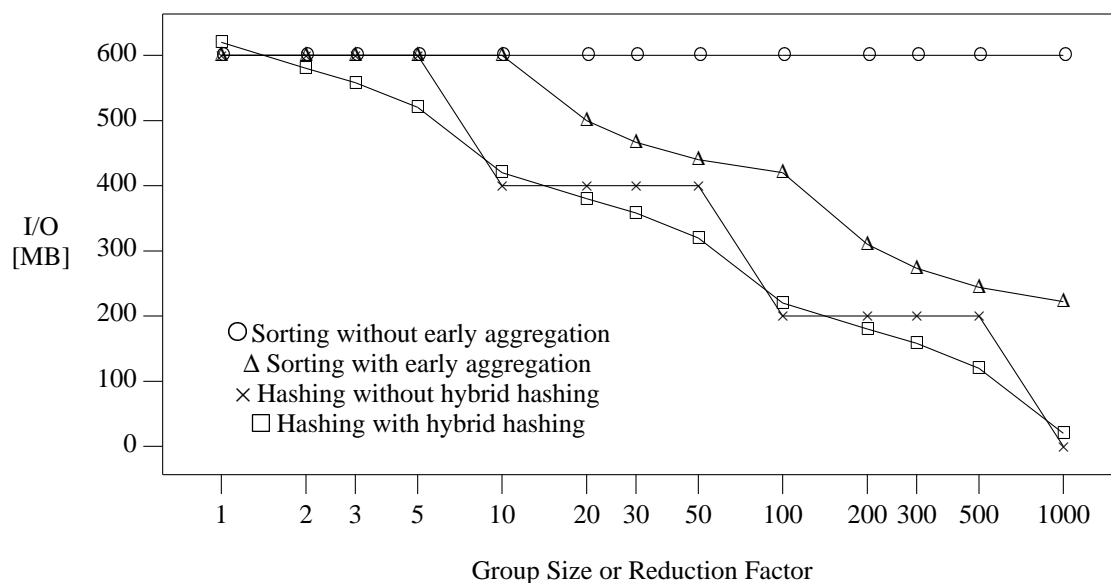


Figure 12. Performance of Sort- and Hash-Based Aggregation.

steps.

Hashing without hybrid hashing improves in steps as the number of partitioning levels can be reduced, with "step" points where $G = R / M / F^i$ for some i . Hybrid hashing exploits all available memory to improve performance, and generally outperforms overflow avoidance hashing. At points where overflow avoidance hashing shows a step, hybrid hashing has no effect and the two hashing schemes have the same performance.

While hash-based aggregation and duplicate removal seem superior in this rough analytical performance comparison, recall that the cost formula for sort-based aggregation does not include the effects of replacement selection or the merge optimizations discussed earlier in the section on sorting; therefore, Figure 12 shows an upper bound for the I/O cost of sort-based aggregation and duplicate removal. Furthermore, since the cost formula for hashing presumes optimal assignments of hash buckets to output partitions, the real costs of sort- and hash-based aggregation will be much more similar than they appear in Figure 12. The important point is that both their costs are logarithmic with the input size, improve with the group size or reduction factor, and are quite similar overall.

4.4. Additional Remarks on Aggregation

Some applications require multi-level aggregation. For example, a report generation language might permit a request like "sum (employee.salary by employee.id by employee.department by employee.division)" to create a report with an entry for each employee and a sum for each department and each division. In fact, specifying such reports conveniently was the driving design goal for the report generation language RPG. In SQL, this requires multiple cursors within an application program, one for each level of detail. This is very undesirable for two reasons. First, the application program performs essentially a three-way join, which should be provided by the database system. Second, the database system more likely than not executes the operations for these cursors independently from one another, resulting in three sort operations on the employee file instead of one.

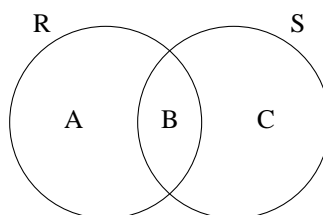
If complex reporting applications are to be supported, the query language should support direct requests (perhaps similar to the syntax suggested above), and the sort operator should be implemented such that it can perform the entire operation in a single sort and final pass over the sorted data. An analogous algorithm based on hashing can be defined; however, if the aggregated data are required in sort order, sort-based aggregation will be the algorithm of choice.

For some applications, exact aggregate functions are not required; reasonably close approximations will do. For example, exploratory (rather than final precise) data analysis is frequently very useful in "approaching" a new set of data [280]. In real-time systems, precision and response time may be reasonable tradeoffs. For database query optimization, approximate statistics are a sufficient basis for selectivity estimation, cost calculation, and comparison of alternative plans. For these applications, faster algorithms can be designed that rely either on a single sequential scan of the data (no run files, no overflow files) or on sampling [5, 135-137].

5. Binary Matching Operations

While aggregation is essential for *condensing* information, there are a number of database operations that *combine* information from two inputs, files, or sets and therefore are essential for database systems' ability to provide more than reliable shared storage and to perform inferences, albeit limited. A group of operators that all do basically the same task are called the one-to-one match operations here because an input item contributes to the output depending on the its match with one other item. The most prominent among these operations is the relational join. Mishra and Eich have recently written a survey of join algorithms [191], which includes an interesting analysis and comparison of algorithms based on the data items from the two inputs that are compared with one another. The other one-to-one match operations are left and right semi-joins, left, right, and symmetric outer-joins, anti-join, left and right anti-semi-join, intersection, union, left and right differences, and anti-difference. Figure 13 shows the basic principle underlying all these operations, namely separation of the matching and non-matching components of two sets, called R and S in the figure, and production of appropriate subsets⁹, possibly after some transformation and combination of records as in the case of a join. Since all these operations require basically the same steps and can be implemented with the same algorithms, it is logical to implement them in one general and efficient module. For simplicity, only join algorithms are discussed here. Moreover, we discuss algorithms for only one join attribute since the algorithms for multi-attribute joins (and their performance) are not different from those for single-attribute joins.

Since set operations such as intersection and difference will be used and must be implemented efficiently for any data model, this discussion is relevant to relational, extensible, and object-oriented database systems alike.



Output	Match on all Attributes	Match on some Attributes
A	Difference	Anti-semi-join
B	Intersection	Join, semi-join
C	Difference	Anti-semi-join
A, B		Left outer join
A, C	Symmetric difference	Anti-join
B, C		Right outer join
A, B, C	Union	Symmetric outer join

Figure 13. Binary One-to-One Matching.

⁹ If the sets R and S have different schemas as in relational joins, it might make sense to think of the set B as two sets B_R and B_S , i.e., the matching elements from R and S . This distinction permits a clearer definition of left semi-join and right semi-join, etc.

Furthermore, binary matching problems occur in some surprising places. Consider an object-oriented database system that uses a table to map logical object identifiers (OID's) to physical locations (record identifiers or RID's). Resolving a set of OID's to RID's can be regarded (as well as optimized and executed) as a semi-join of the mapping table and the set of OID's, and all conventional join strategies can be employed. Another example that can occur in a database management system for any data model is the use of multiple indices in a query: the pointer (OID or RID) lists obtained from the indices must be intersected (for a conjunction) or unioned (for a disjunction) to obtain the list of pointers to items that satisfy the whole query. Moreover, the actual lookup of the items using the pointer list can be regarded as a semi-join of the underlying data set and the list, as in Kooi's thesis and the Ingres product [169, 170] and a recent study by Shekita and Carey [248]. Finally, *path expressions* in object-oriented database systems such as "employee.department.manager.office.location" can frequently be interpreted, optimized, and executed as a sequence of one-to-one match operations using existing join and semi-join algorithms. Thus, even if relational systems were completely abolished and replaced by object-oriented database systems, set matching and join techniques developed in the relational context would continue to be important for the performance of database systems.

Most of today's database systems use only nested loops join and merge-join because an analysis performed in connection with the System R project determined that of all the join methods considered, one of these two always provided either the best or very close to the best performance [35, 36]. However, the System R study did not consider hash join algorithms, which are now regarded as more efficient in many cases.

For the I/O cost formulas given here, we assume that the left and right inputs have R and S pages, respectively, and that the memory size is M pages. We assume that the algorithms are implemented as iterators, and omit the cost of reading stored inputs and writing an operation's output from the cost formulas because both inputs and output may be iterators, i.e., these intermediate results are never written to disk, and because these costs are equal for all algorithms.

5.1. Nested Loops Join Algorithms

The simplest and, in some sense, most direct algorithm for binary matching is the nested loops join: for each item in one input (called the outer input), scan the entire other input (called the inner input) and find matches. The main advantage of this algorithm is its simplicity. Another advantage is that it can also compute a Cartesian product and any Θ -join of two relations, i.e., a join with an arbitrary two-relation comparison predicate. However, Cartesian products are avoided by query optimizers because their outputs tend to contain many data items that will eventually not satisfy a query predicate verified later in the query evaluation plan.

Since the inner input is scanned repeatedly, it must be stored in a file, i.e., a temporary file if the inner input is produced by a complex subplan. This situation does not change the cost of nested loops, it just replaces the first read of the inner input with a write.

Except for very small inputs, the performance of nested loops join is disastrous because the inner input is scanned very often, once for each item in the outer input. There are a number of improvements that can be made to this *naive nested loops join*. First, for one-to-one match operations in which a single match carries all necessary information, e.g., semi-join and intersection, a scan of the inner input can be terminated after the first match for an item of the outer input. Second, instead of scanning the inner input once for each item from the outer input, the inner input can be scanned once for each page of the outer input, an algorithm called *block nested loops join* [161]. Third, the performance can be improved further by filling all of memory except K pages with pages of the outer

input, and using the remaining K pages to scan the inner input and to save pages of the inner input in memory. Finally, scans of the inner input can be made a little faster by scanning the inner input alternately forwards and backwards, thus reusing the last page of the previous scan and therefore saving one I/O per scan. The I/O cost for this version of nested loop join is the product of the number of scans (determined by the size of the outer input) and the cost per scan of the inner input, plus K I/O's because the first inner scan has to scan or save the entire inner input. Thus, the total cost for temporary files or for scanning the inner input repeatedly is $\lceil R / (M - K) \rceil \times (S - K) + K$. This expression is minimized if $K = 1$ and $R \geq S$, i.e., the larger input should be the outer.

If the critical performance measure is not the amount of data read in the repeated inner scans but the number of I/O operations, more than one page should be moved in each I/O, even if more memory has to be dedicated to the inner input and less to the outer input, thus increasing the number of passes over the inner input. If C pages are moved in each I/O on the inner input, and $M - C$ pages for the outer input, the number of I/O's is $\lceil R / (M - C) \rceil \times (S / C) + 1$, which is minimized if $C = M / 2$. In other words, in order to minimize the number of large-chunk I/O operations, the cluster size should be chosen as half the available memory size [128].

Finally, there is the index nested loops join, which uses an index on the inner input's join attribute to replace file scans by index lookups. In principle, each scan of the inner input in naive nested loops join is used to find matches, i.e., to provide associativity. Not surprisingly, since all index structures are designed and used for the same general purpose, any index structure supporting the join predicate (such as eq , \leq , etc.) can be used for index nested loops join. The fastest indices for exact match queries are hash indices, but any index structure can be used, ordered or unordered (hash), single- or multi-attribute, single- or multi-dimensional. Therefore, indices on frequently used join attributes (keys and foreign keys in relational systems) may be useful. Index nested loops join is also used sometimes with indices built on-the-fly, i.e., indices built on intermediate query processing results. For complex queries, multi-way joins are sometimes written as a single module, i.e., a module that performs index look-ups into indices of multiple relations and joins all relations simultaneously. However, it is not clear how such a multi-way join implementation is superior to multiple index nested loops joins.

5.2. Merge-Join Algorithms

The second commonly used join method is the merge-join. It requires that both inputs are sorted on the join attribute. Merging the two inputs is similar to the merge process used in sorting. An important difference, however, is that one of the two merging scans (the one which is advanced on equality, usually called the inner input) must be backed up when both inputs contain duplicates of a join attribute value and the specific one-to-one match operation requires that all matches be found, not just one match. Thus, the control logic for merge-join variants for join and semi-join are slightly different. Some systems include the notion of "value packet," meaning all items with equal (join attribute) values [169, 170]. An iterator's *next* call returns a value packet, not an individual item, which makes the control logic for merge-join much easier. If (or after) both inputs have been sorted, the merge-join algorithm typically does not require any I/O, except when "value packets" are larger than memory¹⁰.

An input may be sorted because a stored database file was sorted, an ordered index was used, an input was sorted explicitly, or the input came from an operation that produced sorted output, e.g., another merge-join. The

¹⁰ See an earlier footnote on multiple granule sizes in the section on the architecture of query execution engines.

last point makes merge-join an efficient algorithm if items from multiple sources are matched on the same join attribute(s) in multiple binary steps because sorting intermediate results is not required for later merge-joins, which led to the concept of *interesting orderings* in the System R query optimizer [239]. Since set operations such as intersection and union can be evaluated using any sort order, as long as the same sort order is present in both inputs, the effect of interesting orderings for one-to-one match operators based on merge-join can always be exploited for set operations.

A combination of nested loops join and merge-join is the *heap-filter merge-join* [107]. It first sorts the smaller, inner input by the join attribute, and saves it in a temporary file. Next, it uses all available memory to create sorted runs from the larger, outer input using replacement selection. As discussed in the section on sorting, there will be about $W = R / (2 \times M) + 1$ such runs for outer input size R . These runs are not written to disk; instead, they are joined immediately with the sorted inner input using merge-join. Thus, the number of scans of the inner input is reduced to about one half when compared to block nested loops. On the other hand, when compared to merge-join, it saves writing and reading temporary files for the larger outer input.

Another derivation of merge-join is the hybrid join used in IBM's DB2 product [55], combining elements from index nested loops join, merge-join, and techniques joining sorted lists of index leaf entries. After sorting the outer input on its join attribute, hybrid join uses a merge algorithm to "join" the outer input with the leaf entries of a pre-existing B-tree index on the join attribute of the inner input. The result file contains entire tuples from the outer input and record identifiers (RID's, physical addresses) for tuples of the inner input. This file is then sorted on the physical locations and the tuples of the inner relation can then be retrieved from disk very efficiently. This algorithm is not entirely new as it is a special combination of techniques explored by Blasgen and Eswaran [35, 36] and Kooi [169]. Blasgen and Eswaran considered the manipulation of RID lists but concluded that either merge-join or nested loops join is the optimal choice in almost all cases; based on this study, only these two algorithms were implemented in System R [4] and subsequent relational database systems. Kooi's optimizer treated an index similarly to a base relation and the lookup of data records from index entries as a join; this naturally permitted joining two indices or an index with a base relation as in hybrid join.

5.3. Hash Join Algorithms

Hash join algorithms are based on the idea of building an in-memory hash table on one input (the smaller one, frequently called the *build input*) and then probing this hash table using items from the other input (frequently called the *probe input*). These algorithms have only recently found greater interest [41, 70-72, 90, 162, 164, 194, 203, 232, 244, 298]. One reason is that they work very fast, i.e., without any temporary files, if the build input does indeed fit into memory, independently of the size of the probe input. However, they require overflow avoidance or resolution methods for larger build inputs, and suitable methods were developed and experimentally verified only in the mid-1980's, most notably in connection with the Grace and Gamma database machine projects [72, 74, 90, 162]

In hash-based join methods, build and probe inputs are partitioned using the same partitioning function, e.g., the join key value modulo the number of partitions. The final join result can be formed by concatenating the join results of pairs of partitioning files. Figure 14 shows the effect of partitioning the two inputs of a binary operation

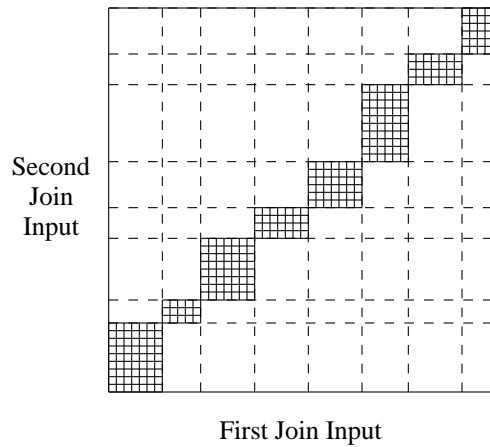


Figure 14. Effect of Partitioning for Join Operations.

such as join into hash buckets and partitions¹¹. Without partitioning, each item in the first input must be compared with each item in the second input; this would be represented by complete shading of the entire diagram. With partitioning, items are grouped into partition files, and only pairs in the series of small rectangles (representing the partitions) must be compared.

If a build partition file is still larger than memory, recursive partitioning is required. Recursive partitioning is used for both build and probe partitioning files using the same hash and partitioning functions. Figure 15 shows how both input files are partitioned together. The partial results obtained from pairs of partition files are concatenated to form the result of the entire match operation. Recursive partitioning stops when the build partition fits into memory. Thus, the recursion depth of partitioning for binary match operators depends only on the size of the build input (which therefore should be chosen to be the smaller input) and is independent of the size of the probe

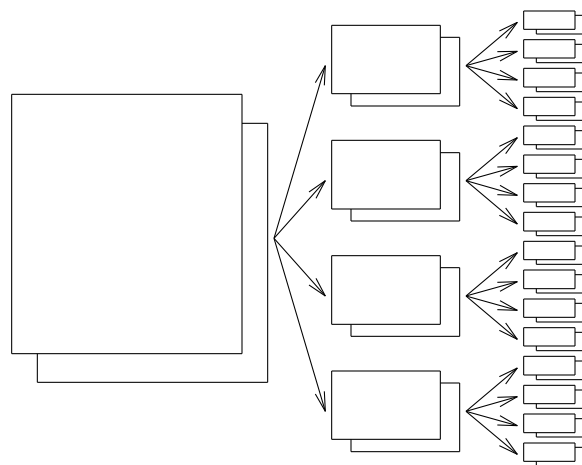


Figure 15. Recursive Partitioning in Binary Operations.

This figure was adapted from a similar diagram by Kitsuregawa et al. [162]. Mishra and Eich recently adapted and generalized it in their survey and comparison of relational join algorithms [191].

input. Compared to sort-based binary matching operators, i.e., variants of merge-join in which the number of merge levels is determined for each input file individually, hash-based binary matching operators are particularly effective when the input sizes are very different [41, 112].

The I/O cost for binary hybrid hash operations can be determined by the number of complete levels (i.e., levels without hash table) and the fraction of the input remaining in memory in the deepest recursion level. For memory size M , cluster size C , partitioning fan-out $F = \lfloor M / C - 1 \rfloor$, build input size R , and probe input size S , the number of complete levels is $L = \left\lceil \log_F (R / M) \right\rceil$, after which the build input partitions should be of size $R' = R / F^L$. The I/O cost for the binary operation is the cost of partitioning the build input divided by the size of the build input and multiplied by the sum of the input sizes. Adapting the cost formula for unary hashing discussed earlier, the total amount of I/O for a recursive binary hash operations is

$$2 \times \left[R \times (L + 1) - F^L \times \left(M - \lceil (R' - M + C) / (M - C) \rceil \times C \right) \right] / R \times (R + S)$$

which can be approximated with $2 \times \log_F (R / M) \times (R + S)$. In other words, the cost of binary hash operations on large inputs is logarithmic; the main difference to the cost of merge-join is that the recursion depth (the logarithm) depends only on one file, the build input, and is not taken for each file individually.

As for all operations based on partitioning, partitioning (hash) value skew is the main danger to effectiveness. When using statistics on hash value distributions to determine which buckets should stay in memory in hybrid hash algorithms, the goal is to avoid as much I/O as possible with the least memory "investment." Thus, it is most effective to retain those buckets in memory with few build items but many probe items or, more formally, the buckets with the smallest value for $r_i / (r_i + s_i)$ where r_i and s_i indicate the total size of a bucket's build and probe items.

5.4. Pointer-Based Joins

Recently, links between data items have found renewed interest, be it in object-oriented systems in the form of object identifiers (OID's) or as access paths for faster execution of relational joins. In a sense, links represent a limited form of precomputed results, similar to indices and join indices in particular, and have the usual cost vs. benefit tradeoff between query performance enhancement and maintenance effort. Shekita and Carey analyzed three pointer-based join methods based on nested loops join, merge-join, and hybrid hash join [248]. Presuming relations R and S , with a pointer to an S tuple embedded in each R tuple, the nested loops join algorithm simply scans through R and retrieves the appropriate S tuple for each R tuple. This algorithm is very reminiscent of unclustered index scans and performs similarly poorly for larger set sizes. Their conclusion on naive pointer-based join algorithms is that "it is unwise for object-oriented database systems to support only pointer-based join algorithms."

The merge-join variant starts with sorting R on the pointers (i.e., according to the disk address they point to) and then retrieves all S items in one elevator pass over the disk, reading each S page at most once. Again, this idea was suggested before for unclustered index scans, and variants similar to heap-filter merge-join [107] and complex object assembly using a window and priority heap of open references [154] can be designed.

The hybrid hash join variant partitions only relation R on pointer values, ensuring that R tuples with S pointers to the same page are brought together, and then retrieves S pages and tuples. Notice that the two relations' roles are fixed by the direction of the pointers, whereas for standard hybrid hash join the smaller relation should be the build input. Differently than standard hybrid hash join, relation S is not partitioned. This algorithm

performs somewhat faster than pointer-based merge-join if it keeps some partitions of R in memory and sorting writes all R tuples into runs before merging them.

Pointer-based join algorithms tend to outperform their standard, value-based counterparts in many situations, in particular if only a small fraction of S actually participates in the join and can be selected effectively using the pointers in R. Historically, due to the difficulty of correctly maintaining pointers (non-essential links), they were rejected as a relational access method in System R [48], and subsequently in basically all other system, perhaps with the exception of Kooi's modified Ingres [169, 170]. However, they were reevaluated and implemented in the Starburst project, both as a test of Starburst's extensibility and as a means of supporting "more object-oriented" modes of operation [126].

5.5. A Rough Performance Comparison

Figure 16 shows an approximate performance comparison using the cost formulas developed above for block nested loops join; merge-join with sorting both inputs, without optimized merging; hash join without hybrid hashing, bucket tuning, or dynamic destaging; and pointer joins with pointers from R to S and from S to R without grouping pointers to the same target page together. This comparison is not precise; its sole purpose is to give a rough idea of the relative performance of the algorithm groups, deliberately ignoring the many tricks used to improve and fine-tune the basic algorithms. The relation sizes vary; S is always ten times larger than R. The memory size is 100 KB, the cluster size is 8 KB, merge fan-in and partitioning fan-out are 10, and the number of R-records per cluster is 20.

It is immediately obvious in Figure 16 that nested loops join is unsuitable for medium-size and large relations, because the cost of nested loops join is proportional to the size of the Cartesian product of the two inputs. Both merge-join (sorting) and hash join have logarithmic cost functions; the sudden rise in merge-join and hash join cost around $R = 1000$ is due to the fact that additional partitioning or merging levels become necessary at that

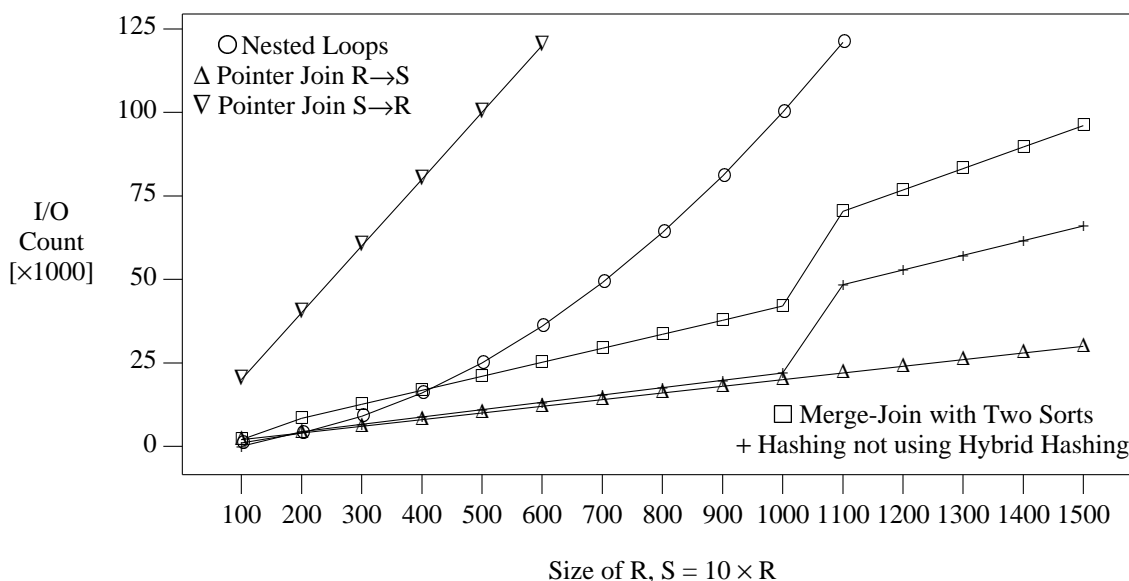


Figure 16. Performance of Alternative Join Methods.

point. The sort-based merge-join is not quite as fast as hash join because the merge levels are determined individually for each file, including the bigger S file, while only the smaller build relation R determines the partitioning depth of hash join. Pointer joins are competitive with their linear cost function, but only when the pointers are embedded in the smaller relation R. When S-records point to R-records, the cost of the pointer join is even higher than for nested loops join.

The important point of Figure 16 is to illustrate that pointer joins can be very efficient or very inefficient, that one-to-one match algorithms based on nested loops join are not competitive for medium-size and large inputs, and that sort- and hash-based algorithms for one-to-one match operations both have logarithmic cost growth. Of course, this comparison is quite naive since it uses only the simplest form of each algorithm. Thus, a comparison among alternative algorithms in a query optimizer must use the precise cost function for the available algorithm variant.

6. Universal Quantification

Universal quantification permits queries such as "find the students who have taken *all* database courses;" the difference to one-to-one match operations is that a student qualifies because his or her transcript matches an entire set of courses, not only one item as in an existentially quantified query (e.g., "find students who have taken a (at least one) database course") that can be executed using a semi-join. In the past, universal quantification has been largely ignored for four reasons. First, typical database applications, e.g., record-keeping and accounting applications, rarely require universal quantification. Second, it can be circumvented using a complex expression involving a Cartesian product. Third, it can be circumvented using complex aggregation expressions. Fourth, there seemed to be a lack of efficient algorithms.

The first reason will not remain true for database systems supporting logic programming, rules, and quantifiers, and algorithms for universal quantification will become more important. The second reason is valid; however, the substitute expressions are very slow to execute because of the Cartesian product. The third reason is also valid, but replacing a universal quantifier may require very complex aggregation clauses that are easy to "get wrong" for the database user. Furthermore, they might be too complex for the optimizer to recognize as universal quantification and to execute with a direct algorithm. The fourth reason is not valid; universal quantification algorithms can be very efficient (in fact, as fast as semi-join, the operator for existential quantification), useful for very large inputs, and easy to parallelize [102, 113]. In the remainder of this section, we discuss sort- and hash-based direct and indirect (aggregation-based) algorithms for universal quantification.

In the relational world, universal quantification is expressed with the universal quantifier in relational calculus and with the division operator in relational algebra. We will explain algorithms for universal quantification using relational terminology. The running example in this section uses the relations *Student* (*student-id*, *name*, *major*), *Course* (*course-no*, *title*), *Transcript* (*student-id*, *course-no*, *grade*) and *Requirement* (*major*, *course-no*) with the obvious key attributes. The query to find the students who have taken all courses can be expressed in relational algebra as

$$\pi_{student-id, course-no} Transcript \div \pi_{course-no} Course.$$

The projection of the *Transcript* relation is called the dividend, the projection of the *Course* relation the divisor, and the result relation the quotient. The quotient attributes are those attributes of the dividend that do not appear in the divisor. The dividend relation semi-joined with the divisor relation and projected on the quotient attributes, in the example the set of *student-id*'s of *Students* who have taken at least one course, is called the set of quotient

candidates here.

Some universal quantification queries seem to require relational division but actually do not. Consider the query for the students who have taken all courses required for their major. This query can be answered with a sequence of one-to-one match operations. A join of *Student* and *Requirement* projected on the *student-id* and *course-no* attributes minus the *Transcript* relation can be projected on *student-id*'s to obtain a set of students who have not taken all their requirements. An anti-semi-join of the *Student* relation with this set finds the students who have satisfied all their requirements. This sequence will have acceptable performance because its required set matching algorithms (join, difference, anti-semi-join) all belong to the family of one-to-one match operations, for which efficient algorithms are available as discussed in the previous section.

Division algorithms differ not only in their performance but also in how they fit into complex queries. Prior to the division, selections on the dividend, e.g., only *Transcript* entries with "A" grades, or on the divisor, e.g., only the database courses, may be required. Restrictions on the dividend can easily be enforced without much effect on the division operation, while restrictions on the divisor can imply a significant difference for the query evaluation plan. Subsequent to the division operation, the resulting quotient relation (e.g., a set of *student-id*'s) may be joined with the *Student* relation to obtain student *names*. Thus, obtaining the quotient in a form suitable for further processing (e.g., join or semi-join with a third relation) can be advantageous.

All universal quantification can be replaced by aggregations. For example, the example query about database courses can be re-stated as "find the students who have taken as many database courses as there are database courses." When specifying the aggregate function, it is important to count only database courses both in the dividend (the *Transcript* relation) and in the divisor (the *Course* relation). Counting only database courses might be easy for the divisor relation, but requires a semi-join of the dividend with the divisor relation to propagate the restriction on the divisor to the dividend if it is not known a priori whether or not referential integrity holds between the dividend's divisor attributes and the divisor, i.e., whether or not there are divisor attribute values in the dividend that cannot be found in the divisor. For example, *course-no*'s in the *Transcript* relation that do not pertain to database courses (and are therefore not in the divisor) must be removed from the dividend by a semi-join with the divisor. In general, if the divisor is the result of a prior selection, any referential integrity constraints known for stored relations will not hold, and must be explicitly enforced using a semi-join. Furthermore, in order to ensure correct counting, duplicates have to be removed from either input if the inputs are projections on non-key attributes.

There are four methods to compute the quotient of two relations, a sort-based and a hash-based direct method, and sort- and hash-based aggregation. Table 3 shows this classification of relational division algorithms. Methods for sort- and hash-based aggregation and the possible sort- or hash-based semi-join have already been discussed, including their variants for inputs larger than memory and their cost functions. Therefore, we focus here on the direct division algorithms.

	Sorting	Hashing
Direct	Naive division	Hash-division
Indirect by semi-join and aggregation	Sorting with duplicate removal, merge-join, sorting with aggregation	Hash-based duplicate removal, hybrid hash join, hash-based aggregation

Table 3. Classification of Relational Division Algorithms.

The sort-based direct method, proposed by Smith and Chang [254] and called *naive division* here, sorts the divisor input on all its attributes and the dividend relation with the quotient attributes as major and the divisor attributes as minor sort keys. It then proceeds with a merging scan of the two sorted inputs to determine which items belong in the quotient. Notice that the scan can be programmed such that it ignores duplicates in either input (in case those had not been removed yet in the sort) as well as dividend items that do not refer to items in the divisor. Thus, neither a preceding semi-join nor explicit duplicate removal step are necessary for naive division. The I/O cost of naive division is the cost of sorting the two inputs plus the cost of repeated scans of the divisor input.

Figure 17 shows two tables, a dividend and a divisor, already sorted for naive division (with the student id's and course no's replace by student name and course title for the sake of readability). Concurrent scans of the "Jane" tuples in the dividend and the entire divisor determines that "Jane" is not part of the quotient because she has not taken the "DB Readings" course. A continuing scan through the "Joe" tuples in the dividend and a new scan of the entire divisor includes "Joe" in the output of the naive division. The fact that "Joe" took some courses in addition to the ones in the divisor is ignored by the naive division algorithm.

The hash-based direct method, called *hash-division*, uses two hash tables, one for the divisor and one for the quotient candidates. While building a the divisor table, a unique sequence number is assigned to each divisor item. After the divisor table has been built, the dividend is consumed. For each quotient candidate, a bit map is kept with one bit for each divisor item. The bit map is indexed with the sequence numbers assigned to the divisor items. If a dividend item does not match with an item in the divisor table, it can be ignored immediately. Otherwise, a quotient candidate is either found or created and the bit corresponding to the matching divisor item is set. When the entire dividend has been consumed, the quotient consists of those quotient candidates for which all bits are set.

Figure 18 shows the two hash tables used in hash-division. The divisor table on the left contains all divisor tuples and associates a unique sequence number with each item. The quotient table on the right contains quotient candidates, obtained by projecting dividend tuples on their quotient attributes, and a bit map for each item indicating for which divisor tuples there has been a dividend tuple. The fact that "Jane" took only one DB course is indicated by the incompletely filled bit map. The AI course does not appear in either hash table because it was immediately determined that there was no AI course in the divisor table.

This algorithm can ignore duplicates in the divisor (using hash-based duplicate removal during insertion into the divisor table) and automatically ignores duplicates in the dividend as well as dividend items that do not refer to items in the divisor (e.g., the AI course in the example). Thus, neither prior semi-join nor duplicate removal are required. However, if both inputs are known to be duplicate-free, the bit maps can be replaced by counters.

Student	Course
Jane	Intro AI
Jane	Intro DB
Joe	DB Readings
Joe	Intro AI
Joe	Intro DB

Course
DB Readings
Intro DB

Figure 17. Sorted Inputs into Naive Division.

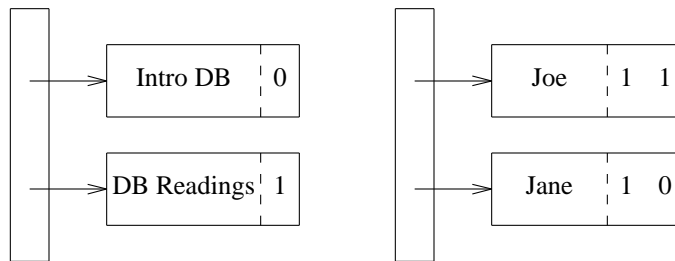


Figure 18. Divisor Table and Quotient Table in Hash-Division.

Furthermore, if referential integrity is known to hold, the divisor table can be omitted and be replaced by a single counter. Hash-division, including these variants, has been implemented in the Volcano query execution engine and has shown better performance than the other three algorithms [102, 113]. In fact, the performance of hash-division is almost equal to a hash-based join or semi-join of dividend and divisor relations (a semi-join corresponds to existential quantification), making universal quantification and relational division realistic operations and algorithms to use in database applications.

The aspect of hash-division that makes it an efficient algorithm is that the set of matches between a quotient candidate and the divisor is represented efficiently using a bit map. Bit maps are one of the standard data structures to represent sets, and just as bit maps can be used for a number of set operations, the bit maps associated with each quotient candidate can also be used for a number of operations similar to relational division. For example, Carlis proposed a generalized division operator called "HAS" that includes relational division as a special case [45]. The hash-division algorithm can easily be extended to compute quotient candidates in the dividend that match a majority or given fraction of divisor items as well as (with one more bit in each bit map) quotient candidates that do or do not match exactly the divisor items.

For real queries containing a division, consider the operation that frequently follows a division. In the example, a user is typically not really interested in *student-id*'s only but in information about the students. Thus, in many cases, relational division results will be used to select items from another relation using a semi-join. The sort-based algorithms produce their output sorted, which will facilitate a subsequent (semi-) merge-join. The hash-based algorithms produce their output in hash order; if overflow occurred, in no predictable order at all. However, both aggregation-based and direct hash-based algorithms use a hash table on the quotient attributes, which may be used immediately for a subsequent (semi-) join. It seems quite straightforward to use the same hash table for the aggregation and a subsequent join as well as to modify hash-division such that it removes quotient candidates from the quotient table that do not belong to the final quotient and then performs a semi-join with a third input relation.

If the two hash tables do not fit into memory, the divisor table or the quotient table or both can be partitioned and individual partitions held on disk for processing in multiple steps. In *divisor partitioning*, the final result consists of those items that are found in all partial results; the final result is the intersection of all partial results. For example, if the *Courses* relation in the example above are partitioned into undergraduate and graduate courses, the final result consists of the students who have taken all undergraduate courses and all graduate courses, i.e., those that can be found in the division result of each partition. In *quotient partitioning*, the entire divisor must be kept in memory for all partitions. The final result is the concatenation (union) of all partial results. For example, if

Transcript items are partitioned by odd and even *student-id*'s, the final results is the union (concatenation) of all students with odd *student-id* who have taken all courses and those with even *student-id* who have taken all courses. If warranted by the input data, divisor partitioning and quotient partitioning can be combined.

Hash-division can be modified into an algorithm for duplicate removal. Consider the problem of removing duplicates from a relation $R(X,Y)$ where X and Y are suitably chosen attribute groups. This relation can be stored using two hash tables, one storing all values of X (similar to the divisor table) and assigning each of them a unique sequence number, the other storing all values of Y and bit maps that indicate which X values have occurred with each Y value. Consider a brief example for this algorithm: Say relation $R(X,Y)$ contains 1,000,000 tuples, but only 100,000 tuples if duplicates were removed. Let X and Y be each 100 B long (total record size 200 B), and assume there are 4,000 unique values of each X and Y . For the standard hash-based duplicate removal algorithm, $100,000 \times 200$ B of memory are needed for duplicate removal without use of temporary files. For the redesigned hash-division algorithm, $2 \times 4,000 \times 100$ B are needed for data values, $4,000 \times 4$ B for unique sequence numbers, and $4,000 \times 4,000$ bits for bit maps. Thus, the new algorithm works efficiently with less than 3 MB of memory while conventional duplicate removal requires slightly more than 19 MB of memory, or seven times more than the duplicate removal algorithm adapted from hash-division. Clearly, choosing attribute groups X and Y to find attribute groups with relatively few unique values is crucial for the performance and memory-efficiency of this new algorithm. Since such knowledge is not available in most systems and queries (even though some efficient and helpful algorithms exist, e.g. [5]), optimizer heuristics for choosing this algorithm might be difficult to design and verify.

To summarize the discussion on universal quantification algorithms, aggregation can be used in systems that lack direct division algorithms, and hash-division performs universal quantification and relational division generally, i.e., it covers cases with duplicates in the inputs and with referential integrity violations, and efficiently, i.e., it permits partitioning and using hybrid hashing techniques similar to hybrid hash join, making universal quantification (division) as fast as existential quantification (semi-join). As will be discussed later, it can also be effectively parallelized.

7. Duality of Sorting and Hashing

In this section,¹² we conclude the discussion of individual query processing by outlining the many existing similarities and dualities of sort- and hash-based query processing algorithms as well as the points where the two types of algorithms differ. The purpose is to contribute to a better understanding of the two approaches and their tradeoffs. We try to discuss the approaches in general terms, ignoring whether the algorithms are used for relational join, union, intersection, aggregation, duplicate removal, or other operations. Where appropriate, however, we indicate specific operations.

Table 4 gives an overview of the features that correspond to one another. Both approaches permit in-memory versions for small data sets and disk-based versions for larger data sets. If a data set fits into memory, quicksort is the sort-based method to manage data sets while classic (in-memory) hashing can be used as hashing technique. It is interesting to note that both quicksort and classic hashing are also used in memory to operate on subsets after "cutting" an entire large data set into pieces. The cutting process is part of the *divide-and-conquer*

¹² Parts of this section have been derived from [112], which also provides experimental evidence for the relative performance of sort- and hash-based query processing algorithms.

Aspect	Sorting	Hashing
In-memory algorithm	Quicksort	Classic Hash
Divide-and-conquer paradigm	Physical division, logical combination	Logical division, physical combination
Large inputs	Single-level merge	Partitioning into overflow files
I/O Pattern	Sequential write, random read	Random write, sequential read
	Fan-in	Fan-out
I/O Optimization	Read-ahead, forecasting	Write-behind
	Double-buffering and striping for merge output	Double-buffering and striping for partitioning input
Very large inputs	Multi-level merge	Recursive overflow resolution
	Number of merge levels	Recursion depth
	Non-optimal final fan-in	Non-optimal hash table size
Optimizations	Merge optimizations	Bucket tuning
Better use of memory	Reverse runs & LRU	Hybrid hashing
	Replacement selection	?
	?	Single input in memory
Aggregation	Aggregation in replacement selection	Aggregation in hash table
Resource sharing	Eager merging	Depth-first partitioning
	Lazy merging	Breadth-first partitioning
Bit vector filtering	For inputs	For inputs and partitions in each recursion level
Interesting orderings	Merge-Join without sorting	N-way joins, hash-merging

Table 4. Duality of Sort- and Hash-Based Algorithms.

paradigm employed for both sorting and hashing. This important similarity of sorting and hashing has been observed before, e.g., by Bratbergsengen [41] and Salzberg [227]. There exists, however, an important difference. In the sort-based algorithms, a large data set is divided into subsets using a physical rule, namely into chunks as large as memory. These chunks are later combined using a logical step, merging. In the hash-based algorithms, the large data set is cut into subsets using a logical rule, by hash values. The resulting partitions are later combined using a physical step, i.e., by simply concatenating the subsets or result subsets. In other words, a single-level merge in a sort algorithm is a dual to partitioning in hash algorithms. Figure 19 illustrates this duality and the opposite directions.

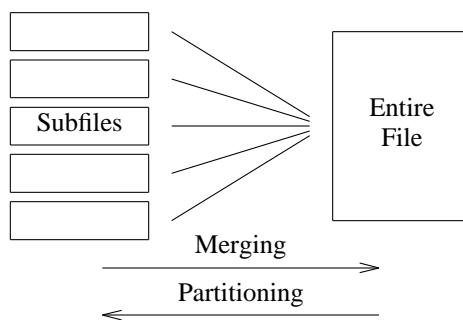


Figure 19. Duality of Partitioning and Merging.

This duality can also be observed in the behavior of a disk arm performing the I/O operations for merging or partitioning. While writing initial runs after sorting them with quicksort, the I/O is sequential. During merging, read operations access the many files being merged and require random I/O capabilities. During partitioning, the I/O operations are random, but when reading a partition later on, they are sequential.

For both approaches, sorting and hashing, the amount of available memory limits not only the amount of data in a basic unit processed using quicksort or classic hashing, but also the number of basic units that can be accessed simultaneously. For sorting, it is well known that merging is limited to the quotient of memory size and buffer space required for each run, called the merge *fan-in*. Similarly, partitioning is limited to the same fraction, called the *fan-out*, since the limitation is encountered while writing partition files.

In order to keep the merge process active at all times, many merge implementations use read-ahead controlled by forecasting, trading reduced I/O delays for a reduced fan-in. In the ideal case, I/O speed and processing (merging) speed match, and I/O delays for both the merge input and output are hidden by read-ahead and double-buffering, as mentioned earlier in the section on sorting. The dual to read-ahead during merging is write-behind during partitioning, i.e., keeping a free output buffer that can be allocated to an output file while the previous page for that file is being written to disk. Both read-ahead in merging and write-behind in partitioning are used to ensure that the processor never has to wait for the completion of an I/O operation. Another dual are double-buffering and striping over multiple disks for the output of sorting and the input of partitioning.

Considering the limitation on fan-in and fan-out, additional techniques must be used for very large data sets. Merging can be performed in multiple levels, each combining multiple runs into larger ones. Similarly, partitioning can be repeated recursively, i.e., partition files are re-partitioned, the results re-partitioned, etc., until the partition files fit into main memory. During merging, the runs grow in each level by a factor equal to the fan-in. For each recursion step, the partition files decrease in size by a factor equal to the fan-out. Thus, the number of levels during merging is equal to the recursion depth during partitioning. There are two exceptions to be made regarding hash value distribution and relative sizes of inputs in binary operations such as join; we ignore those for now and will come back to them later.

If merging is done in the most naive way, i.e., merging all runs of a level as soon as their number reaches the fan-in, the last merge on each level might not be optimal. Similarly, if the highest possible fan-out is used in each partitioning step, the partition file in the deepest recursion level might be smaller than memory, and less than the entire memory is used when processing these files. Thus, in both approaches the memory resources are not used optimally in the most naive versions of the algorithms.

In order to make best use of the final merge (which, by definition, includes all output items and is therefore the most expensive merge), it should proceed with the maximal possible fan-in. Making best use of the final merge can be ensured by merging fewer runs than the maximal fan-in after the end of the input file has been reached (as discussed in the earlier section on sorting). There is no direct dual in hash-based algorithms for this optimization. With respect to memory utilization, the fact that a partition file and therefore a hash table might actually be smaller than memory is the closest to a dual. Utilizing memory more effectively and using less than the maximal fan-out in hashing has been addressed in research on bucket tuning [164].

The development of hybrid hash algorithms [70, 244] was a consequence of the advent of large main memories that had led to the consideration of hash-based join algorithms in the first place. If the data set is only slightly larger than the available memory, e.g., 10% larger or twice as large, much of the input can remain in memory and is never written to a disk-resident partition file. To obtain the same effect for sort-based algorithms, if

the database system's buffer manager is sufficiently smart or receives and accepts appropriate hints, it is possible to retain some or all of the pages of the last run written in memory and thus achieve the same effect of saving I/O operations. This effect can be used particularly easily if the initial runs are written in reverse (descending) order and scanned backward for merging. However, if one does not believe in buffer hints or prefers to absolutely ensure these I/O savings, using a final memory-resident run explicitly in the sort algorithm and merging it with the disk-resident runs can guarantee this effect.

A well-known technique to improve sort performance is to generate runs twice as large as main memory using a priority heap for replacement selection [167], as discussed in the earlier section on sorting. If the runs' sizes are doubled, their number is cut in half. Therefore, merging can be reduced by some amount, namely $\log_F(2) = 1 / \log_2(F)$ merge levels. This optimization for sorting has no direct dual in the realm of hashing.

If two sort operations produce input data for a binary operator such as a merge-join and both sort operators' final merges are interleaved with the join, each final merge can employ only half the memory. In hash-based one-to-one match algorithms, only one of the two inputs resides in and consumes memory beyond a single input buffer, not both as in two final merges interleaved with a merge-join. This difference in the use of the two inputs is a distinct advantage of hash-based one-to-one match algorithms that does not have a dual in sort-based algorithms.

Interestingly, these two differences of sort- and hash-based one-to-one match algorithms cancel each other out. Cutting the number of runs in half (on each merge level, including the last one) by using replacement selection for run generation exactly offsets this disadvantage of sort-based one-to-one match operations.

Run generation using replacement selection has a second advantage over quicksort; this advantage has a direct dual in hashing. If a hash table is used to compute an aggregate function using grouping, e.g., sum of salaries by department, hash table overflow occurs only if the operation's *output* does not fit in memory. Consider, for example, the sum of salaries by department for 100,000 employees in 1,000 departments. If the 1,000 result records fit in memory, classic hashing (without overflow) is sufficient. On the other hand, if sorting based on quicksort is used to compute this aggregate function, the input must fit into memory to avoid temporary files.¹³ If replacement selection is used for run generation, however, the same behavior as with classic hashing is easy to achieve.

If an iterator interface is used for both its input and output, a sort operator can be divided into three distinct phases. First, input items are consumed and sorted into initial runs. Second, intermediate merging reduces the number of runs such that only one final merge step is left. Third, the final merge is performed on demand from the consumer of the sorted data stream. During the first phase, the sort iterator has to share resources, most notable memory and disk bandwidth, with its producer operators in a query evaluation plan. Similarly, the third phase must share resources with the consumers.

In many sort implementations, namely those using *eager merging*, the first and second phase interleave as a merge step is initiated whenever the number of runs on one level becomes equal to the fan-in. Thus, intermediate merge steps must cannot use all resources. In *lazy merging*, which starts intermediate merges only after the input has been consumed and all initial runs have been created, the intermediate merges do not share resources with other operators and can use the entire memory allocated to a query evaluation plan; thus, these merges can be more

¹³ A scheme using quicksort and avoiding temporary I/O in this case can be devised but would be extremely cumbersome; we do not know of any report or system with such a scheme.

effective in lazy than in eager merging.

Hash-based query processing algorithms exhibit similar three phases. First, the first partitioning step executes concurrently with the input operator or operators. Second, intermediate partitioning steps divide the partition files to ensure that they can be processed with hybrid hashing. Third, hybrid and in-memory hash methods process these partition files and produce output passed to the consumer operators. As in sorting, the first and third phases must share resources with other concurrent operations in the same query evaluation plan.

The standard implementation of hash-based query processing algorithms for very large inputs uses recursion, i.e., the original algorithm is invoked for each partition file (or pair of partition files). While conceptually simple, this method has the disadvantage that output is produced before all intermediate partitioning steps are complete. Thus, the operators that consume the output must allocate resources to receive this output, typically memory (e.g., a hash table). Further intermediate partitioning steps will have to share resources with the consumer operators, making them less effective. We call this direct recursive implementation of hash-based partitioning *depth-first* partitioning, and consider its behavior as well as its resource sharing and performance effects a dual to eager merging in sorting. The alternative schedule is *breadth-first* partitioning, which completes each level of partitioning before starting the next one. Thus, hybrid and in-memory hashing are not initiated until all partition files have become small enough to permit hybrid and in-memory hashing, and intermediate partitioning steps never have to share resources with consumer operators. Breadth-first partitioning is a dual to lazy merging, and it is not surprising that they are both equally more effective than depth-first partitioning and eager merging, respectively.

Bit vector filtering, which will be discussed later in more detail, can be used for both sort- and hash-based one-to-one match operations, although it has been used mainly for parallel hybrid hash join to-date. The basic idea of bit vector filtering is a large array of bits initialized by hashing items in the first input of a one-to-one match operator and used to detect items in the second input that cannot possibly have a match in the first input. In effect, bit vector filtering reduces the second input to the items that truly participate in the binary operation plus some "false passes" due to hash collisions in the bit vector. Bit vector filtering is equally effective in merge-join and hybrid hash join for reducing the size of the second input, and bit vector creating and the actual filtering process should be performed at the input side of the two sort operations. In recursive hybrid hash join, bit vector filtering can be used in each recursion level. Moreover, it can be used in both directions, i.e., to reduce the second input using a bit vector based on the first input and to reduce the first input using a bit vector based on the second input. The effectiveness of bit vector filtering increases in deeper recursion levels because the number of distinct data values in each partition decreases, thus reducing the number of hash collisions and false passes for the same size bit vector used in each recursion level. The same effect could be achieved for sort-based binary operations requiring multi-level sorting and merging, although to do so requires switching back and forth between the two sorts for the two inputs after each merge level. Not surprisingly, switching back and forth after each merge level would be the dual to the partitioning process of both inputs in recursive hybrid hash join.

The final entry in Table 4 concerns *interesting orderings* used in the System R query optimizer [239] and presumably other query optimizers as well. A strong argument in favor of sorting and merge-join is the fact that merge-join delivers its output in sorted order; thus, multiple merge-joins on the same attribute can be performed without sorting intermediate join results. For joining three relations, as shown in Figure 20, pipelining data from one merge-join to the next without sorting translates into a 3:4 advantage in the number of sorts compared to two joins on different join keys because the intermediate result O_1 does not need to be sorted. For joining N relations on the same key, only N sorts are required instead of $2 \times N - 2$ for joins on different attributes. Since set

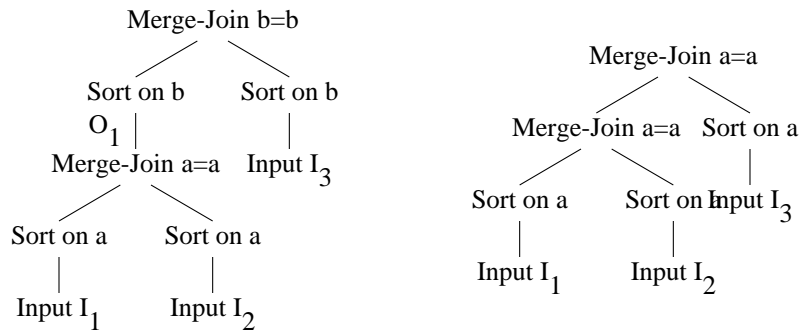


Figure 20. The Effect of Interesting Orderings.

operations such as the union or intersection of N sets can always be performed using a merge-join algorithm without sorting intermediate results, the effect of interesting orderings is even more important for set operations than for relational joins.

Hash-based algorithms tend to produce their outputs in a very unpredictable order, depending on the hash function and on overflow management. In order to take advantage of multiple joins on the same attribute (as well as intersections, etc.) similar to the advantage derived from interesting orderings in sort-based query processing, the equality of attributes has to be exploited during the logical step of hashing, i.e., during partitioning. In other words, such set operations and join queries can be executed effectively by a hash join algorithm that recursively partitions N inputs concurrently. The recursion terminates when $N - 1$ inputs fit into memory and the N^{th} input is used to probe $N - 1$ hash tables. Thus, the basic operation of this N -way join (intersection, etc.) is an N -way join of an N -tuple of partition files, not pairs as in binary hash join with one build and one probe file for each partition. Figure 21 illustrates recursive partitioning for a 3-way join.

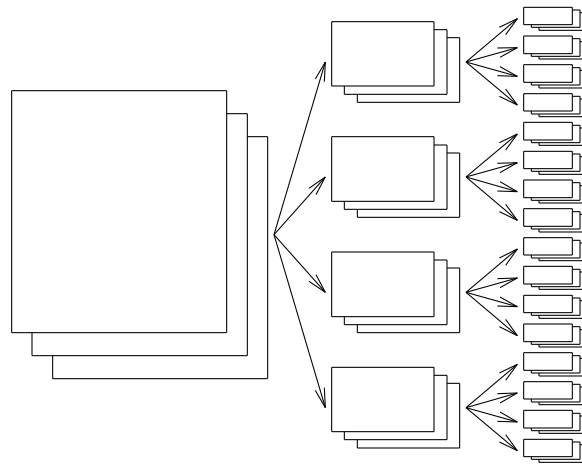


Figure 21. Partitioning in a Multi-Way Hash Join.

However, N-way recursive partitioning is cumbersome to implement, in particular if some of the "join" operations are actually semi-join, outer join, set intersection, union, or difference. Therefore, until a clean implementation method for hash-based N-way matching has been found, it might well be that this distinction, joins on the same or on different attributes, determines the right choice between sort- and hash-based algorithms for complex queries.

Another use of interesting orderings is the interaction of (sorted, B-tree) index scans and merge-join. While it has not been reported explicitly in the literature, the leaves and entries of two hash indices can be merge-joined just like those of two B-trees, provided the same hash function was used to create the indices. For example, it is easy to imagine "merging" the leaves (data pages) of two extendible hash indices [84], even if the key cardinalities and distributions are very different.

In summary, there exist many dualities between sorting using multi-level merging and recursive hash table overflow management. Two special cases exist which favor one or the other, however. First, if two join inputs are of different size (and the query optimizer can reliably predict this difference), hybrid hash join will outperform merge-join because only the smaller of the two inputs will determine what fraction of the input files will have to be written to temporary disk files during partitioning (or how often each record has to be written to disk during recursive partitioning), while each file determines its own disk I/O in sorting. In other words, sorting the larger of two join inputs is more expensive than writing a small fraction of that file to hash overflow files. This performance advantage of hashing grows with the relative size difference, not the absolute sizes, of the two inputs.

Second, if the hash function is very poor, e.g., because of a prior selection on the join attribute or a correlated attribute, hash partitioning can perform very poorly and create significantly higher costs than sorting and merge-join. If the quality of the hash function cannot be predicted or improved (tuned) dynamically, sort-based query processing algorithms are superior because they are less vulnerable to non-uniform data distributions. Since both cases, join of differently-sized files and skewed hash value distributions, are realistic situations in database query processing, we recommend that both sort- and hash-based algorithms be included in a query processing engine and chosen by the query optimizer according to the two cases above. If both cases arise simultaneously, i.e., a join of differently-sized inputs with unpredictable hash value distribution, the query optimizer has to estimate which one poses the greater danger to system performance and predictability and choose accordingly.

The important conclusion from these dualities is that neither the input sizes nor the memory size determine the choice between sort- and hash-based query processing algorithms. Instead, the choice should be governed by the *relative* sizes of the two inputs into binary operators and by the danger of skewed data or hash value distributions. Furthermore, because neither algorithm type outperforms the other in all situations, both should be available in a query execution engine for a choice to be made in each case by the query optimizer.

8. Execution of Complex Query Plans

When multiple operators such as aggregations and joins execute concurrently in a pipelined execution engine, physical resources such as memory and disk bandwidth must be shared by all operators. Thus, optimal scheduling of multiple operators and the division and allocation of resources in a complex plan are important issues.

In earlier relational execution engines, these issues were largely ignored for two reasons. First, only left-deep trees were used for query execution, i.e., the right (inner) input of a binary operator had to be a scan. In other words, concurrent execution of multiple subplans in a single query was not possible. Second, under the assumption that sorting was needed at each step and considering that sorting for non-trivial file sizes requires that the entire input be written to temporary files at least once, concurrency and the need for resource allocation were basically absent. Today's query execution engines consider more join algorithms permit extensive pipelining, e.g., hybrid hash join, and more complex query plans, including bushy trees. Moreover, today's systems support more concurrent users and use parallel processing capabilities. Thus, resource allocation for complex queries is of increasing importance for database query processing.

Some researchers have considered resource contention among multiple query processing operators with the focus on buffer management. The goal in these efforts was to assign disk pages to buffer slots such that the benefit of each buffer slot would be maximized, i.e., the number of I/O operations avoided in the future. Sacco and Schkolnick analyzed several database algorithms and found that their cost functions exhibit steps when plotted over available buffer space, and suggested that buffer space should be allocated at the low end of a step for the least buffer use at a given cost [223, 224]. Chou and DeWitt took this idea further by combining it with separate page replacement algorithms for each relation or scan, following observations by Stonebraker on operating system support for database systems [263], and with load control, calling the resulting algorithm DBMIN [59, 60]. Faloutsos et al. generalized this goal and used the classic economic concepts of decreasing marginal gain and balanced marginal gains for maximal overall gain [85, 198]. Zeller and Gray designed a hash join algorithm that adapts to the current memory and buffer contention each time a new hash table is built [298].

Schneider was the first to systematically examine execution schedules and costs for right-deep trees, i.e., query evaluation plans with multiple binary hash joins for which all build phases proceed concurrently or at least could proceed concurrently (notice that in a left-deep plan, each build phase receives its data from the probe phase of the previous join, limiting left-deep plans to two concurrent joins in different phases) [233, 234]. Among his most interesting findings are that through effective use of bit vector filtering (discussed later in its own subsection), memory requirements for right-deep plans might actually be comparable to those of left-deep plans [235]. This work has recently been extended to bushy plans interpreted and executed as multiple right-deep subplans by Chen et al. [54].

For binary matching iterators to be used in bushy plans, we have identified several concerns. First, some query processing algorithms include a point when all data are in temporary files on disk and no intermediate result data reside in memory. Such "stop" points can be used to switch efficiently between different subplans. For example, if two subplans produce and sort two merge-join inputs, stopping work on the first subplan and switching to the second one should be done when the first sort operator has all its data in sorted runs and only the final merge is left, but no output has been produced yet. Figure 22 illustrates this point in time. Fortunately, this timing can be realized naturally in the iterator implementation of sorting if opening runs for the final merge is done in the first call of the *next* procedure, not at the end of the *open* phase. A similar stop point is available in hash join when

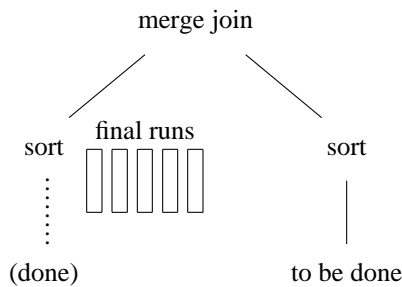


Figure 22. The Stop Point During Sorting.

using overflow avoidance. However, this point does not occur in the same way in hybrid hashing because hybrid hashing produces some output data before the memory contents (output buffers and hash table) can be discarded.

Second, implementations of hybrid hash join and other binary match operations should be parameterized to permit overflow avoidance as a run-time option to be chosen by the query optimizer. This dynamic choice will permit the query optimizer to force a stop point in some operators while using hybrid hash in most operations.

Third, binary operator implementations should include a switch that controls which subplan is initiated first. In Table 1 with algorithm outlines for iterators' *open*, *next*, and *close* procedures, the hash join *open* procedure executes the entire build input plan first before opening the probe input. However, there might be situations in which it would be better to *open* the probe input before executing the build input. If the probe input does not hold any resources such as memory between *open* and *next* calls, initiating the probe input first is not a problem. However, there are situations in which it creates a big benefit, in particular in parallel systems to be discussed later.

Fourth, if multiple operators are active concurrently, memory has to be divided among them. If two sorts produce input data for a merge-join which in turn passes its output into another sort using quicksort, memory should be divided proportionally to the sizes of the three files involved. We believe that for multiples sorts producing data for multiple merge-joins on the same attribute, proportional memory division will also work best. A comprehensive investigation of memory allocation for multiple unary and binary operators will need to consider many special cases for hybrid hash join, its data flow and control logic, and its complex cost function.

Fifth, in recursive hybrid hash join, the recursion levels should be executed level by level. In the most straightforward recursive algorithm, recursive invocation of the original algorithm for each output partition results

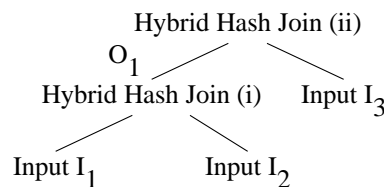


Figure 23. Plan for Three-Way Join.

in depth-first partitioning. The algorithm produces output as soon as the first leaf in the recursion tree is reached. However, if the operator that consumes the output requires memory as soon as it receives input, for example hybrid hash join (ii) in Figure 23 as soon as hybrid hash join (i) produces output, the remaining partitioning operations in the producer operator (hybrid hash join (i)) must share memory with the consumer operator (hybrid hash join (ii)), effectively cutting the partitioning fan-out in the producer in half. Thus, hash-based recursive matching algorithms should proceed in three phases — consuming input and initial partitioning, partitioning into files suitable for hybrid hash join, and final hybrid hash join for all partitions — with phase two completed entirely before phase three commences. We call this breadth-first partitioning as opposed to depth-first partitioning in the most straightforward recursive algorithms. Of course, depth-first partitioning in for the top-most operator in a query evaluation plan will provide faster response time, i.e., earlier delivery of the first data item.

Sixth, the allocation of resources other than memory, e.g., disk bandwidth and disk arms for seeking in partitioning and merging, is an open issue that should be addressed soon because the different improvement rates in CPU and disk speeds will increase the importance of disk performance for overall query processing performance. One possible alleviation of this problem might come from disk arrays configured exclusively for performance, not for reliability. On the other hand, disk arrays might not deliver the entire performance gain the large number of disk drives could provide if it is not possible to access specific disks within an array, particularly during partitioning and merging.

Finally, scheduling bushy trees in multi-processor systems is not entirely understood yet. While all considerations discussed above apply in principle, multi-processors permit truly concurrent execution of multiple subplans in a bushy tree. However, it is a very hard problem to schedule two or more subplans such that their result streams are available at the right times and at the right rates, in particular in light of the unavoidable errors in selectivity and cost estimation during query optimization [61, 146].

The last point, estimation errors, leads us to suspect that plans with 30 (or even 100) joins or other operations cannot be optimized completely before execution. Thus, we suspect that a technique reminiscent of Ingres Decomposition [291, 294] will prove the most effective one. One of the principal ideas of Ingres Decomposition is a repetitive cycle consisting of three steps. First, the next step is selected, e.g., a selection or join. Second, the chosen step is executed into a temporary table. Third, the query is simplified by removing predicates evaluated in the execution step and replacing one range variable (relation) in the query with the new temporary table. The justification and advantage of this approach is that all earlier selectivities are known for each decision because the intermediate results are materialized. The disadvantage is that data flow between operators cannot be exploited, resulting in a significant cost for writing and reading intermediate files. For very complex queries, we suggest modifying Decomposition to decide on and execute multiple steps in each cycle, e.g., 3-9 joins, instead of executing only one selection or join as in Ingres. Such a hybrid approach might very well combine the advantages of a priori optimization, namely in-memory data flow between iterators, and optimization with exactly known intermediate result sizes.

An optimization and execution environment even further tuned for very complex queries would anticipate possible outcomes of executing subplans and provide multiple alternative subsequent plans. Figure 24 shows the structure of such a dynamic plan for a complex query. First, subplan A is executed and statistics about its result are gathered while it is saved on disk. Depending on these statistics, either B or C is executed next. If B is chosen, one of D, E, and F is used to complete the query; or G or H if C had been chosen. Notice that each letter A–H can be an arbitrarily complex subplan, although probably not more than 10 operations due to the limitations of current

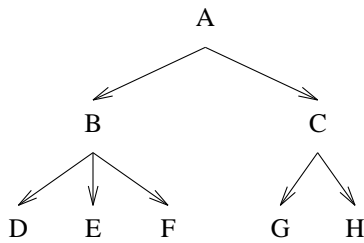


Figure 24. A Decision Tree of Partial Plans.

selectivity estimation methods. Unfortunately, realization of such sophisticated query optimizers will require further research, e.g., into determination of when separate cases are warranted and limitation of the exponential growth in the number of subplans.

Even for moderately complex subplans with about 10 operators, optimal resource allocation during execution time is an important issue as shown by Schneider for right-deep plans [233] and by Chen et al. for segmented right-deep plans [54]. In order to allocate resources such as memory, processors, and disks to operator, processes, and queries in any kind of query plan, a uniform and comparable measure of their value must be defined. We have defined a productivity measure for this purpose that we introduce here in several steps.¹⁴ It is important that this measure permits comparisons for sort- and hash-based query processing algorithms with their multiple phases (initial run generation and merging for sorting, partitioning and in-memory matching for hashing). Furthermore, the measure uses bandwidth for I/O and CPU processing, which may be measured in records or bytes per unit time. Algorithm bandwidth is a measure of the rate at which the algorithm can move records or bytes through the slowest part of the architecture (processor, disk, or network); hence, bandwidth considers both the architectural constraints and algorithm complexity and behavior. CPU performance measures based on MIPS numbers are ignored.

Our model compares the effectiveness of sorting with a fictitious external sort using quicksort for a single page of memory and binary merge levels only. The merge levels used in this sort are called the *required binary merge levels* of any sort and can be calculated as $RC = ld(R)$ for an input of R pages. The effectiveness of a real quicksort, replacement selection, or merge step is measured in how many of the required binary merge levels are performed in that processing step. If a real quicksort may use more memory than one page, say M pages, it makes several binary merge levels unnecessary. The size of the initial runs will be M pages, while the run size in the fictitious sort would reach M pages only after $ld(M)$ merge levels. Therefore, we say that the *replaced binary merge levels* of a quicksort with M pages is $ld(M)$. For initial-run generation based on replacement selection, in which the runs are about twice the size of memory, the replaced binary merge levels are $ld(M) + 1$. In a merge step with fan-in F , the size of the output run is F times larger than the input. In order to achieve the same increase in run size with binary merges, $ld(F)$ binary merge levels would be required. Therefore, we say that the *replaced binary merge levels* of a merge step with fan-in F is $ld(F)$. For hash-based query processing algorithms, the required and replaced binary merge levels are defined similarly as $ld(R)$ for a (build) input of size R , $ld(M)$ for a

¹⁴ This productivity measure has been defined in joined research with Diane Davison.

hash table of size M , and $ld(F)$ for a partitioning step with fan-out F .

The required binary merge levels multiplied with the data volume, i.e., the input size, are called the *required contribution*. The replaced binary merge levels multiplied with the data volume affected by a quicksort or a merge step (not the total input size) are called *performed contributions*. The important aspects of required and performed contributions are that the required contribution is independent of the sort algorithm used and that the performed contributions add up to the required contribution, no matter how a sort is performed, even if the memory allocation and merge fan-in are varied due to changing resource contention in a computer system. Thus, performed contributions are a direct measure of the value of a processing step to the completion of a sort, a database operator such as a merge-join, a process participating in a parallel query evaluation plan, or an entire query evaluation plan.

The performed contribution of an operation divided over its required amount of time defines the *productivity* of the operation. This productivity measure is equal to the replaced binary merge levels multiplied with the bandwidth of the operation. Thus, productivity relates an operation's value to the length of time it requires that resources be allocated to the operation.

For resource allocation, *marginal productivity* is defined for each resource as the additional productivity gained by one additional resource unit. It is assumed that marginal productivities are *positive*, i.e., that an additional unit of a resource increases the replaced binary merge levels or improves an operations bandwidth. This might require that some resources be "bundled," e.g., each disk drive "comes with" some amount of memory for read-ahead and write-behind buffer space. Otherwise, increasing the number of disk drives might decrease the amount of memory available for run generation or hash tables, thus decreasing the number of replaced binary merge levels. Moreover, it is assumed that marginal productivities are *monotonically decreasing*, i.e., that adding two unit of a resource improves an operations productivity by at most twice as much as one unit. Considering the logarithmic nature of replaced binary merge levels and the linear performance effect of CPU's and disks, this assumption seems justified.

Marginal gains can be used to perform resource allocations to operations, processes, and query plans. Among two contenders for a resource such as a unit of memory, the resource is allocated to the requestor which can derive the higher (marginal) gain from it. Thus, optimal resource allocation is achieved if the marginal gains are *balanced* among all contenders for each resource. The open issues to be resolved for this resource allocation scheme are (i) how interdependent allocation of different resources is, i.e., whether or not query processing algorithms can benefit from an additional unit of one resource such as processing power without an additional unit of another resource such as memory; (ii) how much resource allocation is affected by the logarithmic approximation of algorithms' costs; and (iii) whether or not resources can be "bundled" such that the two assumptions on marginal gains are justified. We are currently working on experimental answers to these questions.

9. Mechanisms for Parallel Query Execution

Considering that all high-performance computers today employ some form of parallelism in their processing hardware, it seems obvious that software written to manage large data volumes ought to be able to exploit parallel execution capabilities. In fact, we believe that five years from now it will be argued that a database management system without parallel query execution will be as handicapped in the market place as one without indices.

The goal of parallel algorithms and systems is to obtain speedup and scaleup, and speedup results are frequently used to demonstrate the accomplishments of a design and its implementation. Speedup considers additional hardware resources for a constant problem size; linear speedup is considered optimal. In other words, N times as many resources should solve a constant-size problem in $1/N$ of the time. Speedup can also be expressed as parallel efficiency, i.e., a measure of how close a system comes to linear speedup. For example, if solving a problem takes 1400 seconds on a single machine and 100 seconds on 16 machines, the speedup is slightly less than linear. The parallel efficiency is $(1 \times 1400) / (16 \times 100) = 87.5\%$.

A third measure for the success of a parallel algorithm based on Amdahl's law is the fraction of the sequential program for which linear speedup was attained, defined by $p = f \times s / d + (1 - f) \times s$ for sequential execution time s , parallel execution time p , and degree of parallelism d . For the example above, this fraction is $f = ((1400 - 100)/1400) / ((16 - 1)/16) = 99.05\%$. Notice that this measure is give much higher values than the parallel efficiency calculated earlier.

An alternative measure for a parallel system's design and implementation is scaleup in which the problem size is altered with the resources. Linear scaleup is achieved when N times as many resources can solve a problem with N times as much data in the same amount of time. Scaleup can also be expressed using parallel efficiency, but since speedup and scaleup are different, it should always be clearly indicated which parallel efficiency measure is being reported.

For query processing problems involving sorting or hashing in which multiple merge or partitioning levels are expected, the speedup can frequently be more than linear, or super-linear. Consider a sorting problem that requires two merge levels in a single machine. If multiple machines are used, the sort problem can be partitioned such that each machine sorts a fraction of the entire data amount. Such partitioning will, in a good implementation, result in linear speedup. If, in addition, each machine has its own memory such that the total memory in the system grows with the size of the machine, fewer than two merge levels will suffice, making the speedup super-linear.

9.1. Parallel vs. Distributed Database Systems

It might be useful to start the discussion of parallel and distributed query processing with a distinction of the two concepts. In the database literature, "distributed" usually implies "locally autonomous," i.e., each participating system is a complete database management system in itself, with access control, meta-data (catalogs), query processing, etc. In other words, each node in a distributed database management system can function entirely on its own, whether or not the other nodes are present or accessible. Each node performs its own access control, and cooperation of each node in a distributed transaction is voluntary. Examples of distributed (research) systems are R* [124, 279], distributed Ingres [83, 264], and SDD-1 [21, 219]. There are now several commercial distributed relational database management systems. Ozsu and Valduriez have discussed distributed database systems in much more detail [206, 207]. If the cooperation among multiple database systems is only limited, the system can be called a "federated" database system [250].

In parallel systems, on the other hand, there is only one locus of control. In other words, there is only one database management system that divides individual queries into fragments and executes the fragments in parallel. Access control to data is independent of where data objects currently reside in the system. The query optimizer and the query execution engine typically assume that all nodes in the system are available to participate in efficient execution of complex queries, and participation of nodes in a given transaction is either presumed or controlled by a global resource manager, but is not based on voluntary cooperation as in distributed systems. There are several parallel research prototypes, e.g., Gamma [72, 74], Bubba [39, 40], Grace [90, 162], and Volcano [104, 111], and products, e.g., Tandem's NonStop SQL [81, 297] and Teradata's TBC/1012 [196, 276].

Both distributed database systems and parallel systems have been designed in various kinds, which may create some confusion. Distributed systems can be either homogeneous, meaning that all participating database management systems are of the same type (the hardware and the operating system may even be of the same types), or heterogeneous, meaning that multiple database management systems work together using standardized interfaces but are internally different.¹⁵ Furthermore, distributed systems may employ parallelism, e.g., by pipelining datasets between nodes with the receiver already working on some items while the producer is still sending more. Parallel systems can be based on shared-memory (also called shared-everything), shared-disk (multiple processors sharing disks but not memory), distributed-memory (without sharing disks, also called shared-nothing), or hierarchical computer architectures. Stonebraker compared the first three alternatives using several aspects of database management, and came to the conclusion that distributed memory is the most promising database management system platform [265]. Each of these approaches has advantages and disadvantages; our belief is that the hierarchical architecture consisting of multiple clusters, each with multiple CPU's and disks and a large shared memory (see Figure 24), is the most general of these architectures and should be the target architecture for new database software development [109].

9.2. Forms of Parallelism

There are several forms of parallelism that are interesting to designers and implementors of database query processing systems. *Inter-query* parallelism is a direct result of the fact that most database management systems can service multiple requests concurrently. In other words, multiple queries (transactions) can be executing concurrently within a single database management system. In this form of parallelism, resource contention is of great concern, in particular contention for memory and disk arms.

The other forms of parallelism are all based on the use of algebraic operations on sets for database query processing, e.g., selection, join, and intersection. The theory and practice of exploiting other "bulk" types such as lists for parallel database query execution is only now developing. *Inter-operator* parallelism is basically pipelining, or parallel execution of different operators in a single query. For example, the iterator concept discussed earlier has also been called "synchronous pipelines" [210]; there is no reason not to consider asynchronous pipelines in which operators work independently connected by a buffering mechanism to provide flow control.

Inter-operator parallelism can be used in two forms, either to execute producers and consumers in pipelines, called *vertical inter-operator* parallelism here, or to execute independent subtrees in a complex, bushy query

¹⁵ In some organizations, two different database management systems may run on the same (fairly large) computer. Their interactions could be called non-distributed heterogeneous. However, since the rules governing such interactions are the same as for distributed heterogeneous systems, the case is usually ignored in research and system design.

evaluation plan concurrently, called *horizontal inter-operator* or *bushy* parallelism here. A simple example for bushy parallelism is a merge-join receiving its input data from two sort processes. The main problem with bushy parallelism is that it is hard or impossible to ensure that the two subplans start generating data at the right time and generate them at the right rates. Note that the right time does not necessarily mean the same time, e.g., for the two inputs of a hash join, and that the right rates are not necessarily equal, e.g., if two inputs of a merge-join have different sizes. Therefore, bushy parallelism presents too many open research issues and is hardly used in practice at this time.

The final form of parallelism in database query processing is *intra-operator* parallelism in which a single operator in a query plan is executed in multiple processes, typically on disjoint pieces of the problem and disjoint subsets of the data. This form, also called parallelism based on *fragmentation* or *partitioning*, is enabled by the fact that query processing focuses on sets. If the underlying data represented sequences, for example time series in a scientific database management system, partitioning into subsets to be operated upon independently would not be feasible or would require additional synchronization when putting the independently obtained results together.

Both vertical inter-operator parallelism and intra-operator parallelism are used in database query processing to obtain higher performance. Beyond the obvious opportunities for speedup and scaleup that these two concepts offer, they both have significant problems. Pipelining does not easily lend itself to load balancing because each process or processor in the pipeline is loaded proportionally to the amount of data it has to process. This amount cannot be chosen by the implementor or the query optimizer, and cannot be predicted very well. For intra-operator, partitioning-based parallelism, load balance and performance are optimal if the partitions are all of equal size; however, this can be hard to achieve as discussed earlier for partitioning as a hash table overflow management method.

9.3. Implementation Strategies

The purpose of the query execution engine is to provide mechanisms for query execution from which the query optimizer can choose — the same applies for the means and mechanisms for parallel execution. There are two general approaches to parallelizing a query execution engine, which we call the *bracket* and *operator models* and which are used, for example, in the Gamma and Volcano systems, respectively.

In the bracket model, there is a generic process template that can receive and send data and can execute exactly one operator at any point of time. A schematic diagram of a template process is shown in Figure 25 with two possible operators, join and aggregation. The code that makes up the generic template initiates the operator which then controls execution; network I/O on the receiving and sending sides is performed as a service to the operator on its request and initiation, and is implemented as procedures to be called by the operator. The number of inputs that can be active at any point of time is limited to two since there are only unary and binary operators in most database systems. The operator is surrounded by generic template code, which shields it from its environment, for example the operator(s) that produce its input and consume its output. For parallel query execution, many templates are executed concurrently in the system, using one process per template. Because each operator is written with the implicit assumption that this operator controls all activities in its process, it is not possible to execute two operators in one process without resorting to some thread or co-routine facility i.e., a second implementation level of the process concept.

In a query processing system using the bracket model, operators are coded in such a way that network I/O is their only means of obtaining input and delivering output (with the exception of scan and store operators). The

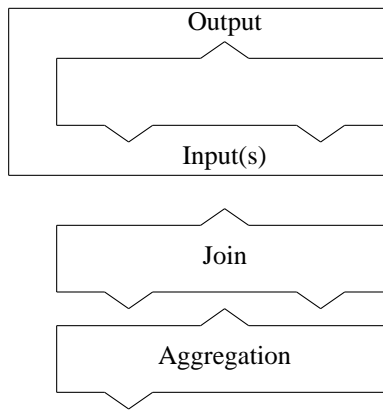


Figure 25. Bracket Model of Parallelization.

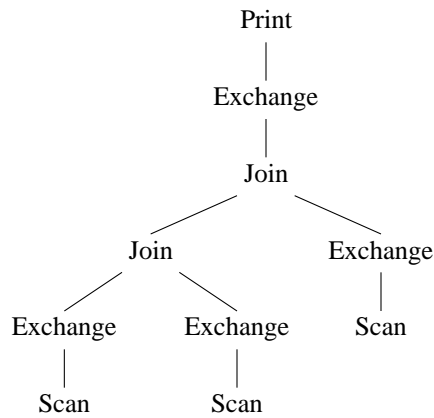


Figure 26. Operator Model of Parallelization.

reason is that each operator is its own locus of control and network flow control must be used to coordinate multiple operators, e.g., to match two operators' speed in a producer-consumer relationship. Unfortunately, this coordination requirement also implies that passing a data item from one operator to another always involves expensive inter-process communication system calls, even in the cases when an entire query is evaluated on a single machine (and could therefore be evaluated in a single process, without interprocess communication and operating system involvement) or when data do not need to be repartitioned among nodes in a network. An example for the latter is the three-way join query "joinCselAselB" in the Wisconsin Benchmark [76] which uses the same join attribute for both two-way joins. Thus, in queries with multiple operators (meaning almost all queries), interprocess communication and its overhead are mandatory rather than optional.

An alternative to the bracket model is the operator model. Figure 26 shows a possible parallelization of a three-way join plan using the operator model, i.e., by inserting "parallelism" operators into a sequential plan, called *exchange* operator in the Volcano system [104]. The exchange operator is an iterator like all other operators in the system with *open*, *next*, and *close* procedures; therefore, the other operators are entirely unaffected by the presence

of exchange operators in a query evaluation plan. Since it does not contribute to data manipulation but provides query processing control, we call it a *meta-operator*. Figure 27 shows the processes created by the exchange operators in the previous figure, with each circle representing a process. Note that this set of processes is only one possible parallelization, which makes sense if the joins are on the same join attributes. Furthermore, the degrees of data parallelism, i.e., the number of processes in each process group, can be controlled using an argument to the exchange operator.

There is no reason to assume that the two models differ significantly in their performance if implemented with similar care. Both models can be implemented with a minimum of control overhead and can be combined with any partitioning scheme for load balancing. The only difference with respect to performance is that the operator model permits multiple data manipulation operators such as join in a single process, i.e., operator synchronization and data transfer between operators with a single procedure call without operating system involvement. The important advantages of the operator model are that it permits easy parallelization of an existing sequential system as well as development and maintenance of operators and algorithms in a familiar and relatively simple single-process environment [111].

The bracket and operator models both provide pipelining and partitioning as part of pipelined data transfer between process groups. For most algebraic operators used in database query processing, these two forms of parallelism are sufficient. However, not all operations can be easily supported by these two models. For example, in a transitive closure operator, newly inferred data is equal to input data in its importance and role for creating further data. Thus, to parallelize a single transitive closure operator, the newly created data must also be partitioned like the input data. Neither bracket nor operator model immediately allow for this need. Hence, for transitive closure operators, intra-operator parallelism based on partitioning requires that the processes exchange data among themselves outside of the stream paradigm.

The transitive closure operator is not the only operation for which this restriction holds. Other examples include the complex object assembly operator described by Keller et al. [154] and operators for numerical optimizations as might be used in scientific databases. Both models, the bracket model and the operator model, could be

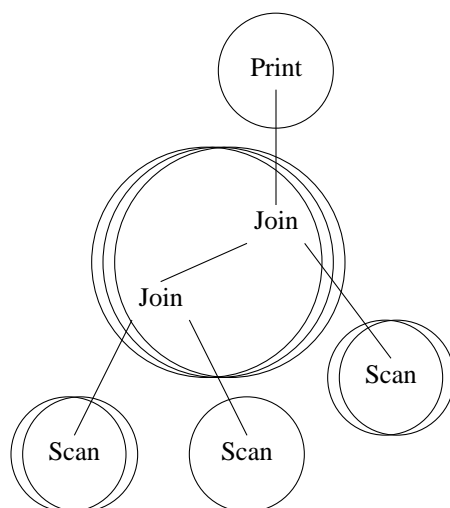


Figure 27. Processes Created by Exchange Operators.

extended to provide a general and efficient solution to intra-operator data exchange for intra-operator parallelism.

9.4. Load Balancing and Skew

For optimal speedup and scaleup, pieces of the processing load must be assigned carefully to individual processors and disks to ensure equal completion times for all pieces. In inter-operator parallelism, operators must be grouped to ensure that no one processor becomes the bottleneck for an entire pipeline. Balanced processing loads are very hard to achieve because intermediate set sizes cannot be anticipated with accuracy and certainty in database query optimization. Thus, no existing or proposed query processing engine relies solely on inter-operator parallelism. In intra-operator parallelism, data sets must be partitioned such that the processing load is nearly equal for each processor. Notice that in particular for binary operations such as join, equal processing loads can be different from equal-sized partitions.

There are several research efforts developing techniques to avoid skew or to limit the effects of skew in parallel query processing, e.g., [10, 79, 141, 165, 174, 175, 203, 241, 284, 285, 288, 289]. However, all of these methods have their drawbacks, for example additional requirements for local processing to determine quantiles.

Skew management methods can be divided into basically two groups. First, *skew avoidance* methods rely on determining suitable partitioning rules before data is exchanged between processing nodes or processes. For range-partitioning, quantiles can be determined or estimated from sampling the data set to be partitioned, from catalog data, e.g., histograms, or from a preprocessing step. Histograms kept on permanent base data have only limited use for intermediate query processing results, in particular if the partitioning attribute or a correlated attribute has been used in a prior selection or matching operation. However, for stored data they may be very beneficial. Sampling implies that the entire population is available for sampling because the first memory load of an intermediate result may be a very poor sample for partitioning decisions. Thus, sampling might imply that the data flow between operators be halted and an entire intermediate result be materialized on disk to ensure proper sampling and subsequent partitioning. However, if such a halt is required anyway for processing a large set, it can be used for both purposes. For example, while creating and writing initial run files without partitioning in a parallel sort, quantiles can be determined or estimated and used in a combined partitioning and merging step.

Second, *skew resolution* repartitions some or all of the data if an initial partitioning has resulted in skewed loads. Repartitioning is relatively easy in shared-memory machines, but can also be done in distributed-memory architectures, albeit at the expense of more network activity. Skew resolution can be based on both re-hashing in hash partitioning and quantile adjustment in range partitioning. Since hash partitioning tends to create fairly even loads and network bandwidth will increase in the near future within distributed-memory machines as well as in local- and wide-area networks, skew resolution is a reasonable method for cases in which a prior processing step cannot be exploited to gather the information necessary for skew avoidance as in the sort example above.

In their recent research into sampling for load balancing, Naughton et al. have shown that stratified random sampling can be used, i.e., samples are selected randomly not from the entire, distributed data set but from each local data set at each site, and that even small sets of samples ensure reasonably balanced loads [79, 241]. Their definition of skew is the quotient of sizes of the largest partition and the average partition, i.e., the sum of sizes of all partitions divided by the degree of parallelism. In other words, a skew of 1.0 indicates a perfectly even distribution. Figure 28 shows the required sample sizes per partition for various skew limits, degrees of parallelism, and confidence levels. For example, to ensure a maximal skew of 1.5 among 1,000 partitions with 95% confidence, 110 random samples must be taken at each site. Thus, relatively small samples suffice for reasonably safe skew

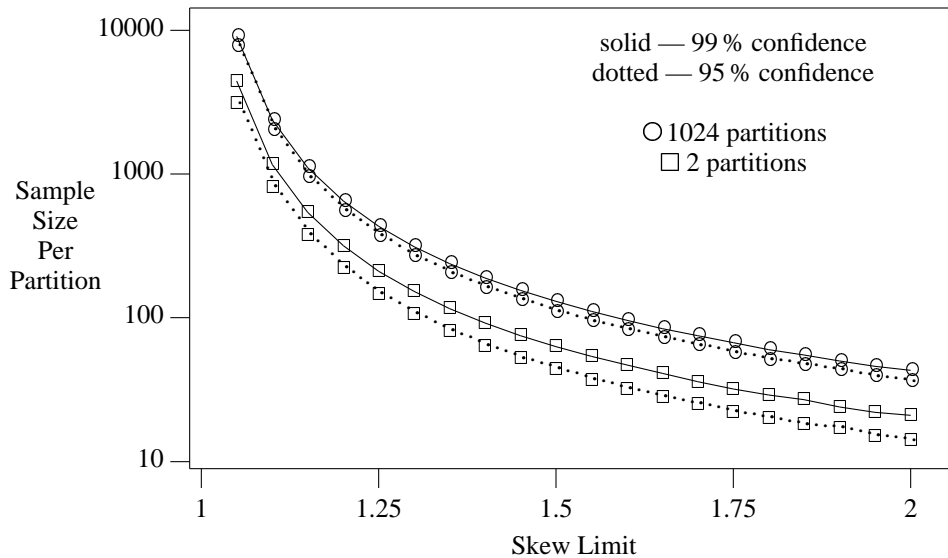


Figure 28. Skew Limit, Confidence, and Sample Size per Partition.

avoidance and load balancing, making precise methods unnecessary. Typically, only tens of samples per partition are needed, not several hundreds of samples at each site.

For allocation of active processing elements, i.e., CPUs and disks, the bandwidth considerations discussed briefly in the section on sorting can be generalized for parallel processes. In principal, all stages of a pipeline should be sized such that they all have bandwidths proportional to their respective data volumes in order to ensure that no stage in the pipeline become a bottleneck and slows the other ones down. The latency almost unavoidable in data transfer between pipeline stages should be hidden by the use of buffer memory equal in size to the product of bandwidth and latency.

9.5. Tuning a Parallel System

Beyond skew, the major impediments to effective use of parallelism in database query processing are *control overhead*, in particular initialization delays, *interference* on lower software and hardware levels, and *synchronization delays*.

Control overhead was one of the major obstacles to obtaining linear speedups in early database machine designs, for example DIRECT [37]. The main issue here is that the number of control messages should be linear to the number of processes, not with the number of data pages or items. If distribution of initiation messages is a problem and broadcasting is not available or not practical, a linear number of messages can be executed in a logarithmic number of single node-to-node messages using a tree-shaped propagation scheme, as proposed for example for a highly parallel implementation of the Gamma database machine [92, 93].

Another possible approach to reducing control overhead is to use not only primed processes (i.e., a pool of processes waiting until work packets arrive), but primed processes with primed connections. Instead of allocating processes and then establishing connections between them, a pool of process groups is used, each group already connected in a pattern typically found in parallel query processing. This idea clearly has a lot of promise for massively parallel systems but it requires determining the size of process groups and suitable "typical" connections.

Furthermore, the query optimizer must be designed to exploit such preconnected process groups.

Interference can occur both on a hardware level and in lower software levels. The prime example for hardware interference is bus contention in shared-memory multi-processors, the reason for the limited scalability of single-bus shared-memory architectures. Lower software levels can also introduce bottlenecks and contention, for example a buffer manager with its residency lookup table or a bit map for disk page allocation. A study of software interference in a more general sense has been discussed by Fontenot [88].

As an example, Figure 29, taken from [114], shows the effect of tuning and removing interference among processes on a shared-memory machine executing a parallel sort of 100 MB (10^6 records of 100 B) with an equal number of processors and disks in each measurement. The time measurements are shown using solid lines and refer to the labels on the left. The speedups are shown with dashed lines and refer to the labels on the right. The initial times and speedups are marked with \square 's while the final ones are marked with Δ 's. The ideal, linear speedup is shown by the dotted line. It is immediately obvious from the solid lines that the final times are significantly lower than the initial ones, demonstrating the effect of the tuning measures. For two to eight processors and disks, the observed performance improvements by a factor slightly more than two are largely due to increased cluster size and reduced I/O cost. Beyond eight processors and disks, the dashed lines indicate that the modifications and adjustments also improved the speedup which had been completely unsatisfactory with the initial software. For sixteen processors and disks, the fully tuned software performed about $3\frac{1}{2}$ times better than the original version. A comparison of the dashed and dotted lines shows very close to linear speedup with the fully tuned software. Thus, tuning improved the parallel behavior as well as the absolute performance. The main reason why the speedup could be improved was the removal of a latch contention bottleneck from the buffer manager, i.e., interference of multiple processes using a shared resource.

Synchronization delays can also be induced by hardware and software. Software delays may occur if multiple processes in a pipeline operate with different throughput and require the standard producer-consumer synchronization. Synchronizing multiple caches can introduce significant delays, in particular since today's cache sizes can almost be thought of as making a shared-memory machine into a distributed-memory machine. Thus,

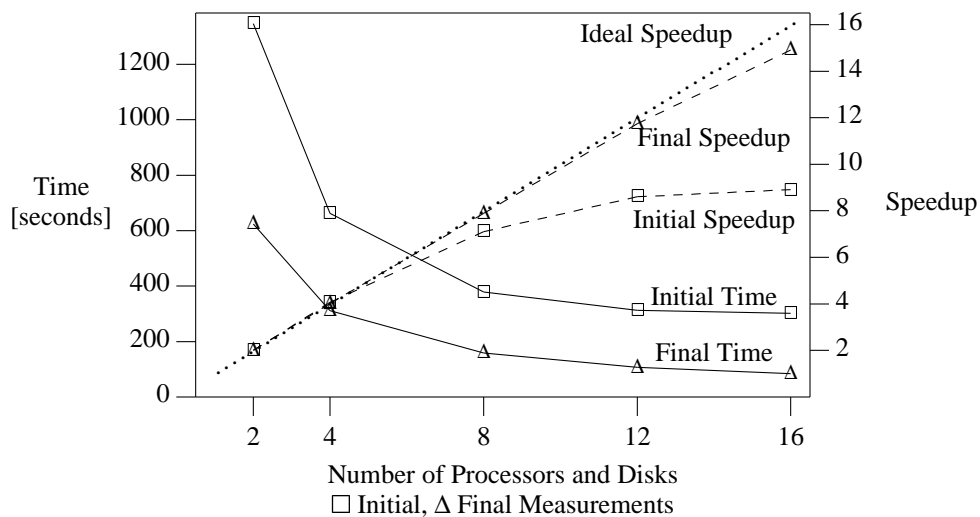


Figure 29. Tuning Effectiveness for a Parallel Algorithm.

effective load balancing through a single run-queue, usually considered one of the advantages of shared-memory machines, might be counter-productive if processes (and their cache residency set) tend to migrate too much or too often. Our recommendation is to make processes "sticky," i.e., to try but not to force process-to-processor affinity [114].

9.6. Architectures and Architecture-Independence

¹⁶Many database research projects have investigated hardware architectures for parallelism in database systems. Stonebraker compared shared-nothing (distributed-memory), shared-disk (distributed-memory with multiported disks), and shared-everything (shared-memory) architectures for database use based on a number of issues including scalability, communication overhead, locking overhead, and load balancing [265]. His conclusion was that shared-everything excels in none of the points considered, shared-disk introduces too many locking and buffer coherency problems, and that shared-nothing has the significant benefit of scalability to very high degrees of parallelism. Therefore, he concluded that overall shared-nothing is the preferable architecture for database system implementation.

Bhide and Stonebraker compared architectural alternatives for transaction processing [27, 28] and concluded that a shared-everything (shared-memory) design provides best performance. To achieve high performance, reliability, and scalability, Bhide suggested considering shared-nothing (distributed-memory) machines with shared-everything parallel nodes. The same idea is mentioned in equally general terms by Pirahesh et al. [210] and Boral et al. [40], but none of these authors elaborate on the idea's generality or potential. Kitsuregawa and Ogawa's new database machine SDC uses multiple shared-memory nodes (plus custom hardware such as the Omega network and a hardware sorter) [165], although the effect of the hardware design on operators other than join is not evaluated in [165].

Customized parallel hardware was investigated but largely abandoned after Boral and DeWitt's influential analysis [38] that compared CPU and I/O speeds and their trends and concluded that I/O, not processing, is the most likely bottleneck in future high-performance query execution. Subsequently, both Boral and DeWitt embarked on new database machine projects, Bubba and Gamma, that executed customized software on standard processors with local disks [40, 74]. For scalability and availability, both projects used distributed-memory hardware with single-CPU nodes, and investigated scaling questions for very large configurations.

The XPRS system, on the other hand, has been based on shared memory [134, 267, 268]. Its designers believe that modern bus architectures can handle up to 2,000 transactions per second. Shared-memory architectures provide automatic load balancing and faster communication than shared-nothing machines and are equally reliable and available for most errors, i.e., media failures, software, and operator errors [115]. However, we believe that attaching 250 disks to a single machine as necessary for 2,000 transactions per second [267] requires significant special hardware, e.g., channels or I/O processors, and it is quite likely that the investment for such hardware can have greater impact on overall system performance if spent on general-purpose CPU's or disks. Without such special hardware, the performance limit for shared-memory machines is probably much lower than 2,000 transactions per second. Furthermore, there already are applications that require larger storage and access capacities.

¹⁶ Much of this section has been derived from [109].

Richardson et al. [213] performed an analytical study of parallel join algorithms on multiple shared-memory "clusters" of CPU's. They assumed a group of clusters connected by a global bus with multiple microprocessors and shared memory in each cluster. Disk drives were attached to the busses within clusters. However, their analysis suggested that the best performance is obtained by using only one cluster, i.e., a shared-memory architecture. We contend that their results are due to their parameter settings, in particular small relations (typically 100 pages of 32 KB), slow CPU's (e.g., 5 μ sec for a comparison, about 2–5 MIPS), a slow global network (a bus with typically 100 Mbit/sec), and a modest number of CPU's in the entire system (128). It would be very interesting to see the analysis with larger relations (e.g., 1–10 GB), more and faster CPU's (e.g., 1,000 \times 50 MIPS), and a faster network e.g., a modern hypercube or mesh with hardware routing. In such machines, multiple clusters could be the better choice. On the other hand, communication between clusters will remain a significant expense. Wong and Katz developed the concept of "local sufficiency" [292] that might provide guidance in declustering and replication to reduce data movement between nodes. Other work on declustering and limiting declustering includes [65, 86, 94, 138, 140].

Finally, there are several hardware designs that attempt to overcome the shared-memory scaling problem, e.g., the DASH project [3], the Wisconsin Multicube [98], and the Paradigm project [57]. However, these designs follow the traditional separation of operating system and application program. They rely on page or cache-line faulting and do not provide typical database concepts such as read-ahead and dataflow. Lacking separation of mechanism and policy in these designs almost makes it imperative to implement dataflow and flow control for database query processing within the query execution engine. At this point, none of these hardware designs has been experimentally tested for database query processing.

New software systems designed to exploit parallel hardware should be able to exploit both the advantages of shared memory, namely efficient communication, synchronization, and load balancing, and of distributed memory, namely scalability to very high degrees of parallelism and reliability and availability through independent failures. Figure 30 shows a general hierarchical architecture, which we believe combines these advantages. The important

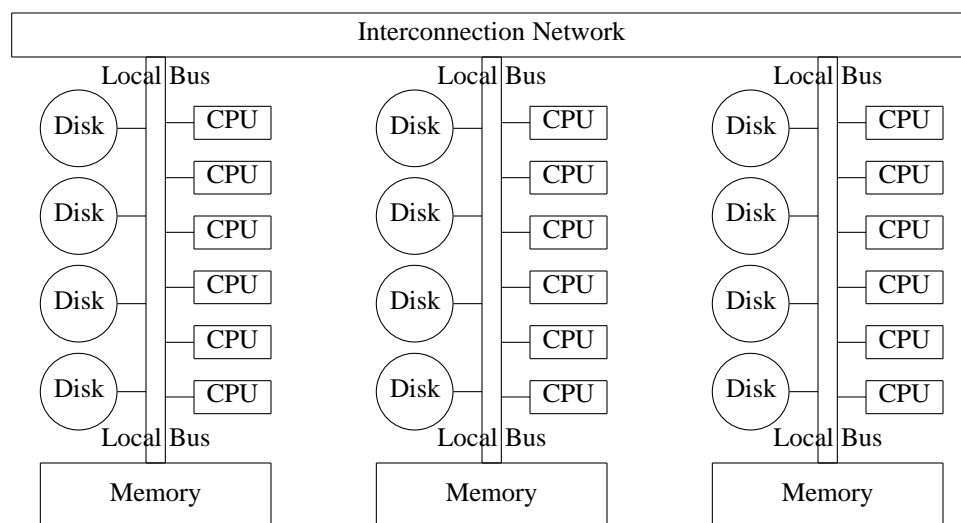


Figure 30. A Hierarchical-Memory Architecture.

point is the combination of local busses within shared-memory parallel machines and a global interconnection network between machines. The diagram is only a very general outline of such an architecture; many details are deliberately left out and unspecified. The network could be implemented using a bus such as an ethernet, a ring, a hypercube, a mesh, or a set of point-to-point connections. The local busses may or may not be split into code and data or by address range to obtain less contention and higher bus bandwidth and hence higher scalability limits for the use of shared memory. Design and placement of caches, disk controllers, terminal connections, and local- and wide-area network connections are also left open. Tape drives or other backup devices would be connected to local busses.

Modularity is a very important consideration for such an architecture, i.e., the ability to add, remove, and replace individual units. For example, it should be possible to replace all CPU boards with upgraded models without having to replace memories or disks. Considering that new components will change communication demands, e.g., faster CPU's might require more local bus bandwidth, it is also important that the allocation of boards to local busses can be changed. For example, it should be easy to reconfigure a machine with 4×16 CPU's into one with 8×8 CPU's.

Beyond the effect of faster communication and synchronization, this architecture can also have a significant effect on control overhead, load balancing, and resulting response time problems. Investigations in the Bubba project at MCC demonstrated that large degrees of parallelism may reduce performance unless load imbalance and overhead for startup, synchronization, and communication can be kept low [65]. For example, when placing 100 CPU's either in 100 nodes or in 10 nodes of 10 CPU's each, it is much faster to distribute query plans to all CPU's and much easier to achieve reasonably balanced loads in the second case than in the first case. Within each shared-memory parallel node, load imbalance can be dealt with either by compensating allocation of resources, e.g., memory for sorting or hashing, or by relatively efficient reassignment of data to processors.

Many of today's parallel machines are built as one of the two extreme cases of this hierarchical design: a distributed-memory machine uses single-CPU nodes, while a shared-memory machine consists of a single node. Software designed for this hierarchical architecture will run on either conventional design as well as a genuinely hierarchical machine, and will allow exploring tradeoffs in the range of alternatives in between. The most recent version of Volcano's exchange operator is designed for hierarchical memory, demonstrating that the operator model of parallelization also offers architecture- and topology-independent parallel query evaluation [109]. In other words, the parallelism operator is the only operator that needs to "understand" the underlying architecture, while all data manipulation operators can be implemented without concern for parallelism, data distribution, and flow control.

10. Parallel Algorithms

In the previous section, mechanisms for parallelizing a database query execution engine were discussed. In this section, individual algorithms and their special cases for parallel execution are considered in more detail. Parallel database query processing algorithms are typically based on partitioning an input using range- or hash-partitioning. Either form of partitioning can be combined with sort- and hash-based query processing algorithms; in other words, the choices of partitioning scheme and local algorithm are almost always entirely orthogonal.

When building a parallel system, there is sometimes a question whether it is better to parallelize a slower sequential algorithm with better speedup behavior or a fast sequential algorithm with inferior speedup behavior. The answer to this question depends on the design goal and the planned degree of parallelism. In the few single-

user database systems in use, the goal has been to minimize response time; for this goal, a slow algorithm with linear speedup implemented on highly parallel hardware might be the right choice. In multi-user systems, the goal typically is to minimize resource consumption in order to maximize throughput. For this goal, only the best sequential algorithms should be parallelized. For example, Boral and DeWitt concluded that parallelism is no substitute for effective and efficient indices [38]. For a new parallel algorithm with impressive speedup behavior, the question of whether or not the underlying sequential algorithm is the most efficient choice should always be considered.

10.1. Parallel Selections and Updates

Since disk I/O is a performance bottleneck in many systems, it is natural to parallelize it. Typically, either asynchronous I/O or one process per participating I/O device is used, be it a disk or an array of disks under a single controller. If a selection attribute is also the partitioning attribute, fewer than all disks will contain selection results, and the number of processes and activated disks can be limited. Notice that parallel selection can be combined very effectively with local indices, i.e., indices covering the data of a single disk or node. In general, it is most efficient to maintain indices close to the stored data sets, i.e., on the same node in a parallel database system.

For updates of partitioning attributes in a partitioned data set, items may need to move between disks and sites, just as items may move if a clustering attribute is updated. Thus, updates of partitioning attributes may require setting up data exchange with producers and consumers in order to maintain the consistency of the partitioning. The fact that updating partitioning attributes is more expensive is one reason why immutable (or nearly immutable) identifiers or keys are usually used as partitioning attributes.

10.2. Parallel Sorting

Since sorting is the most expensive operation in many of today's database management systems, much research has been dedicated to parallel sorting [13, 15, 31, 105, 147, 163, 183, 190, 228]. There are two dimensions along which parallel sorting methods can be classified: the number of their parallel inputs (e.g., scan or subplans executed in parallel) and the number of parallel outputs (consumers) [105]. As sequential input or output restrict the throughput of parallel sorts, we assume a multiple-input multiple-output parallel sort here, and further assume that the input items are partitioned randomly with respect to the sort attribute and that the output items should be range-partitioned and sorted.

Considering that data exchange is expensive, both in terms of communication and synchronization delays, each data item should be exchanged only once between processes. Thus, most parallel sort algorithms consist of a local sort and a data exchange step. If the data exchange step is done first, quantiles must be known to ensure load balancing during the local sort step. Such quantiles can be obtained from histograms in the catalogs or by sampling. It is not necessary that the quantiles be precise; a reasonable approximation will suffice.

If the local sort is done first, the final local merging should pass data directly into the data exchange step. On each receiving site, multiple sorted streams must be merged during the data exchange step. One of the possible problems is that all producers of sorted streams first produce low key values, limiting performance by the speed of the first (single!) consumer, then all producers switch to the next consumer, etc.

If a different partitioning strategy than range-partitioning is used, sorting with subsequent partitioning is not guaranteed to be deadlock-free in all situations. Deadlock will occur if (i) multiple consumers feed multiple producers, and (ii) each producer produces a sorted stream and each consumer merges multiple sorted streams, and

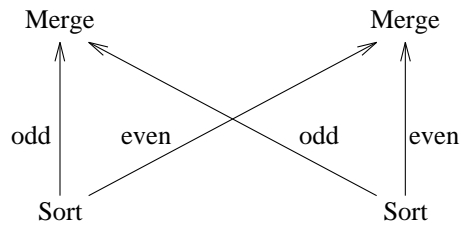


Figure 31. Scenario with Possible Deadlock.

(iii) some key-based partitioning rule is used other than range partitioning, i.e., hash partitioning, and (iv) flow control is enabled, and (v) the data distribution is particularly unfortunate. Figure 31 shows the scenario with two producers and two consumers. Presume that the left producer produces the stream 1, 3, 5, 7, ..., 999, 1002, 1004, 1006, 1008, ..., 2000 while the right producer produces 2, 4, 6, 8, ..., 1000, 1001, 1003, 1005, 1007, ..., 1999. The merging consumers must receive the first item from each producer before they can create their first output item and remove additional items from their input buffers. However, the producers will need to produce 500 items each (and insert them into one consumer's input buffer, all 500 for one consumer) before they will send their first item to the other consumer. The data exchange buffer needs to hold 1000 items at one point of time, 500 on each side of Figure 31, minus some items that are in output and input buffers currently being filled or emptied. If flow control is enabled and the exchange buffer (flow control slack) is less than 500 items, deadlock will occur.

The reason deadlock can occur in this situation is that the producers need to ship data in the order obtained from their input subplan while the consumers need to receive data in sorted order as required by the merge. Thus, there are two sides which both require absolute control over the order in which data pass over the process boundary. If the two requirements are incompatible, an unbounded buffer is required to ensure freedom from deadlock.

In order to avoid deadlock, it must be ensured that one of the conditions outlined earlier is not satisfied. The second condition is the easiest to avoid, and should be focussed on. If the receiving processes do not perform a merge, i.e., the individual input streams are not sorted, deadlock cannot occur because the slack given in the flow control must be somewhere, either at some producer or some consumer or several of them, and the process holding the slack can continue to process data, thus preventing deadlock.

Our recommendation is to avoid the above situation, i.e., to ensure that such query plans are never generated by the optimizer. Consider for which purposes such a query plan would be used. The typical scenario is that

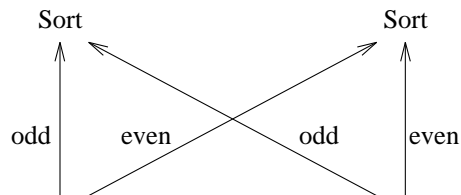


Figure 32. Deadlock-Free Scenario.

multiple processes perform a merge join of two inputs, and each (or at least one) input is sorted by several producer processes. An alternative scenario that avoids the problem is shown in Figure 32. Result data are partitioned and sorted as in the previous scenario. The important difference is that the consumer processes do not merge multiple sorted incoming streams.

An interesting parallel sorting method with balanced communication and without the possibility of deadlock in spite of local sort followed by data exchange (if the data distribution is known a priori) is to sort locally only by the position within the final partition and then exchange data guaranteeing a balanced data flow. This method might be best seen in an example: Consider 10 partitions with key values from 0 to 999 in a uniform distribution. The goal is to have all key values between 0 to 99 sorted on site 0, between 100 and 199 sorted on site 1, etc. First, each partition is sorted locally at its original site, without data exchange, on the the last two digits only, ignoring the first digit. Thus, each site has a sequence like 200, 301, 401, 902, 2, 603, 804, 605, 105, 705, ... 999, 399. Now each site sends data to its correct final destination. Notice that each site sends data simultaneously to all other sites, creating a balanced data flow among all producers and consumers. While this method seems elegant, its problem is that it requires fairly detailed distribution information to ensure the desired balanced data flow.

In shared-memory machines, memory must be divided over all concurrent sort processes. Thus, the more processes are active, the less memory each one can get. The importance of this memory division is the limitation it puts on the size of initial runs and on the fan-in in each merge process. In other words, large degrees of parallelism may impede performance because they increase the number of merge levels. Figure 33 shows how the number of merge levels grows with increasing degrees of parallelism, i.e., decreasing memory per process and merge fan-in. For input size R , total memory size M , and P parallel processes, the merge depth L is $L = \log_{M/P-1}((R/P)/(M/P)) = \log_{M/P-1}(R/M)$.

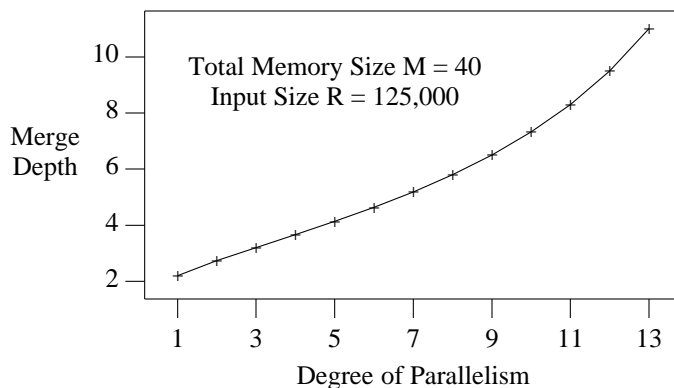


Figure 33. Merge Depth as a Function of Parallelism.

The optimal degree of parallelism must be determined considering the tradeoff between parallel processing and large fan-ins, somewhat similar to the tradeoff between fan-in and cluster size. Extending this argument using the duality of sorting and hashing, too much parallelism in hash-partitioning on shared-memory machines can also be detrimental, both for aggregation and for binary matching.

10.3. Parallel Aggregation and Duplicate Removal

Parallel algorithms for aggregation and duplicate removal are best divided into a local step and a global step. First, duplicates are eliminated locally and then data are partitioned to detect and remove duplicates from different original sites. For aggregation, local and global aggregate functions may differ. For example, to perform a global count, the local aggregation counts while the global aggregation sums local counts into a global count.

For local hash-based aggregation, a special technique might improve performance. Instead of creating overflow files locally on hash table overflow, items are moved directly to their final site. Hopefully, this site can aggregate them immediately into the local hash table because a similar item already exists. In many recent distributed-memory machines, it is faster to ship an item to another site than to do a local disk I/O.¹⁷ The advantage of this technique is that disk I/O is required only when the aggregation output size does not fit into the aggregate memory available on all machines, while the standard local aggregation-exchange-global aggregation scheme requires local disk I/O if any local output size does not fit into a local memory. The difference between the two is determined by the degree to which the original input is already partitioned (usually not at all), making this technique very beneficial.

10.4. Parallel Joins and Other Binary Matching Operations

Binary matching operations such as join, semi-join, outer join, intersection, union, and difference are different than the previous operations exactly because they are binary. For bushy parallelism, i.e., a join for which two subplans create the two inputs independently from one another in parallel, we might consider symmetric hash join algorithms. Instead of differentiating build and probe inputs, the symmetric hash join uses two hash tables, one for each input. When a data item (or packet of items) arrives, the join algorithm first determines which input it came from, and then joins the new data items with the hash table built from the other input as well as inserting the new data items into its hash table such that data items from the other input arriving later can be joined correctly. Such a symmetric hash join algorithm has been used in XPRS, a shared-memory high-performance extensible-relational database system [134, 267, 268]. The advantage of symmetric matching algorithms is that they are independent of the data rates of the inputs; their disadvantage is that they require that both inputs fit in memory.

For parallelizing a single binary matching operation, there are basically two techniques, called here *symmetric partitioning* and *fragment-and-replicate*. In both cases, the global result is the union (concatenation) of all local results. Some algorithms exploit the topology of certain architectures, e.g., ring- or cube-based communication networks [10, 201].

In the symmetric partitioning methods, both inputs are partitioned on the attributes relevant to the operation (i.e., the join attribute for joins or all attributes for set operations), and then the operation is performed at each site.

¹⁷ In fact, some distributed-memory vendors attach disk drives not to the primary processing nodes but to special "I/O nodes" because network delay is negligible compared to I/O time, e.g. in Intel's iPSC/2 and its subsequent parallel architectures.

Both the Gamma and the Teradata database machines use this method. Notice that the partitioning method (usually hashed) and the local join method are independent of each other; Gamma and Grace use hash joins while Teradata uses merge-join.

In the fragment-and-replicate methods, one input is partitioned and the other one is broadcast to all sites. Typically, the larger input is partitioned by not moving it at all, i.e., the existing partitions are processed at their locations prior to the binary matching operation. Fragment-and-replicate methods were considered the join algorithms of choice in early distributed database systems such as R*, SDD-1, and distributed Ingres because communication costs overshadowed local processing costs, and it was cheaper to send a small input to a small number of sites than to partition both a small and a large input.

A technique for reducing network traffic during join processing in distributed database systems uses redundant semi-joins [21, 58, 99], and the idea can also be used in distributed-memory parallel systems. For example, consider the join on a common attribute A of relations R and S stored on two different nodes in a network, say r and s . The semi-join method transfers a duplicate-free projection of R on A to s , performs a semi-join there to determine the items in S that actually participate in the join result, and ships these items to r for the actual join. In other words, based on the relational algebra law that

$$R \text{ JOIN } S = R \text{ JOIN } (S \text{ SEMIJOIN } \pi_A R),$$

cost savings of not shipping all of S were realized at the expense of projecting and shipping the $R.A$ -column and executing the semi-join. Of course, this idea can be used symmetrically to reduce R or S or both, and all operations (projection, duplicate removal, semi-join, and final join) can be executed in parallel on both r and s or on more than two nodes using the parallel join strategies discussed earlier in this subsection. Furthermore, there are probabilistic variants of this idea that use bit vector filtering instead of semi-joins, discussed later in its own subsection.

Roussopoulos and Kang recently showed that symmetric semi-joins are particularly useful [221]. Using the equalities (for a join of relations R and S on attribute A)

$$\begin{aligned} R \text{ JOIN } S &= R \text{ JOIN } \left[S \text{ SEMIJOIN } \pi_A R \right] \\ &= \left[R \text{ SEMIJOIN } \pi_A \left[S \text{ SEMIJOIN } \pi_A R \right] \right] \text{ JOIN } \left[S \text{ SEMIJOIN } \pi_A R \right] \end{aligned} \quad (\text{a})$$

$$= \left[R \overline{\text{SEMIJOIN}} \pi_A \left[S \overline{\text{SEMIJOIN}} \pi_A R \right] \right] \text{ JOIN } \left[S \text{ SEMIJOIN } \pi_A R \right], \quad (\text{b})$$

they designed a four-step procedure to compute the join of two relations stored at two sites. First, the first relation's join attribute column $R.A$ is sent duplicate-free to the other relation's site, s . Second, the first semi-join is computed at s , and either the matching values (term (a) above) or the non-matching values (term (b) above) of the join column $S.A$ are sent back to the first site, r . The choice between (a) and (b) is made based on the number of matching and non-matching¹⁸ values of $S.A$. Third, site r determines which items of R will participate in the join $R \text{ JOIN } S$, i.e., $R \text{ SEMIJOIN } S$. Fourth, both input sites send exactly those items that will participate in the join $R \text{ JOIN } S$ to the site that will compute the final result, which may or may not be one of the two input sites. Of

¹⁸ $\overline{\text{SEMIJOIN}}$ stands for the anti-semi-join, which determines those items in the first input that do *not* have a match in the second input.

course, this two-site algorithm can be used across any number of sites in a parallel query evaluation system.

Typically, each data item is exchanged only once across the interconnection network in a parallel algorithm. However, for parallel systems with small communication overhead, in particular for shared-memory systems, and in parallel processing systems with processors without local disk, it may be useful to spread each overflow file over all available nodes and disks in the system. The disadvantage of the scheme may be communication overhead; however, the advantages of load balancing and cumulative bandwidth while reading a partition file have led to the use of this scheme both in the Gamma and SDC database machines, called *bucket spreading* in the SDC design [74, 165].

For parallel non-equi-joins, a symmetric fragment-and-replicate method has been proposed by Stamos and Young [261]. As shown in Figure 34, processors are organized into rows and columns. One input relation is partitioned over rows and partitions are replicated within each row, while the other input is partitioned and replicated over columns. Each item from one input "meets" each item from the other input at exactly one site, and the global join result is the concatenation of all local joins.

Avoiding partitioning as well as broadcasting for many joins can be accomplished with a physical database design that considers frequently performed joins and distributes and replicates data over the nodes of a parallel or distributed system such that many joins already have their input data suitably partitioned. Katz and Wong formalized this notion as *local sufficiency* [153, 292]; more recent research on the issue was performed in the Bubba project [65].

For joins in distributed systems, a third class of algorithms, called *fetch-as-needed*, was explored. The idea of these algorithms is that one site performs the join by explicitly requesting (fetching) only those items from the other input needed to perform the join [67, 287]. If one input is very small, fetching only the necessary items of the larger input might seem advantageous. However, this algorithm is a particularly poor implementation of a semi-join technique discussed above. Instead of requesting items or values one by one, it seems better to first project all join attribute values, ship (stream) them across the network, perform the semi-join using any local binary matching algorithm, and then stream exactly those items that will be required for the join back to the first site. The difference between the semi-join technique and fetch-as-needed is that the semi-join scans the first input twice, once to extract the join values and once to perform the real join, while fetch-as-needed needs to work on each data

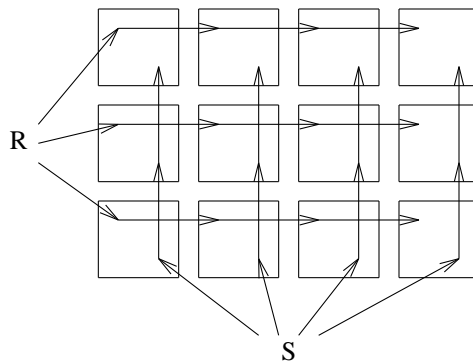


Figure 34. Symmetric Fragment-and-Replicate Join.

item only once.

10.5. Parallel Universal Quantification

In our previous discussion on sequential universal quantification, we discussed four algorithms for universal quantification or relational division, namely naive division (a direct, sort-based algorithm), hash-division (direct, hash-based), and sort- and hash-based aggregation (indirect) algorithms, which might require semi-joins and duplicate removal in the inputs.

For naive division, pipelining can be used between the two sort operators and the division operator. However, both quotient partitioning and divisor partitioning can be employed as described below for hash-division.

For algorithms based on aggregation, both pipelining and partitioning can be applied immediately using standard techniques for parallel query execution. While partitioning seems to be a promising approach, it has an inherent problem due to the possible need for a semi-join. Recall that in the example for universal quantification using transcript and course relations, the join attribute in the semi-join (*course-no*) is different than the grouping attribute in the subsequent aggregation (*student-id*). Thus, the *Transcript* relation has to be partitioned twice, once for the semi-join and once for the aggregation.

For hash-division, pipelining has only limited promise because the entire division is performed within a single operator. However, both partitioning strategies discussed earlier for hash table overflow can be employed for parallel execution, i.e., quotient partitioning and divisor partitioning [102, 113].

For hash-division with quotient partitioning, the divisor table must be *replicated* in the main memory of all participating processors. After replication, all local hash-division operators work completely independently of each other. Clearly, replication is trivial for shared-memory machines, in particular since a single copy of the divisor table can be shared without synchronization among multiple processes once it is complete.

When using divisor partitioning, the resulting partitions are processed in parallel instead of in phases as discussed for hash table overflow. However, instead of tagging the quotient items with phase numbers, processor network addresses are attached to the data items, and the collection site divides the set of all incoming data items over the set of processor network addresses. In the case that the central collection site is a bottleneck, the collection step can be decentralized using quotient partitioning.

11. Non-Standard Query Processing Algorithms

In this section, we briefly review the query processing needs of data models and database systems for non-standard applications. In many cases, the logical operators defined for new data models can use existing algorithms, e.g., for intersection. The reason is that for processing, bulk data types such as array, set, bag (multi-set), or list are represented as sequences similar to the streams used in the query processing techniques discussed earlier, and the algorithms to manipulate these bulk types are equal to the ones used for sets of tuples, i.e., relations. However, some algorithms are genuinely different from the algorithms we have surveyed so far. In this section, we review operators for nested relations, temporal and scientific databases, object-oriented databases, and more meta-operators for additional query processing control.

There are several reasons for integrating these operators into an algebraic query processing system. First, it permits efficient data transfer from the database to the application embodied in these operators. The interface between database operators is designed to be as efficient as possible; the same efficient interface should also be used for applications. Second, operator implementors can take advantage of the control provided by the meta-

operators. For example, an operator for a scientific application can be implemented in a single-process environment and later parallelized with the exchange operator. Third, query optimization based on algebraic transformation rules can cover all operators, including operations that are normally considered database application code. For example, using algebraic optimization tools such as the EXODUS and Volcano optimizer generators [100, 108], optimization rules that can move an unusual database operator in a query plan are easy to implement. For example, for a sampling operator, a rule might permit transforming an algebra expression to query a sample instead of sampling a query result.

11.1. Nested Relations

Nested relations, or Non-First-Normal-Form (NF²) relations, permit relation-valued attributes in addition to atomic values such as integers and strings used in the normal or "flat" relational model. For example, in an order processing application, the set of individual line items on each order could be represented as a nested relation, i.e., as part of an order tuple. Figure 35 shows a NF² relation with two tuples with two and three nested tuples and the equivalent normalized relations, which we call the master and detail relations. Nested relations can be used for all one-to-many relationships but are particularly well-suited for the representation of "weak entities" in the Entity-Relationship (ER) Model [51], i.e., entities whose existence and identification depends on another entity as for order entries in Figure 35. In general, nested sub-tuples may include relation-valued attributes, with arbitrary nesting depth. The advantages of the NF² model are that component relationships can be represented more naturally than in the fully normalized model, many frequent join operations can be avoided, and structural information can be used for physical clustering. Its disadvantage is the added complexity, in particular in storage management and query processing.

Several algebras for nested relations have been defined, e.g. [78, 231]. Our discussion here focuses not on the conceptual design of NF² algebras but on algorithms to manipulate nested relations, which are unfortunately not very well documented to-date.

Order -No	Customer -No	Date	Items	
			Part-No	Count
110	911	910902	4711	8
			2345	7
112	912	910902	9876	3
			2222	1
			2357	9

Order-No	Customer-No	Date
110	911	910902
112	912	910902

Order-No	Part-No	Quantity
110	4711	8
110	2345	7
112	9876	3
112	2222	1
112	2357	9

Figure 35. Nested Relation and Equivalent Flat Relations.

Two operations required in NF^2 database systems are operations that transform a NF^2 relation into a normalized relation with atomic attributes only, and vice versa. The first operation is frequently called *unnest* or *flatten*; the opposite direction is called the *nest* operation. The unnest operation can be performed in a single scan over the NF^2 relation that includes the nested subtuples; both normalized relations in Figure 35 and their join can be derived readily enough from the NF^2 relation. The nest operation requires grouping of tuples in the detail relation and a join with the master relation. Grouping and join can be implemented using any of the algorithms for aggregate functions and binary matching discussed earlier, i.e., sort- and hash-based sequential and parallel methods.

All operations defined for flat relations can also be defined for nested relations, in particular selection, join, and set operations (union, intersection, difference). For selections, additional power is gained with selection conditions on subtuples and sets of subtuples using set comparisons or existential or universal quantification. In principle, since the nested relation is a relation, any relational calculus and algebra expression should be permitted for it. In the example of Figure 35, there may be a selection of orders in which the ordered quantity of all items is more than 100, which is a universal quantification. The algorithms for selections with quantifier are similar to the ones discussed earlier for flat relations, e.g., relational semi-join and division, but are easier to implement because the grouping process built into the flat-relational algorithms is inherent in the nested tuple structure.

For joins, similar considerations apply. Matching algorithms discussed earlier can be used in principle. They may be more complex if the join predicate involves subrelations, and algorithm combinations may be required that are derived from a flat-relation query over flat relations equivalent to the NF^2 query over the nested relations. However, there should be some performance improvements possible if the grouping of values in the nested relations can be exploited, as for example in the join algorithms described by Rosenthal et al. & [217].

11.2. Temporal and Scientific Database Management

For a variety of reasons, management and manipulation of statistical, temporal, and scientific data are gaining interest in the database research community. Most work on temporal databases has focused on its semantics and representation in data models and query languages [188, 256]; some work has considered special storage structures, e.g. [1, 180, 218, 242], algebraic operators, e.g., temporal joins [238], and optimization of temporal queries, e.g. [120, 237]. While logical query algebras require extensions to accommodate time, only some storage structures and algorithms, e.g., multi-dimensional indices, differential files, and versioning, and the need for approximate selection and matching (join) predicates are new in the query execution algorithms for temporal databases.

A number of operators can be identified that both add functionality to database systems used to process scientific data and fit into the database query processing paradigm. Schneider et al. considered algorithms for join predicates that express proximity, i.e., join predicates of the form $R.A - c_1 \leq S.B \leq R.A + c_2$ [75]. Such join predicates operations are very different from the usual use of relational join. They do not reestablish relationships based on identifying keys but match data values that express a dimension in which distance can be defined, in particular time. Traditionally, such join predicates have been considered non-equi-joins and were evaluated by a variant of nested loops join. However, such "band joins" can be executed much more efficiently by a variant of merge-join that keeps a "window" of inner-relation tuples in memory or by a variant of hash join that uses range-partitioning and assigns some build tuples to multiple partition files. A similar partitioning model must be used for parallel execution, requiring multi-cast for some tuples. Clearly, these variants of merge-join and hash join will outperform nested loops for large inputs unless the band is so wide that the join result approaches the Cartesian product.

For storage and management of the massive amounts of data resulting from scientific experiments, database techniques are very desirable. Operators for processing time series in scientific databases are based on an interpretation of a stream between operators not as a set of items (as in most database applications) but as a sequence in which the order of items in a stream has semantic meaning. For example, data reduction using interpolation as well as extrapolation can be performed within the stream paradigm. Similarly, digital filtering [130] also fits the stream processing protocol very easily. Interpolation, extrapolation, and digital filtering with a single algorithm (physical operator) were implemented in the Volcano system to verify this fit, including their optimization and parallelization [290]. Another promising candidate is visualization of single-dimensional arrays such as time series.

Problems that do not fit the stream paradigm, e.g., many matrix operations such as transformations used in linear algebra, Laplace or Fast Fourier Transform, and slab (multi-dimensional sub-array) extraction, are not as easy to integrate into database query processing systems. Some of them seem to fit better into the storage management sub-system rather than the algebraic query execution engine. For example, slab extraction has been integrated into the NetCDF storage and access software [47, 212].

11.3. Object-Oriented Database Systems

Research into query processing for extensible and object-oriented systems has been growing rapidly in the last few years. Most proposals or implementations use algebras for query processing, e.g. [2, 63, 101, 155, 192, 245-247, 273, 275, 283, 295]. These algebras resemble relational algebra in the sense that they focus on bulk data types but are generalized to support operations on arrays, lists, etc., user-defined operations (methods) on instances, heterogeneous bulk types, and inheritance. The use of algebras permits several important conclusions. First, naive execution models that execute programs as if all data were in memory are not the only alternative. Second, data manipulation operators can be designed and implemented that go beyond data retrieval and permit some amount of data reduction, aggregation, and even inference. Third, algebraic execution techniques including the stream paradigm and parallel execution can be used in object-oriented data models and database systems. Fourth, algebraic optimization techniques will continue to be useful.

Associative operations are an important part in all object-oriented algebras because they permit reducing large amounts of data to the interesting subset of the database suitable for further consideration and processing. Thus, set processing and matching algorithms as discussed earlier in this survey will be found in object-oriented systems, implemented in such a way that they can operate on heterogeneous sets. The challenge for query optimization is to map a complex query involving complex behavior and complex object structures to primitives available in a query execution engine. Translating an initial request with abstract data types and encapsulated behavior coded in a computationally complete language into an internal form that both captures the entire query's semantics and allows effective query optimization is still an open research issue [68, 101].

Beyond associative indices discussed earlier, object-oriented systems can also benefit from special relationship indices, i.e., indices that contain condensed information about inter-object references. In principle, these index structures are similar to join indices [282] but can be generalized to support multiple levels of referencing. Examples for indices in object-oriented database systems include the work of Maier and Stein in the Gemstone object-oriented database system product [186], Bertino et al. in the Orion project [24-26] and by Kemper et al. in the GOM project [157-159]. At this point, it is too early to decide which index structures will be the most useful because the entire field of query processing in object-oriented systems is still developing rapidly, from query

languages to algebra design, algorithm repertoire, and optimization techniques. Other areas of intense current research interest are buffer management and clustering of objects on disk.

One of the big performance penalties in object-oriented database systems is the expense of "pointer chasing" (using OID references) which may involve object faults, also called "goto's on disk." In order to reduce I/O cost, some systems use what amounts to main memory databases or map the entire database into virtual memory. For systems with an explicit database on disk and an in-memory buffer, there are various techniques to detect object faults; some commercial object-oriented database systems use hardware mechanisms originally perceived and implemented for virtual-memory systems. While such hardware support makes fault detection faster, it does not address the problem of expensive I/O operations. In order to reduce actual I/O cost, read-ahead and planned buffering must be used. Palmer and Zdonik recently proposed keeping access patterns or sequences and activating read-ahead if accesses equal or similar to a stored pattern are detected [208]. Another recent proposal for efficient assembly of complex objects uses a window (a small set) of open references and resolves, at any point of time, the most convenient one by fetching this object or component from disk, which has shown dramatic improvements in disk seek times and makes complex object retrieval more efficient and more independent of object clustering [154].

11.4. More Meta-Operators

The exchange operator used for parallel query processing is not a normal operator in the sense that it does not manipulate, select, or transform data. Instead, the exchange operator provides control of query processing in a way orthogonal to what a query does and what algorithms it uses. Therefore, we call it a meta-operator. There are several other meta-operators that can be used in database query processing, and we survey them briefly in this section.

In query processing systems, dataflow is usually paced or driven from the top, the consumer. The left-most diagram of Figure 36 shows the control flow of normal iterators. The data flow of all diagrams in Figure 36 is assumed to be upward. However, in real-time systems that capture data from experiments, this approach may not be realistic because the data source, e.g., a satellite receiver, has to be able to unload data as they arrive. In such systems, data-driven operators, shown in the second diagram of Figure 36, might be more appropriate. To combine the algorithms implemented and used for query processing with such real-time data capture requirements, one could design data *flow translation* meta-operators. The first such operator which we call the *active scheduler* can be used between a demand-driven producer and a data-driven consumer. In this case, neither operator will schedule the other; therefore, an *active scheduler* that demands items from the producer and forces them onto the consumer will glue these two operators together. An active scheduler schematic is shown in the third diagram of

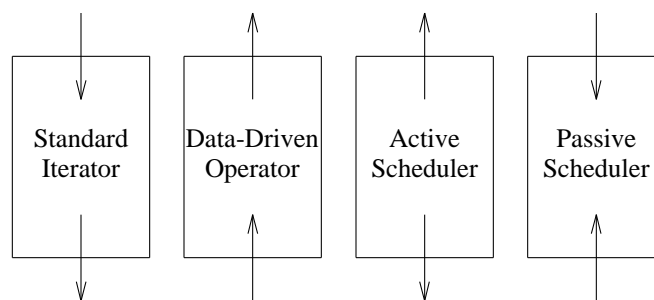


Figure 36. Operators, Schedulers, and Control Flow.

Figure 36. The opposite case, a data-driven producer and a demand-driven consumer, has two operators, each trying to schedule the other one. A second flow control operator, called the *passive scheduler*, can be built that accepts procedure calls from either neighbor and resumes the other neighbor in a co-routine fashion to ensure that the resumed neighbor will eventually demand the item the scheduler just received. The final diagram of Figure 36 shows the control flow of a passive scheduler. (Notice that this case is similar to the bracket model of parallel operator implementations discussed earlier in which an operating system or networking software layer had to be placed between data manipulation operators and perform buffering and flow control.)

For very complex queries, it might be useful to break the data flow between operators at some point, for two reasons. First, if too many operators run in parallel, contention for memory or temporary disks might be too intense, and none of the operators will run as efficiently as possible. A long series of hybrid hash joins in a right-deep query plan illustrates this situation. Second, due to the inherent error in selectivity estimation during query optimization [146, 187], it might be worthwhile to execute only a subset of a plan, verify the correctness of the estimation, and then resume query processing with another few steps. After a few processing steps have been performed, their result size and other statistical properties such as minimum and maximum and approximate number of duplicate values can be easily determined while saving the result on temporary disk.

In principle, this was done in Ingres' original optimization method called *Decomposition*, except that Ingres performed only one operation at a time before optimizing the remaining query [291, 294]. We propose alternating more slowly between optimization and execution, i.e., to perform a "reasonable" number of steps between optimizations, where reasonable may be three to ten selections and joins depending on errors and error propagation in selectivity estimation. Stopping the data flow and resuming after additional optimization could very well turn out to be the most reliable technique for very large, complex queries.

Implementation of this technique could be embodied in a new meta-operator, which we call the *stop-and-go* operator, but further research is required to develop the techniques used for placing stop-and-go operators in complex plans. An initial implementation of a stop-and-go operator that uses preoptimized alternative subplans is the *choose-plan* operator implemented in Volcano and first described in 1989 [103]. In its current implementation, it executes zero or more subplans and then invokes a decision function provided by the optimizer that decides which of multiple equivalent plans to execute depending on intermediate result statistics, current system load, and runtime values of query parameters unknown at optimization time. We plan on studying optimization and creation of *dynamic query evaluation plans* for very complex queries using the Volcano optimizer generator.

The final function that can be encapsulated in a meta-operator is an operator that passes the result of a common subexpression to multiple consumers, as mentioned briefly in the section of the architecture of query execution engines. The problem is that multiple consumers, typically demand-driven and demand-driving their inputs, will request items of the common subexpression result at different times or rates. The two standard solutions are either to execute the common subexpression into a temporary file and let each consumer scan this file at will, or to determine which consumer will require be the first to require the result of the common subexpression, to execute the common subexpression as part of this consumer, and to create a file with the common subexpression result as a by-product of the first consumer's execution. Instead, we suggest a new meta-operator, which we call the *split* operator, to be placed at the top of the common subexpression's plan and which can serve multiple consumers at their own paces. It automatically performs buffering to account for different paces, uses temporary disk space if the discrepancies are too wide, and is suitably parameterized to permit both standard solutions described above.

12. Additional Techniques for Performance Improvement

In this section, we consider some additional techniques that have been proposed in the literature or used in real systems, and that have not been discussed in earlier sections of this survey. In particular, we consider precomputation, data compression, surrogate processing, bit vectors, and specialized hardware. Recently proposed techniques that have not been fully developed are not discussed here, e.g., "racing" equivalent plans and terminating the ones that seem not competitive after some small amount of time.

12.1. Precomputation and Derived Data

It is trivial to answer a query for which the answer is already known — therefore, precomputation of frequently requested information is an obvious idea. The problem with keeping preprocessed information in addition to base data is that it is redundant and must be invalidated or maintained upon updates to the base data.

Precomputation and of derived data such as relational views are duals. Thus, concepts and algorithms designed for one will typically work well for the other. The main difference is the database user's view: precomputed data are typically used after a query optimizer has determined that they can be used to answer a user query against the base data, while derived data are known to the user and can be queried without regard to the fact that they actually must be derived at run-time from stored based data. Not surprisingly, since derived data are likely to be referenced and requested by users and application programs, precomputation of derived data has been investigated, both for relational and object-oriented data models.

Indices are the simplest form of precomputed data since they are a redundant and, in a sense, precomputed selection. They represent a compromise between a non-redundant database and one with complex precomputed data because they can be maintained relatively efficiently.

The next more sophisticated form of precomputation are inversions as provided in System R's "0th" prototype [48], view indices as analyzed by Roussopoulos [220], two-relation *join indices* as proposed by Valduriez [282], or domain indices as used in the ANDA project (called *VALTREE* there) [77] in which all occurrences of one domain (e.g., part number) are indexed together and each index entry contains a relation identification with each record identifier. With join or domain indices, join queries can be answered very fast, typically faster than using multiple single-relation indices. On the other hand, single-selection selections and updates may be slightly slower if there are more entries for each indexed key.

For binary operators, there is a spectrum of possible levels of precomputations¹⁹, explored predominantly for joins. The simplest form of precomputation in support of binary operations are individual indices, e.g., clustering B-trees that ensure and maintain sorted relations. On the other extreme are completely materialized join results. Intermediate levels are pointer-based joins [248] (discussed earlier in the section on matching) and join indices [282]. For each form of precomputed result, the required redundant data structures must be maintained each time the underlying base data are updated, and larger retrieval speedup might be paid for with larger maintenance overhead.

Babb explored storing only results of outer join, but not the normalized base relations, in the content-addressable file store (CAFS), and called this encoding join normal form [7]. Larson et al. investigated storing and maintaining materialized views in relational database systems [33, 34, 177, 189, 277, 293]. Their hope was to

¹⁹ This paragraph was written using ideas and notes by José Blakeley.

speed relational query processing by using derived data, possibly without storing all base data, and ensuring that their maintenance overhead would be less than their benefits in faster query processing. For example, Blakeley demonstrated that for a single join there exists a large range of retrieval and update mixes in which materialized views outperform both join indices and hybrid hash join [34]. This investigation should be extended, however, for more complex queries, e.g., three- and four-way joins, and for queries in object-oriented systems and emerging database applications.

Hanson compared query modification (i.e., query evaluation from base relations) against the maintenance costs of materialized views, and considered in particular the cost of immediate vs. deferred updates [131]. His results indicate that for modest update rates, materialized views provide better system performance. Furthermore, for modest selectivities of the view predicate, deferred view maintenance using differential files [242] outperforms immediate maintenance of materialized views. However, Hanson also did not include multi-way joins in his study.

Sellis analyzed caching of results in a query language called Quel+ (which is a subset of Postquel [270]) over a relational database with procedural (QUEL) fields [240]. He also considered the case of limited space on secondary storage used for caching query results, and replacement algorithms for query results in the cache when the space becomes insufficient.

Links between records (pointers of some sort, e.g., record, tuple, or object identifiers) are another form of precomputation. Links are particularly effective for system performance if they are combined with clustering (assignment of records to pages). Database systems for the hierarchical and network models have used physical links and clustering, but supported basically only queries and operations that were "precomputed" in this way. Some researchers tried to overcome this restriction by building relational query engines on top of network systems, e.g. [52, 216, 296]. However, with performance improvements in the relational world, these efforts seem to have been abandoned. With the advent of extensible and object-oriented database management systems, combining links and ad-hoc query processing might become a more interesting topic again. A recent effort for an extensible-relational system are Starburst's pointer-based joins discussed earlier [126, 248].

In order to ensure good performance for its extensive rule processing facilities, Postgres uses precomputation and caching of the action parts of production rules [266, 270, 271]. For automatic maintenance of such derived data, persistent "invalidation locks" are stored for detection of invalid data after updates to the base data. A performance evaluation of this caching and invalidation scheme is still outstanding.

Finally, the Cactis project focused on maintenance of derived data in object-oriented environments [142]. The conclusions of this project include that incremental maintenance coupled with a fairly simple adaptive clustering algorithm is an efficient way to propagate updates to derived data.

12.2. Data Compression

A number of researchers have investigated the effect of compression on database systems and their performance [106, 184, 222, 243]. There are two types of compression in database systems. First, the amount of redundancy can be reduced by prefix- and suffix-truncation, in particular in indices, and by use of encoding tables (e.g., color combination "9" means "red car with black interior"). Second, compression schemes can be applied to attribute values, e.g., adaptive Huffman coding or Ziv-Lempel methods [19, 178]. This type of compression can be exploited most effectively in database query processing if all attributes of the same domain use the same encoding, e.g., the "Part-No" attributes of datasets representing parts, orders, shipments, etc. because common encodings permit comparisons without decompression.

Most obviously, compression can reduce the amount of disk space required for a given data set. Disk space savings has a number of ramifications on I/O performance. First, the reduced data space fits into a smaller physical disk area; therefore, the seek distances and seek times are reduced. Second, more data fit into each disk page, track, and cylinder, allowing more intelligent clustering of related objects into physically near locations. Third, the unused disk space can be used for disk shadowing to increase reliability, availability, and I/O performance [30]. Fourth, compressed data can be transferred faster to and from disk. In other words, data compression is an effective means to increase disk bandwidth (not by increasing physical transfer rates but by increasing the information density of transferred data) and to relieve the I/O bottleneck found in many high-performance database management systems [38]. Fifth, in distributed database systems and in client-server situations, compressed data can be transferred faster across the network than uncompressed data. Uncompressed data require either more network time or a separate compression step. Finally, retaining data in compressed form in the I/O buffer allows more records to remain in the buffer, thus increasing the buffer hit rate and reducing the number of I/O's. The last three points are actually more general. They apply to the entire storage hierarchy of tape, disk, controller caches, local and remote main memories, and CPU caches.

For query processing, compression can be exploited far beyond improved I/O performance because decompression can often be delayed until a relatively small data set is presented to the user or an application program. First, exact-match comparisons can be performed on compressed data. Second, projection and duplicate removal can be performed without decompressing data. The situation for aggregation is a little more complex since the attribute on which arithmetic is performed typically must be decompressed. Third, neither the join attributes nor other attributes need to be decompressed for most joins. Since keys and foreign keys are from the same domain, and if compression schemes are fixed for each domain, a join on compressed key values will give the same results as a join on normal, decompressed key values. It might seem unusual to perform a merge-join in the order of compressed values, but it nonetheless is possible and will produce correct results.

There are a number of benefits from processing compressed data. First, materializing output records is faster because records are shorter and less copying is required. Second, for inputs larger than memory, more records fit into memory. In hybrid hash join, for instance, the fraction of the file that can be retained in the hash table and thus be joined without any I/O is larger. During sorting, the number of records in memory and thus the number of records per run is larger, leading to fewer runs and possibly fewer merge levels. Third, and very interestingly, skew is less likely to be a problem. The goal of compression is to represent the information with as few bits as possible. Therefore, each bit in the output of a good compression scheme has close to maximal information content, and bit columns seen over the entire file are unlikely to be skewed. Furthermore, bit columns will not be correlated. Thus, the compressed key values can be used to create a hash value distribution that is almost guaranteed to be uniform, i.e., optimal for hashing in memory and partitioning to overflow files as well as to multiple processors in parallel join algorithms.

12.3. Surrogate Processing

Another very useful technique in query processing is the use of surrogates for intermediate results. A surrogate is a reference to a data item, be it a logical object identifier (OID) used in object-oriented systems or a physical record identifier (RID) or location. Instead of keeping a complete record in memory, only the fields that are used immediately are kept and the remainder replaced by a surrogate, which has in principle the same effect as compression.

The simplest case in which surrogate processing can be exploited is in avoiding copying. Consider a relational join; when two items satisfy the join predicate, a new tuple is created from the two original ones. Instead of copying the data fields, it is possible to create only a pair of RID's or pointers to the original records if they are kept in memory. If a record is 50 times larger than a RID, e.g., 8 B vs. 400 B, the effort spent on copying bytes is reduced by that factor.

Copying is already a major part of the CPU time spent in many query processing systems, but it is becoming more expensive for two reasons. First, many modern CPU designs and implementations are optimized for an impressive number of instructions per second but do not provide the performance improvements in mundane tasks such as moving bytes from one memory location to another [205]. Second, many modern computer architectures employ multiple CPU's accessing shared memory over one bus because this design permits fast and inexpensive parallelism. Although alleviated by local caches, bus contention is the major bottleneck and limitation to scalability in shared-memory parallel machines. Therefore, reductions in memory-to-memory copying in database query execution engines permits higher useful degrees of parallelism in shared-memory machines.

A second example for surrogate processing was mentioned earlier in connection with indices. To evaluate a conjunction with multiple clauses, each of which is supported by an index, it might be useful to perform an intersection of RID-lists to reduce the number of records needed before actual data are accessed.

A third case is the use of indices and RID's to evaluate joins, for example in the query processing techniques used in Ingres [169, 170] and IBM's hybrid join [55] discussed earlier in the section on binary matching.

Surrogate processing has also been used in parallel systems, in particular distributed-memory implementations, to reduce network traffic. For example, Lorie and Young used RID's to reduce the communication time in parallel sorting by sending (sort key, RID) pairs to a central site, which determines each record's global rank, and then re-partitioning and merging records very quickly by their rank alone without further data comparisons [183].

Berra et al. considered indexing and retrieval organizations for very large (relational) knowledge bases and databases [23, 62]. They employed three techniques, concatenated code words (CCW's), superimposed code words (SCW's), and transformed inverted lists (TIL's). TIL's are normal index structures for all attributes of a relation that permit answering conjunctive queries by bitwise *anding*. CCW's and SCW's use hash values of all attributes of a tuple and either concatenate such hash values or bitwise *or* them together. The resulting code words are then used as keys in indices. In their particular architecture, Berra et al. consider associative memory and optical computing techniques to search efficiently through such indices, although conventional software techniques could be used as well.

Techniques based on hash values and bit patterns have also been used in other contexts; we discuss bit vector filtering in the next subsection.

12.4. Bit Vector Filtering

In parallel systems, bit vectors have been used for what we call here "probabilistic semi-joins." Consider a relational join to be executed on a distributed-memory machine with repartitioning of both input relations on the join attribute. It is clear that communication effort could be reduced if only the tuples that actually contribute to the join result, i.e., those with a match in the other relation, needed to be shipped across the network. To accomplish this, distributed database systems were designed to make extensive use of semi-joins, e.g., SDD-1 [21].

A faster alternative to semi-joins, which, as discussed earlier, require basically the same computational effort as natural joins, is the use of bit vectors [6]. A bit vector with N bits is initialized with zeroes, and all items in the first (preferably the smaller) input are hashed on their join key to $0, \dots, N - 1$. For each item, one bit in the bit vector filter is set to one; hashing collisions are ignored. After the first join input has been exhausted, the bit vector is used to filter the second input. Data items of the second input are hashed on their join key value, and only items for which the bit is set to one can possibly participate in the join. There is some chance for false passes in the case of collisions, i.e., items of the second input pass the bit vector filter although they actually do not participate in the join, but if the bit vector is sufficiently large, the number of false passes is very small. For one-to-one match operations other than join, e.g., outer join and union, bit vectors can also be used but the algorithm must be modified to ensure that items that do not pass the bit vector are properly included in the operation's output stream.

In general, if the number of bits is about twice the number of items in the first input, bit vectors are very effective. If many more bits are available, the bit vector can be split into multiple subvectors or multiple bits can be set for each item using multiple hash functions, reducing the number of false passes. Babb analyzed the use of multiple bit vectors in detail [6].

The Gamma relational database machine demonstrated the effectiveness of bit vector filtering in relational join processing on distributed-memory hardware [72-74, 92]. While Gamma used bit vector filtering basically only for joins, it is equally applicable to all other one-to-one match operators, including semi-join, outer join, intersection, union, and difference. For operators that include non-matching items in their output, e.g., outer joins and unions, part of the result can be obtained before network transfer, based solely on the bit vector. For parallel relational division (universal quantification), bit vector filtering can be used on the divisor attributes to eliminate most of the dividend items that do not pertain to any divisor item. Thus, our earlier assessment that universal quantification can be performed as fast as existential quantification (a semi-join of dividend and divisor relations) even extends to special techniques used to boost join performance.

Bit vector filtering can also be exploited in sequential systems. Consider a merge-join with sort operations on both inputs. If the bit vector is built based on the input of the first sort, i.e., the bit vector is completed when all data have reached the first sort operator. This bit vector can then be used to reduce the input into the second sort operator on the (presumably larger) second input. Depending on how the sort operation is organized into phases, it might even be possible to create a second bit vector from the second merge-join input and use it to reduce the first join input while it is being merged.

For sequential hash joins, bit vectors can be used in two ways. First, they can be used to filter items of the probe input using a bit vector created from items of the build input. This use of bit vectors is analogous to bit vector usage in parallel systems and for merge-join. Second, bit vectors can be used for each partition in each recursion level. In the Volcano query processing system, the operator implementing hash join, intersection, etc. uses the space used as anchor for each bucket's linked list for a small bit vector filter after the bucket has been spilled to an overflow file. Only those items from the probe input that pass the bit vector filter are written to the probe overflow file. This technique is used in each recursion level of overflow resolution. Thus, during recursive partitioning, relatively small and efficient bit vector filters can be used repeatedly and at increasingly finer granularity to remove items from the probe input that do not contribute to the join result. Bit vectors could also be used to remove items from the build input using bit vectors created from the probe input; however, since the probe input is presumed the larger input and hash collisions in the bit vector would make the filter less effective, it may or may not be an effective technique.

In order to find and exploit a dual in the realm of sorting and merge-join to bit vector filtering in each recursion level of recursive hash join, sorting of multiple inputs must be divided into individual merge levels. In other words, for a merge-join of inputs R and S , the sort activity should switch back and forth between R and S , level by level, creating and using a new bit vector filter in each merge level. Unfortunately, even with a sophisticated sort implementation that supports this use of bit vector filters in each merge level, recursive hybrid hash join will make more effective use of bit vector filters because the inputs are partitioned, thus reducing the number of distinct values in each partition in each recursion level.

12.5. Specialized Hardware

Specialized hardware was considered by a number of researchers, e.g., hardware sorters and logic-per-track selection. A relatively recent survey of database machine research is given by Su [274]. Most of this research was abandoned after Boral and DeWitt's influential analysis [38] that compared CPU and I/O speeds and their trends. They concluded that I/O is most likely the bottleneck in future high-performance query execution, not processing. Therefore, they recommended moving from research on custom processors to techniques for overcoming the I/O bottleneck, e.g., by use of parallel readout disks, disk caching and read-ahead, and indexing to reduce the amount of data to be read for a query. Other investigations also came to the conclusion that parallelism is no substitute for effective storage structures and query execution algorithms [69, 195]. An additional very strong argument against custom VLSI processors is that microprocessor speed is currently improving so rapidly that it is likely that, by the time a small organization has designed, fabricated, and tested a special hardware component and integrated it into a larger hardware and software system, the next generation of general-purpose CPU's is available and can execute database functions programmed in a high-level language at the same speed as the specialized hardware component. Furthermore, it is not clear what specialized hardware would be most beneficial to design, in particular in light of today's directions towards extensible database systems and emerging database application domains. Therefore, we do not favor specialized database hardware modules beyond general-purpose processing, storage, and communication hardware dedicated to executing database management software.

Summary and Outlook

Database management systems provide three essential groups of services. First, they maintain both data and associated meta-data in order to make databases self-contained and self-explanatory, at least to some extent, and to provide data independence. Second, they support facilities for data sharing among multiple users as well as prevention and recovery of failures and data loss. Third, they raise the level of abstraction for data manipulation above the primitive access commands provided by file systems with more or less sophisticated matching and inference mechanisms, commonly called the query language or query processing facility. We have surveyed execution algorithms and software architectures used in providing this third essential service.

Query processing has been explored extensively in the last 20 years in the context of relational database management systems, and is slowly gaining interest in the research community for extensible and object-oriented systems. This is a very encouraging development, because if these new systems have increased modeling power over previous data models and database management systems but cannot execute even simple requests efficiently, they will never gain widespread use and acceptance. Databases will continue to manage massive amounts of data; therefore, efficient query and request execution will continue to represent both an important research direction and an important criterion in investment decisions in the "real world." In other words, new database management systems should provide greater modeling power (this is widely accepted and intensely pursued) but also competitive

or better performance than previous systems. We hope that this survey will contribute to the use of efficient and parallel algorithms for query processing tasks in new database management systems.

A large set of query processing algorithms has been developed for relational systems. Sort- and hash-based techniques have been used for physical storage design, for associative index structures, for algorithms for unary and binary matching operations such as aggregation, duplicate removal, join, intersection, and division, and for parallel query processing using hash-partitioning or range-partitioning. Additional techniques such as precomputation and compression have been shown to provide substantial performance benefits when manipulating large volumes of data. Many of the existing algorithms will continue to be useful for extensible and object-oriented systems, and many can easily be generalized from sets of tuples to more general pattern matching functions. Some emerging database applications will require new operators, however, both for translation between alternative data representations and for actual data manipulation.

The most promising aspect of current research into database query processing for new application domains is that the concept of multiple operators, each performing a part of the required data manipulation and each passing an intermediate result to the next operator, is versatile enough to meet the new challenges. This concept permits specification of database queries and requests in a logical algebra as well as concise representation of database programs in a physical algebra. Furthermore, it allows algebraic optimizations of requests, i.e., optimizing translations of logical into physical expressions. Finally, it permits pipelining between operators to exploit parallel computer architectures, and partitioning of stored data and intermediate results for most operators, in particular for operators on sets.

We can hope that much of the existing relational technology for query optimization and parallel execution will remain relevant, and that research into extensible optimization and parallelization will make a significant impact on future database applications. In fact, for database management systems to become acceptable for new application domains, their performance must at least match those of file systems. Automatic optimization and parallelization may be crucial contributions to achieving this goal, in addition to the query execution techniques surveyed here.

Acknowledgements

José Blakeley, Rick Cole, Cathy Brand, Diane Davison, David Helman, Ann Linville, Bill McKenna, Gail Mitchell, Shengsong Ni, Barb Peters, Leonard D. Shapiro, the students of CU's "Readings in Database Systems" (Fall 1991), and David Maier's weekly reading group at the Oregon Graduate Institute gave many valuable comments on earlier drafts of this survey. — This paper is based on research partially supported by the National Science Foundation with grants IRI-8996270, IRI-8912618, IRI-9006348, and IRI-9116547, DARPA with contract DAAB07-91-C-Q518, Texas Instruments, Digital Equipment Corp., Intel Supercomputer Systems Division, Sequent Computer Systems, ADP, and the Oregon Advanced Computing Institute (OACIS).

References

- [1] I. Ahn and R. Snodgrass, "Partitioned storage for temporal databases", *Information Systems* 13, 4 (1988), 369.
- [2] J. Albert, "Algebraic Properties of Bag Data Types", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [3] D. P. Anderson, S. Tzou and G. S. Graham, "The DASH Virtual Memory System", Technical Report 88/461, UC Berkeley CS Division, November 1988.
- [4] M. M. Astrahan, M. W. Blasgen, D. D. Chamberlin, K. P. Eswaran, J. N. Gray, P. P. Griffiths, W. F. King, R. A. Lorie, P. R. McJones, J. W. Mehl, G. R. Putzolu, I. L. Traiger, B. W. Wade and V. Watson, "System

- R: A Relational Approach to Database Management”, *ACM Trans. on Database Systems* 1, 2 (June 1976), 97. Reprinted in M. Stonebraker, Readings in Database Systems, Morgan-Kaufman, San Mateo, CA, 1988.
- [5] M. M. Astrahan, M. Schkolnick and K. Y. Whang, “Approximating the number of unique values of an attribute without sorting”, *Information Systems* 12, 1 (1987), 11.
- [6] E. Babb, “Implementing a Relational Database by Means of Specialized Hardware”, *ACM Trans. on Database Systems* 4, 1 (March 1979), 1.
- [7] E. Babb, “Joined Normal Form: A Storage Encoding for Relational Databases”, *ACM Trans. on Database Systems* 7, 4 (December 1982), 588.
- [8] R. A. Baeza-Yates and P. A. Larson, “Performance of B+-Trees with Partial Expansions”, *IEEE Trans. on Knowledge and Data Eng.* 1, 2 (June 1989), 248.
- [9] F. Bancilhon and R. Ramakrishnan, “An Amateur’s Introduction to Recursive Query Processing Strategies”, *Proc. ACM SIGMOD Conf.*, Washington, DC, May 1986, 16. Reprinted in M. Stonebraker, Readings in Database Systems, Morgan-Kaufman, San Mateo, CA, 1988.
- [10] C. K. Baru and O. Frieder, “Database Operations in a Cube-Connected Multicomputer System”, *IEEE Trans. on Computers* 38, 6 (June 1989), 920.
- [11] D. S. Batory, T. Y. Leung and T. E. Wise, “Implementation Concepts for an Extensible Data Model and Data Language”, *ACM Trans. on Database Systems* 13, 3 (September 1988), 231.
- [12] D. S. Batory, J. R. Barnett, J. F. Garza, K. P. Smith, K. Tsukuda, B. C. Twichell and T. E. Wise, “GENESIS: An Extensible Database Management System”, *IEEE Trans. on Software Eng.* 14, 11 (November 1988), 1711.
- [13] B. Baugsto and J. Greipsland, “Parallel Sorting Methods for Large Data Volumes on a Hypercube Database Computer”, *Proc. Sixth Int’l Workshop on Database Machines*, Deauville, France, June 19-21, 1989.
- [14] R. Bayer and E. McCreighton, “Organisation and Maintenance of Large Ordered Indices”, *Acta Informatica* 1 (1972).
- [15] M. Beck, D. Bitton and W. K. Wilkinson, “Sorting Large Files on a Backend Multiprocessor”, *IEEE Trans. on Computers* 37 (1988), 769.
- [16] B. Becker, H. W. Six and P. Widmayer, “Spatial Priority Search: An Access Technique for Scaleless Maps”, *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 128.
- [17] D. A. Beckley, M. W. Evans and V. K. Raman, “Multikey Retrieval from K-d Trees and Quad-Trees”, *Proc. ACM SIGMOD Conf.*, Austin, TX, May 1985, 291.
- [18] N. Beckmann, H. P. Kriegel, R. Schneider and B. Seeger, “The R*-tree: An Efficient and Robust Access Method for Points and Rectangles”, *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 322.
- [19] T. Bell, I. H. Witten and J. G. Cleary, “Modelling for Text Compression”, *ACM Computing Surveys* 21, 4 (December 1989), 557.
- [20] P. A. Bernstein and N. Goodman, “Concurrency Control in Distributed Database Systems”, *ACM Computing Surveys* 13, 2 (June 1981), 185.
- [21] P. A. Bernstein, N. Goodman, E. Wong, C. L. Reeve and J. B. Rothnie, “Query Processing in a System for Distributed Databases (SDD-1)”, *ACM Trans. on Database Systems* 6, 4 (December 1981), 602.
- [22] P. A. Bernstein, V. Hadzilacos and N. Goodman, *Concurrency Control and Recovery in Database Systems*, Addison-Wesley, Reading, MA, 1987.
- [23] P. B. Berra, S. M. Chung and N. I. Hachem, “Computer Architecture for a Surrogate File to a Very Large Data/Knowledge Base”, *IEEE Computer* 20, 3 (March 1987), 25.
- [24] E. Bertino and W. Kim, “Indexing Techniques for Queries on Nested Objects”, *IEEE Trans. on Knowledge and Data Eng.* 1, 2 (June 1989), 196.
- [25] E. Bertino, “Optimization of Queries Using Nested Indices”, *Lecture Notes in Computer Science* 416 (March 1990), 44, Springer Verlag.
- [26] E. Bertino, “An Indexing Technique for Object-Oriented Databases”, *Proc. IEEE Conf. on Data Eng.*, Kobe, Japan, April 1991.
- [27] A. Bhide, “An Analysis of Three Transaction Processing Architectures”, *Proc. Int’l. Conf. on Very Large Data Bases*, Long Beach, CA, August 1988, 339.
- [28] A. Bhide and M. Stonebraker, “A Performance Comparison of Two Architectures for Fast Transaction Processing”, *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1988, 536.
- [29] D. Bitton and D. J. DeWitt, “Duplicate Record Elimination in Large Data Files”, *ACM Trans. on Database Systems* 8, 2 (June 1983), 255.
- [30] D. Bitton and J. Gray, “Disk Shadowing”, *Proc. Int’l. Conf. on Very Large Data Bases*, Long Beach, CA, August 1988, 331.

- [31] D. Bitton Friedland, “Design, Analysis, and Implementation of Parallel External Sorting Algorithms”, *Computer Sciences Technical Report 464* (January 1982), University of Wisconsin — Madison.
- [32] J. A. Blakeley, P. A. Larson and F. W. Tompa, “Efficiently Updating Materialized Views”, *Proc. ACM SIGMOD Conf.*, Washington, DC, May 1986, 61.
- [33] J. A. Blakeley, N. Coburn and P. A. Larson, “Updating Derived Relations: Detecting Irrelevant and Autonomously Computable Updates”, *ACM Trans. on Database Systems* 14, 3 (September 1989), 369.
- [34] J. A. Blakeley and N. L. Martin, “Join Index, Materialized View, and Hybrid Hash-Join: A Performance Analysis”, *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1990, 256.
- [35] M. Blasgen and K. Eswaran, “On the Evaluation of Queries in a Relational Database System”, *IBM Research Report RJ 1745*, San Jose, CA, April 8, 1976.
- [36] M. Blasgen and K. Eswaran, “Storage and Access in Relational Databases”, *IBM Systems Journal* 16, 4 (1977).
- [37] H. Boral, D. DeWitt, D. Friedland, N. Jarrell and W. Wilkinson, “Implementation of the Database Machine DIRECT”, *IEEE Trans. on Software Eng.* 8, 6 (November 1982), 533.
- [38] H. Boral and D. J. DeWitt, “Database Machines: An Idea Whose Time Has Passed? A Critique of the Future of Database Machines”, *Proc. Int’l. Workshop on Database Machines*, Munich, 1983.
- [39] H. Boral, “Parallelism in Bubba”, *Proc. Int’l. Symp. on Databases in Parallel and Distributed Systems*, Austin, TX, December 1988, 68.
- [40] H. Boral, W. Alexander, L. Clay, G. Copeland, S. Danforth, M. Franklin, B. Hart, M. Smith and P. Valduriez, “Prototyping Bubba, A Highly Parallel Database System”, *IEEE Trans. on Knowledge and Data Eng.* 2, 1 (March 1990), 4.
- [41] K. Bratbergsengen, “Hashing Methods and Relational Algebra Operations”, *Proc. Int’l. Conf. on Very Large Data Bases*, Singapore, August 1984, 323.
- [42] H. L. Bremers, “Hash Partitioning Performance Improved By Exploiting Skew and Dealing with Duplicates”, *M.S. Thesis, University of Colorado at Boulder*, 1991.
- [43] M. J. Carey, D. J. DeWitt, J. E. Richardson and E. J. Shekita, “Object and File Management in the EXODUS Extensible Database System”, *Proc. Int’l. Conf. on Very Large Data Bases*, Kyoto, Japan, August 1986, 91.
- [44] M. J. Carey, E. Shekita, G. Lapis, B. Lindsay and J. McPherson, “An Incremental Join Attachment for Starburst”, *Proc. Int’l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 662.
- [45] J. V. Carlis, “HAS: A Relational Algebra Operator, or Divide Is Not Enough to Conquer”, *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1986, 254.
- [46] J. L. Carter and M. N. Wegman, “Universal Classes of Hash Functions”, *Journal of Computers and System Science* 18, 2 (1979), 143.
- [47] U. P. Center, “NetCDF User’s Guide, An Interface for Data Access”, *NCAR Technical Note TS-334+1A*, Boulder, CO, April 1991. Version 1.11.
- [48] D. D. Chamberlin, M. M. Astrahan, M. W. Blasgen, J. N. Gray, W. F. King, B. G. Lindsay, R. Lorie, J. W. Mehl, T. G. Price, F. Putzolo, P. G. Selinger, M. Schkolnik, D. R. Slutz, I. L. Traiger, B. W. Wade and R. A. Yost, “A History and Evaluation of System R”, *Communications of the ACM* 24, 10 (October 1981), 632. Reprinted in M. Stonebraker, *Readings in Database Systems*, Morgan-Kaufman, San Mateo, CA, 1988.
- [49] D. D. Chamberlin, M. M. Astrahan, W. F. King, R. A. Lorie, J. W. Mehl, T. G. Price, M. Schkolnik, P. G. Selinger, D. R. Slutz, B. W. Wade and R. A. Yost, “Support for Repetitive Transactions and Ad Hoc Queries in System R”, *ACM Trans. on Database Systems* 6, 1 (March 1981), 70.
- [50] E. Chang and R. Katz, “Exploiting Inheritance and Structure Semantics for Effective Clustering and Buffering in an Object-Oriented DBMS”, *Proc. ACM SIGMOD Conf.*, Portland, OR, May-June 1989, 348.
- [51] P. P. Chen, “The Entity Relationship Model - Toward a Unified View of Data”, *ACM Trans. on Database Systems* 1, 1 (March 1976), 9. Reprinted in M. Stonebraker, *Readings in Database Systems*, Morgan-Kaufman, San Mateo, CA, 1988.
- [52] H. Chen and S. M. Kuck, “Combining Relational and Network Retrieval Methods”, *Proc. ACM SIGMOD Conf.*, Boston, MA, June 1984, 131.
- [53] P. M. Chen and D. A. Patterson, “Maximizing Performance in a Striped Disk Array”, *Proc. 17th Annual Int’l Symp. on Computer Architecture, ACM SIGARCH Computer Architecture News* 18, 2 (June 1990), 322.
- [54] M. Chen, M. Lo, P. S. Yu and H. C. Young, “Using Segmented Right-Deep Trees for the Execution of Pipelined Hash Joins”, *to appear in Proc. VLDB Conf.*, Vancouver, BC, Canada, 1992.

- [55] J. Cheng, D. Haderle, R. Hedges, B. Iyer, T. Messinger, C. Mohan and Y. Wang, "An Efficient Hybrid Join Algorithm: Design, Prototype, Modelling and Measurement", *Proc. IEEE Conf. on Data Eng.*, Kobe, Japan, April 1991.
- [56] J. R. Cheng and A. R. Hurson, "Effective Clustering of Complex Objects in Object-Oriented Databases", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 22.
- [57] D. R. Cheriton, H. A. Goosen and P. D. Boyle, "Paradigm: A Highly Scalable Shared-Memory Multicomputer", *IEEE Computer* 24, 2 (February 1991), 33.
- [58] D. M. Chiu and Y. C. Ho, "A Methodology for Interpreting Tree Queries Into Optimal Semi-Join Expressions", *Proc. ACM SIGMOD Conf.*, Santa Monica, CA, May 1980, 169.
- [59] H. T. Chou, "Buffer Management of Database Systems", *Ph.D. Thesis*, May 1985.
- [60] H. T. Chou and D. J. DeWitt, "An Evaluation of Buffer Management Strategies for Relational Database Systems", *Proc. Int'l. Conf. on Very Large Data Bases*, Stockholm, Sweden, August 1985, 127. Reprinted in M. Stonebraker, Readings in Database Systems, Morgan-Kaufman, San Mateo, CA, 1988.
- [61] S. Christodoulakis, "Implications of Certain Assumptions in Database Performance Evaluation", *ACM Trans. on Database Systems* 9, 2 (June 1984), 163.
- [62] S. M. Chung and P. B. Berra, "A Comparison of Concatenated and Superimposed Code Word Surrogate Files for Very Large Data/Knowledge Bases", *Lecture Notes in Computer Science* 303 (April 1988), 364, Springer Verlag.
- [63] S. Cluet, C. Delobel, C. Lecluse and P. Richard, "Reloops, an Algebra Based Query Language for an Object-Oriented Database System", *Proc. First Int'l. Conf. on Deductive and Object-Oriented Databases*, Kyoto, Japan, December 4-6, 1989.
- [64] D. Comer, "The Ubiquitous B-Tree", *ACM Computing Surveys* 11, 2 (June 1979).
- [65] G. Copeland, W. Alexander, E. Boughter and T. Keller, "Data Placement in Bubba", *Proc. ACM SIGMOD Conf.*, Chicago, IL, June 1988, 99.
- [66] G. V. Cormack, "Data Compression In a Database System", *Communications of the ACM* 28, 12 (December 1985), 1336.
- [67] D. Daniels and P. Ng, "Distributed Query Compilation and Processing in R*", *IEEE Database Eng.* 5, 3 (September 1982).
- [68] S. Daniels, G. Graefe, T. Keller, D. Maier, D. Schmidt and B. Vance, "Query Optimization in Revelation, an Overview", *IEEE Database Eng.* 14, 2 (June 1991).
- [69] D. J. DeWitt and P. B. Hawthorn, "A Performance Evaluation of Database Machine Architectures", *Proc. Int'l. Conf. on Very Large Data Bases*, Cannes, France, September 1981, 199.
- [70] D. J. DeWitt, R. Katz, F. Olken, L. Shapiro, M. Stonebraker and D. Wood, "Implementation Techniques for Main Memory Database Systems", *Proc. ACM SIGMOD Conf.*, Boston, MA, June 1984, 1.
- [71] D. J. DeWitt and R. H. Gerber, "Multiprocessor Hash-Based Join Algorithms", *Proc. Int'l. Conf. on Very Large Data Bases*, Stockholm, Sweden, August 1985, 151.
- [72] D. J. DeWitt, R. H. Gerber, G. Graefe, M. L. Heytens, K. B. Kumar and M. Muralikrishna, "GAMMA - A High Performance Dataflow Database Machine", *Proc. Int'l. Conf. on Very Large Data Bases*, Kyoto, Japan, August 1986, 228. Reprinted in M. Stonebraker, Readings in Database Systems, Morgan-Kaufman, San Mateo, CA, 1988.
- [73] D. J. DeWitt, S. Ghandeharizadeh and D. Schneider, "A Performance Analysis of the GAMMA Database Machine", *Proc. ACM SIGMOD Conf.*, Chicago, IL, June 1988, 350.
- [74] D. J. DeWitt, S. Ghandeharizadeh, D. Schneider, A. Bricker, H. I. Hsiao and R. Rasmussen, "The Gamma Database Machine Project", *IEEE Trans. on Knowledge and Data Eng.* 2, 1 (March 1990), 44.
- [75] D. DeWitt, J. Naughton and D. Schneider, "An evaluation of non-equijoin algorithms", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [76] D. J. DeWitt, "The Wisconsin Benchmark: Past, Present, and Future", in *Database and Transaction Processing Systems Performance Handbook*, J. Gray (editor), Morgan-Kaufman, San Mateo, CA, 1991.
- [77] A. Deshpande and D. Van Gucht, "An Implementation for Nested Relational Databases", *Proc. Int'l. Conf. on Very Large Data Bases*, Long Beach, CA, August 1988, 76.
- [78] V. Deshpande and P. A. Larson, *An Algebra for Nested Relations With Support for Nulls and Aggregates*, Computer Science Dept., Univ. of Waterloo, Waterloo, Ontario, Canada, April 1991.
- [79] D. Dewitt, J. Naughton and D. Schneider, "Parallel Sorting on a Shared-Nothing Architecture using Probabilistic Splitting", *Proc. Int'l. Conf. on Parallel and Distributed Information Systems*, Miami Beach, FL, December 1991.
- [80] W. Effelsberg and T. Haerder, "Principles of Database Buffer Management", *ACM Trans. on Database Systems* 9, 4 (December 1984), 560.

- [81] S. Englert, J. Gray, R. Kocher and P. Shah, "A Benchmark of NonStop SQL Release 2 Demonstrating Near-Linear Speedup and Scaleup on Large Databases", *Tandem Computer Systems Technical Report 89.4* (May 1989).
- [82] R. Epstein, "Techniques for Processing of Aggregates in Relational Database Systems", *UCB/Electronics Research Lab. Memorandum M79/8* (February 1979), University of California.
- [83] R. Epstein and M. Stonebraker, "Analysis of Distributed Data Base Processing Strategies", *Proc. Int'l. Conf. on Very Large Data Bases*, Montreal, Canada, October 1980, 92.
- [84] R. Fagin, J. Nievergelt, N. Pippenger and H. R. Strong, "Extendible Hashing: A Fast Access Method for Dynamic Files", *ACM Trans. on Database Systems* 4, 3 (September 1979), 315.
- [85] C. Faloutsos, R. Ng and T. Sellis, "Predictive load control for flexible buffer allocation", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [86] M. T. Fang, R. C. T. Lee and C. C. Chang, "The Idea of Declustering and Its Applications", *Proc. Int'l. Conf. on Very Large Data Bases*, Kyoto, Japan, August 1986, 181.
- [87] S. Finkelstein, M. Schkolnick and P. Tiberio, "Physical Database Design for Relational Databases", *ACM Trans. on Database Systems* 13, 1 (March 1988).
- [88] M. Fontenot, "Software Congestion, Mobile Servers, and the Hyperbolic Model", *IEEE Trans. on Software Eng.* 15, 8 (August 1989), 947.
- [89] J. C. Freytag and N. Goodman, "On the Translation of Relational Queries into Iterative Programs", *ACM Trans. on Database Systems* 14, 1 (March 1989), 1.
- [90] S. Fushimi, M. Kitsuregawa and H. Tanaka, "An Overview of The System Software of A Parallel Relational Database Machine GRACE", *Proc. Int'l. Conf. on Very Large Data Bases*, Kyoto, Japan, August 1986, 209.
- [91] H. Garcia-Molina and K. Salem, "The Impact of Disk Striping on Reliability", *Princeton University Computer Science Technical Report*, January 1988.
- [92] R. Gerber, "Dataflow Query Processing using Multiprocessor Hash-Partitioned Algorithms", *Ph.D. Thesis*, October 1986.
- [93] R. H. Gerber and D. J. DeWitt, "The Impact of Hardware and Software Alternatives on the Performance of the Gamma Database Machine", *Computer Sciences Technical Report 708* (July 1987), University of Wisconsin — Madison.
- [94] S. Ghandeharizadeh and D. J. DeWitt, "Hybrid-Range Partitioning Strategy: A New Declustering Strategy for Multiprocessor Database Machines", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 481.
- [95] S. Ghandeharizadeh and D. J. DeWitt, "A Multiuser Performance Analysis of Alternative Clustering Strategies", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1990, 466.
- [96] S. Ghandeharizadeh, L. Ramos, Z. Asad and W. Qureshi, "Object Placement in Parallel Hypermedia Systems", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [97] G. A. Gibson, L. Hellerstein, R. M. Karp, R. H. Katz and D. A. Patterson, "Failure Correction Techniques for Large Disk Arrays", *Third Int'l Conf. on Architectural Support for Programming Languages and Operating Systems*, April 1989, 123.
- [98] J. R. Goodman and P. J. Woest, "The Wisconsin Multicube: A New Large-Scale Cache-Coherent Multiprocessor", *Computer Sciences Technical Report 766* (April 1988), University of Wisconsin — Madison.
- [99] M. G. Gouda and U. Dayal, "Optimal Semijoin Schedules for Query Processing in Local Distributed Database Systems", *Proc. ACM SIGMOD Conf.*, Ann Arbor, MI, April-May 1981, 164.
- [100] G. Graefe and D. J. DeWitt, "The EXODUS Optimizer Generator", *Proc. ACM SIGMOD Conf.*, San Francisco, CA, May 1987, 160.
- [101] G. Graefe and D. Maier, "Query Optimization in Object-Oriented Database Systems: A Prospectus", in *Advances in Object-Oriented Database Systems*, vol. 334, K. R. Dittrich (editor), Springer-Verlag, September 1988, 358.
- [102] G. Graefe, "Relational Division: Four Algorithms and Their Performance", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1989, 94.
- [103] G. Graefe and K. Ward, "Dynamic Query Evaluation Plans", *Proc. ACM SIGMOD Conf.*, Portland, OR, May-June 1989, 358.
- [104] G. Graefe, "Encapsulation of Parallelism in the Volcano Query Processing System", *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 102.
- [105] G. Graefe, "Parallel External Sorting in Volcano", *CU Boulder Computer Science Technical Report 459*, 1990.

- [106] G. Graefe and L. D. Shapiro, “Data Compression and Database Performance”, *Proc. ACM/IEEE-CS Symp. on Applied Computing*, Kansas City, MO, April 1991.
- [107] G. Graefe, “Heap-Filter Merge Join: A New Algorithm for Joining Medium-Size Inputs”, *IEEE Trans. on Software Eng.* 17, 9 (September 1991), 979.
- [108] G. Graefe and W. J. McKenna, “The Volcano Optimizer Generator”, *submitted for publication*, 1991. Also CU Boulder Computer Science Technical Report 563.
- [109] G. Graefe and D. L. Davison, “Encapsulation of Parallelism and Architecture-Independence in Extensible Database Query Processing”, *submitted for publication*, 1991. Also CU Boulder Computer Science Technical Report 559.
- [110] G. Graefe and H. L. Bremers, “Exploiting Skew to Improve Hybrid Hash Join Performance”, *in preparation*, 1992.
- [111] G. Graefe, “Volcano, An Extensible and Parallel Dataflow Query Processing System”, *to appear in IEEE Trans. on Knowledge and Data Eng.*, 1992. A more detailed version is available as CU Boulder Computer Science Technical Report 481, July 1990.
- [112] G. Graefe, A. Linville and L. D. Shapiro, Sort versus Hash Revisited, Also CU Boulder Computer Science Technical Report 534, July 1991, 1992.
- [113] G. Graefe and R. L. Cole, “Fast Algorithms for Universal Quantification in Large Databases”, *in preparation*, 1992.
- [114] G. Graefe and S. S. Thakkar, “Tuning a Parallel Database Algorithm on a Shared-Memory Multiprocessor”, *to appear in Software—Practice and Experience*, . Also CU Boulder Computer Science Technical Report 470, 1990.
- [115] J. Gray, “A Census of Tandem System Availability Between 1985 and 1990”, *Tandem Computers Technical Report 90.1*, January 1990.
- [116] J. Gray, B. Horst and M. Walker, “Parity Striping of Disc Arrays: Low-Cost Reliable Storage with Acceptable Throughput”, *Proc. Int’l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 148.
- [117] J. Gray and A. Reuter, *Transaction Processing: Concepts and Techniques*, Morgan-Kaufman, San Mateo, CA, 1991.
- [118] O. Guenther and J. Bilmes, “Tree-Based Access Methods for Spatial Databases: Implementation and Performance Evaluation”, *IEEE Trans. on Knowledge and Data Eng.* 3, 3 (September 1991), 342.
- [119] L. Guibas and R. Sedgewick, “A Dichromatic Framework for Balanced Trees”, *Proc. 19th Symp. on the Foundations of Computer Science*, 1978.
- [120] H. Gunadhi and A. Segev, “A Framework for Query Optimization in Temporal Databases”, *Proc. Fifth Int’l. Conf. on Statistical and Scientific Database Management*, April 1990.
- [121] O. Gunther and E. Wong, “A Dual Space Representation for Geometric Data”, *Proc. Int’l. Conf. on Very Large Data Bases*, Brighton, England, August 1987, 501.
- [122] O. Gunther, “The Design of the Cell Tree: An Object-Oriented Index Structure for Geometric Databases”, *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1989, 598.
- [123] A. Guttman, “R-Trees: A Dynamic Index Structure for Spatial Searching”, *Proc. ACM SIGMOD Conf.*, Boston, MA, June 1984, 47. Reprinted in M. Stonebraker, *Readings in Database Systems*, Morgan-Kaufman, San Mateo, CA, 1988.
- [124] L. M. Haas, P. G. Selinger, E. Bertino, D. Daniels, B. Lindsay, G. Lohman, Y. Masunaga, C. Mohan, P. Ng, P. Wilms and R. Yost, *R*: A Research Project on Distributed Relational DBMS*, IBM Research Division, San Jose CA, October 1982.
- [125] L. Haas, J. Freytag, G. Lohman and H. Pirahesh, “Extensible Query Processing in Starburst”, *Proc. ACM SIGMOD Conf.*, Portland, OR, May-June 1989, 377.
- [126] L. Haas, W. Chang, G. Lohman, J. McPherson, P. F. Wilms, G. Lapis, B. Lindsay, H. Pirahesh, M. J. Carey and E. Shekita, “Starburst Mid-Flight: As the Dust Clears”, *IEEE Trans. on Knowledge and Data Eng.* 2, 1 (March 1990), 143.
- [127] T. Haerder and A. Reuter, “Principles of Transaction-Oriented Database Recovery”, *ACM Computing Surveys* 15, 4 (December 1983). Reprinted in M. Stonebraker, *Readings in Database Systems*, Morgan-Kaufman, San Mateo, CA, 1988.
- [128] R. B. Hagmann, “An Observation on Database Buffering Performance Metrics”, *Proc. Int’l. Conf. on Very Large Data Bases*, Kyoto, Japan, August 1986, 289.
- [129] M. Hammer and A. Chan, “Index Selection in a Self-Adaptive Data Base Management System”, *Proc. ACM SIGMOD Conf.*, 1976, 1.
- [130] R. W. Hamming, *Digital Filters*, Prentice-Hall, Englewood Cliffs, NJ, 1977.

- [131] E. N. Hanson, "A Performance Analysis of View Materialization Strategies", *Proc. ACM SIGMOD Conf.*, San Francisco, CA, May 1987, 440.
- [132] L. Harada, M. Nakano, M. Kitsuregawa and M. Takagi, "Query Processing Method for Multi-Attribute Clustered Relations", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 59.
- [133] A. Henrich, H. W. Six and P. Widmayer, "The LSD tree: spatial access to multi-dimensional point and non-point objects", *Proc. Int'l. Conf. on Very Large Data Bases*, Amsterdam, The Netherlands, 1989, 45.
- [134] W. Hong and M. Stonebraker, "Optimization of Parallel Query Execution Plans in XPRS", *Proc. Int'l. Conf. on Parallel and Distributed Information Systems*, Miami Beach, FL, December 1991.
- [135] W. Hou, G. Ozsoyoglu and B. Taneja, "Statistical Estimators for Relational Algebra Expressions", *Proc. SIGACT News-SIGMOD Symp. on Principles of Database Systems*, Austin, TX, March 1988, 276.
- [136] W. C. Hou and G. Ozsoyoglu, "Statistical Estimators for Aggregate Relational Algebra Queries", *ACM Trans. on Database Systems* 16, 4 (December 1991), 600.
- [137] W. C. Hou, G. Ozsoyoglu and E. Dogdu, "Error-Constrained COUNT Query Evaluation in Relational Databases", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 278.
- [138] H. I. Hsiao and D. J. DeWitt, "Chained Declustering: A New Availability Strategy for Multiprocessor Database Machines", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1990, 456.
- [139] H. Hsiao and D. DeWitt, "A Performance Study of Three High Availability Data Replication Strategies", *Proc. Int'l. Conf. on Parallel and Distributed Information Systems*, Miami Beach, FL, December 1991.
- [140] K. A. Hua and C. Lee, "An Adaptive Data Placement Scheme for Parallel Database Computer Systems", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 493.
- [141] K. Hua and C. Lee, "Handling Data Skew in Multicomputer Database Systems Using Partition Tuning", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [142] S. E. Hudson and R. King, "Cactis: A Self-Adaptive, Concurrent Implementation of an Object-Oriented Database Management System", *ACM Trans. on Database Systems* 14, 3 (September 1989), 291.
- [143] A. Hutflesz, H. W. Six and P. Widmayer, "Twin Grid Files: Space Optimizing Access Schemes", *Proc. ACM SIGMOD Conf.*, Chicago, IL, June 1988, 183.
- [144] A. Hutflesz, H. W. Six and P. Widmayer, "The Twin Grid File: A Nearly Space Optimal Index Structure", *Lecture Notes in Computer Science* 303 (April 1988), 352, Springer Verlag.
- [145] A. Hutflesz, H. W. Six and P. Widmayer, "The R-File: An Efficient Access Structure for Proximity Queries", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1990, 372.
- [146] Y. E. Ioannidis and S. Christodoulakis, "On the Propagation of Errors in the Size of Join Results", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 268.
- [147] B. R. Iyer and D. M. Dias, "System Issues in Parallel Sorting for Database Systems", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1990, 246.
- [148] H. V. Jagadish, "A Compression Technique to Materialize Transitive Closure", *ACM Trans. on Database Systems* 15, 4 (December 1990), 558.
- [149] H. V. Jagadish, "A Retrieval Technique for Similar Shapes", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 208.
- [150] M. Jarke and J. Koch, "Query Optimization in Database Systems", *ACM Computing Surveys* 16, 2 (June 1984), 111.
- [151] A. Jhingran, "Precomputation in a Complex Object Environment", *Proc. IEEE Conf. on Data Eng.*, Kobe, Japan, April 1991.
- [152] S. Kao, "DECIDES: An Expert System Tool for Physical Database Design", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1986, 671.
- [153] R. H. Katz and E. Wong, "Resolving Conflicts in Global Storage Design Through Replication", *ACM Trans. on Database Systems* 8, 1 (March 1983), 110.
- [154] T. Keller, G. Graefe and D. Maier, "Efficient Assembly of Complex Objects", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 148.
- [155] T. Keller, G. Graefe and D. Maier, "Efficient Query Processing in Object-Oriented Database Management Systems: Revealing and Assembly", *submitted for publication*, 1992.
- [156] A. Kemper and M. Wallrath, "An Analysis of Geometric Modeling in Database Systems", *ACM Computing Surveys* 19, 1 (March 1987), 47.
- [157] A. Kemper and G. Moerkotte, "Advanced Query Processing in Object Bases Using Access Support Relations", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 290.
- [158] A. Kemper and G. Moerkotte, "Access Support in Object Bases", *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 364.

- [159] A. Kemper, C. Kilger and G. Moerkotte, "Function Materialization in Object Bases", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 258.
- [160] B. W. Kernighan and D. N. Ritchie, *The C Programming Language*, Prentice-Hall, Englewood Cliffs, NJ, 1978.
- [161] W. Kim, "A New Way to Compute the Product and Join of Relations", *Proc. ACM SIGMOD Conf.*, Santa Monica, CA, May 1980, 179.
- [162] M. Kitsuregawa, H. Tanaka and T. Motooka, "Application of Hash to Data Base Machine and Its Architecture", *New Generation Computing* 1, 1 (1983).
- [163] M. Kitsuregawa, W. Yang and S. Fushimi, "Evaluation of 18-Stage Pipeline Hardware Sorter", *Proc. Sixth Int'l Workshop on Database Machines*, Deauville, France, June 19-21, 1989.
- [164] M. Kitsuregawa, M. Nakayama and M. Takagi, "The effect of bucket size tuning in the dynamic hybrid GRACE hash join method", *Proc. Int'l. Conf. on Very Large Data Bases*, Amsterdam, The Netherlands, 1989, 257.
- [165] M. Kitsuregawa and Y. Ogawa, "Bucket Spreading Parallel Hash: A New, Robust, Parallel Hash Join Method for Skew in the Super Database Computer (SDC)", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 210.
- [166] A. Klug, "Equivalence of Relational Algebra and Relational Calculus Query Languages Having Aggregate Functions", *Journal of the ACM* 29, 3 (July 1982), 699.
- [167] D. Knuth, *The Art of Computer Programming, Vol. III: Sorting and Searching*, Addison-Wesley, Reading, MA, 1973.
- [168] C. P. Kolovson and M. Stonebraker, "Segment Indexes: Dynamic Indexing Techniques for Multi-Dimensional Interval Data", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 138.
- [169] R. P. Kooi, "The Optimization of Queries in Relational Databases", *Ph.D. Thesis, Case Western Reserve University*, September 1980.
- [170] R. P. Kooi and D. Frankforth, "Query Optimization in Ingres", *IEEE Database Eng.* 5, 3 (September 1982), 2.
- [171] H. P. Kriegel and B. Seeger, "Multidimensional Dynamic Hashing Is Very Efficient for Nonuniform Record Distributions", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1987, 10.
- [172] H. P. Kriegel and B. Seeger, "PLOP-Hashing: A Grid File without Directory", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1988, 369.
- [173] R. Krishnamurthy, H. Boral and C. Zaniolo, "Optimization of Nonrecursive Queries", *Proc. Int'l. Conf. on Very Large Data Bases*, Kyoto, Japan, August 1986, 128.
- [174] M. S. Lakshmi and P. S. Yu, "Effectiveness of Parallel Joins", *IEEE Trans. on Knowledge and Data Eng.* 2, 4 (December 1990), 410.
- [175] M. S. Lakshmi and P. S. Yu, "Effect of Skew on Join Performance in Parallel Architectures", *Proc. Int'l. Symp. on Databases in Parallel and Distributed Systems*, Austin, TX, December 1988, 107.
- [176] S. Lanka and E. Mays, "Fully Persistent B+-trees", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 426.
- [177] P. Larson and H. Yang, "Computing Queries from Derived Relations", *Proc. Int'l. Conf. on Very Large Data Bases*, Stockholm, Sweden, August 1985, 259.
- [178] D. A. Lelewer and D. S. Hirschberg, "Data Compression", *ACM Computing Surveys* 19, 3 (September 1987), 261.
- [179] J. Li, D. Rotem and H. Wong, "A New Compression Method with Fast Searching on Large Data Bases", *Proc. Int'l. Conf. on Very Large Data Bases*, Brighton, England, August 1987, 311.
- [180] D. Lomet and B. Salzberg, "The Performance of a Multiversion Access Method", *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 353.
- [181] D. B. Lomet and B. Salzberg, "The hB-Tree: A Multiattribute Indexing Method with Good Guaranteed Performance", *ACM Trans. on Database Systems* 15, 4 (December 1990).
- [182] R. A. Lorie and J. F. Nilsson, "An Access Specification Language for a Relational Database Management System", *IBM Journal of Research and Development* 23, 3 (May 1979), 286.
- [183] R. A. Lorie and H. C. Young, "A low communication sort algorithm for a parallel database machine", *Proc. Int'l. Conf. on Very Large Data Bases*, Amsterdam, The Netherlands, 1989, 125.
- [184] C. A. Lynch and E. B. Brownrigg, "Application of Data Compression to a Large Bibliographic Data Base", *Proc. Int'l. Conf. on Very Large Data Bases*, Cannes, France, September 1981, 435.
- [185] L. F. Mackert and G. M. Lohman, "Index Scans Using a Finite LRU Buffer: A Validated I/O Model", *ACM Trans. on Database Systems* 14, 3 (September 1989), 401.

- [186] D. Maier and J. Stein, "Indexing in an Object-Oriented DBMS", *Proc. Int'l Workshop on Object-Oriented Database Systems*, Pacific Grove, CA, September 1986, 171.
- [187] M. V. Mannino, P. Chu and T. Sager, "Statistical Profile Estimation in Database Systems", *ACM Computing Surveys* 20, 3 (September 1988).
- [188] L. E. McKenzie and R. T. Snodgrass, "Evaluation of Relational Algebras Incorporating the Time Dimension in Databases", *ACM Computing Surveys* 23, 4 (December 1991).
- [189] C. Medeiros and F. Tompa, "Understanding the Implications of View Update Policies", *Proc. Int'l. Conf. on Very Large Data Bases*, Stockholm, Sweden, August 1985, 316.
- [190] J. Menon, "A Study of Sort Algorithms for Multiprocessor Database Machines", *Proc. Int'l. Conf. on Very Large Data Bases*, Kyoto, Japan, August 1986, 197.
- [191] P. Mishra and M. H. Eich, "Join Processing in Relational Databases", *ACM Computing Surveys* 24, 1 (March 1992), 63.
- [192] B. Mitschang, "Extending the relational algebra to capture complex objects", *Proc. Int'l. Conf. on Very Large Data Bases*, Amsterdam, The Netherlands, 1989, 297.
- [193] A. Motro, "An Access Authorization Model for Relational Databases Based on Algebraic Manipulation of View Definitions", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1989, 339.
- [194] M. Nakayama, M. Kitsuregawa and M. Takagi, "Hash-Partitioned Join Method Using Dynamic Destaging Strategy", *Proc. Int'l. Conf. on Very Large Data Bases*, Long Beach, CA, August 1988, 468.
- [195] P. M. Neches, "Hardware Support for Advanced Data Management Systems", *IEEE Computer* 17, 11 (November 1984), 29.
- [196] P. M. Neches, "The Ynet: An Interconnect Structure for a Highly Concurrent Data Base Computer System", *Proc. 2nd Symp. on the Frontiers of Massively Parallel Computation*, Fairfax, October 1988.
- [197] L. Neugebauer, "Optimization and Evaluation of Database Queries Including Embedded Interpolation Procedures", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 118.
- [198] R. Ng, C. Faloutsos and T. Sellis, "Flexible Buffer Allocation Based on Marginal Gains", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 387.
- [199] F. Olken and D. Rotem, "Rearranging Data to Maximize the Efficiency of Compression", *Journal of Computer and System Sciences* 38, 2 (1989), 405.
- [200] E. Omiecinski, "Incremental File Reorganization Schemes", *Proc. Int'l. Conf. on Very Large Data Bases*, Stockholm, Sweden, August 1985, 346.
- [201] E. Omiecinski and E. Lin, "Hash-Based and Index-Based Join Algorithms for Cube and Ring Connected Multicomputers", *IEEE Trans. on Knowledge and Data Eng.* 1, 3 (September 1989), 329.
- [202] E. Omiecinski and P. Scheuermann, "A Parallel Algorithm for Record Clustering", *ACM Trans. on Database Systems* 15, 4 (December 1990), 599.
- [203] E. Omiecinski, "Performance Analysis of A Load Balancing Relational Hashing Join Algorithm for a Shared-Memory Multiprocessor", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [204] K. Ono and G. M. Lohman, "Measuring the Complexity of Join Enumeration in Query Optimization", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 314.
- [205] J. Ousterhout, "Why Aren't Operating Systems Getting Faster as Fast as Hardware?", *USENIX Summer Conference*, Anaheim, CA, June 1990.
- [206] M. T. Ozsu and P. Valduriez, "Distributed Database Systems: Where Are We Now?", *IEEE Computer* 24, 8 (August 1991), 68.
- [207] M. T. Ozsu and P. Valduriez, *Principles of Distributed Database Systems*, Prentice-Hall, Englewood Cliffs, NJ, 1991.
- [208] M. Palmer and S. Zdonik, "FIDO: A Cache that Learns to Fetch", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [209] D. A. Patterson, G. Gibson and R. H. Katz, "A Case for Redundant Arrays of Inexpensive Disks (RAID)", *Proc. ACM SIGMOD Conf.*, Chicago, IL, June 1988, 109.
- [210] H. Pirahesh, C. Mohan, J. Cheng, T. S. Liu and P. Selinger, "Parallelism in Relational Data Base Systems: Architectural Issues and Design Approaches", *Proc. Int'l. Symp. on Databases in Parallel and Distributed Systems*, Dublin, Ireland, July 1990, 4.
- [211] X. Qian and G. Wiederhold, "Incremental Recomputation of Active Relational Expressions", *IEEE Trans. on Knowledge and Data Eng.* 3, 3 (September 1991), 337.
- [212] R. K. Rew and G. P. Davis, "The Unidata netCDF: Software for Scientific Data Access", *Sixth Int'l. Conf. on Interactive Information and Processign Systems for Meteorology, Oceanography, and Hydrology*, Anaheim, CA, February 1990.

- [213] J. P. Richardson, H. Lu and K. Mikkilineni, "Design and Evaluation of Parallel Pipelined Join Algorithms", *Proc. ACM SIGMOD Conf.*, San Francisco, CA, May 1987, 399.
- [214] J. E. Richardson and M. J. Carey, "Programming Constructs for Database System Implementation in EXODUS", *Proc. ACM SIGMOD Conf.*, San Francisco, CA, May 1987, 208.
- [215] J. T. Robinson, "The K-D-B-Tree: A Search Structure For Large Multidimensional Dynamic Indices", *Proc. ACM SIGMOD Conf.*, Ann Arbor, MI, April-May 1981, 10.
- [216] A. Rosenthal and D. S. Reiner, "Querying Relational Views of Networks", in *Query Processing in Database Systems*, W. K. D. S. R. D. S. Batory (editor), Springer, Berlin, 1985, 109.
- [217] A. Rosenthal, C. Rich and M. Scholl, "Reducing Duplicate Work in Relational Join(s): A Modular Approach Using Nested Relations", *ETH Technical Report*, Zurich, Switzerland, June 1991.
- [218] D. Rotem and A. Segev, "Physical Organization of Temporal Data", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1987, 547.
- [219] J. B. Rothnie, P. A. Bernstein, S. Fox, N. Goodman, M. Hammer, T. A. Landers, C. Reeve, D. W. Shipman and E. Wong, "Introduction to a System for Distributed Databases (SDD-1)", *ACM Trans. on Database Systems* 5, 1 (March 1980).
- [220] N. Roussopoulos, "An Incremental Access Method for ViewCache: Concept, Algorithms, and Cost Analysis", *ACM Trans. on Database Systems* 16, 3 (September 1991), 535.
- [221] N. Roussopoulos and H. Kang, "A Pipeline N-way Join Algorithm Based on the 2-way Semijoin Program", *IEEE Trans. on Knowledge and Data Eng.* 3, 4 (December 1991), 486.
- [222] S. S. Ruth and P. J. Keutzer, "Data compression for business files", *Datamation* 18 (September 1972), 62.
- [223] G. M. Sacco and M. Schkolnik, "A Mechanism for Managing the Buffer Pool in a Relational Database System Using the Hot Set Model", *Proc. Int'l. Conf. on Very Large Data Bases*, Mexico City, Mexico, September 1982, 257.
- [224] G. M. Sacco and M. Schkolnik, "Buffer Management in Relational Database Systems", *ACM Trans. on Database Systems* 11, 4 (December 1986), 473.
- [225] G. Sacco, "Index Access with a Finite Buffer", *Proc. Int'l. Conf. on Very Large Data Bases*, Brighton, England, August 1987, 301.
- [226] K. Salem and H. Garcia-Molina, "Disk Striping", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1986, 336.
- [227] B. Salzberg, *File Structures: An Analytic Approach*, Prentice-Hall, Englewood Cliffs, NJ, 1988.
- [228] B. Salzberg, A. Tsukerman, J. Gray, M. Stewart, S. Uren and B. Vaughan, "FastSort: An Distributed Single-Input Single-Output External Sort", *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 94.
- [229] B. Salzberg, "Merging Sorted Runs Using Large Main Memory", *Acta Informatica* 27 (1990), 195, Springer.
- [230] H. Samet, "The Quadtree and Related Hierarchical Data Structures", *ACM Computing Surveys* 16, 2 (June 1984), 187.
- [231] H. J. Schek and M. H. Scholl, "The relational model with relation-valued attributes", *Information Systems* 11, 2 (1986), 137.
- [232] D. Schneider and D. DeWitt, "A Performance Evaluation of Four Parallel Join Algorithms in a Shared-Nothing Multiprocessor Environment", *Proc. ACM SIGMOD Conf.*, Portland, OR, May-June 1989, 110.
- [233] D. Schneider, "Complex Query Processing in Multiprocessor Database Machines", *Ph.D. Thesis, Computer Sciences Technical Report 965*, 1990.
- [234] D. A. Schneider and D. J. DeWitt, "Tradeoffs in Processing Complex Join Queries via Hashing in Multiprocessor Database Machines", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 469.
- [235] D. A. Schneider, "Bit Filtering and Multi-Way Join Query Processing", *unpublished manuscript*, Palo Alto, CA, 1991.
- [236] B. Seeger and P. A. Larson, "Multi-Disk B-trees", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 436.
- [237] A. Segev and H. Gunadhi, "Event-join optimization in temporal relational databases", *Proc. Int'l. Conf. on Very Large Data Bases*, Amsterdam, The Netherlands, 1989, 205.
- [238] A. Segev, "Query Processing Algorithms for Temporal Intersection Joins", *Proc. IEEE Conf. on Data Eng.*, Kobe, Japan, April 1991.
- [239] P. G. Selinger, M. M. Astrahan, D. D. Chamberlin, R. A. Lorie and T. G. Price, "Access Path Selection in a Relational Database Management System", *Proc. ACM SIGMOD Conf.*, Boston, MA, May-June 1979, 23. Reprinted in M. Stonebraker, *Readings in Database Systems*, Morgan-Kaufman, San Mateo, CA, 1988.

- [240] T. K. Sellis, "Efficiently Supporting Procedures in Relational Database Systems", *Proc. ACM SIGMOD Conf.*, San Francisco, CA, May 1987, 278.
- [241] S. Seshadri and J. F. Naughton, "Sampling Issues in Parallel Database Systems", *Proc. Int'l. Conf. on Extending Database Technology*, Vienna, Austria, March 1992.
- [242] D. Severance and G. Lohman, "Differential Files: Their Application to the Maintenance of Large Databases", *ACM Trans. on Database Systems* 1, 3 (September 1976).
- [243] D. G. Severance, "A practitioner's guide to data base compression", *Information Systems* 8, 1 (1983), 51.
- [244] L. D. Shapiro, "Join Processing in Database Systems with Large Main Memories", *ACM Trans. on Database Systems* 11, 3 (September 1986), 239.
- [245] G. M. Shaw and S. B. Zdonik, "An object-oriented query algebra", in *Proc. of the 2nd Intl. Workshop on Database Programming Languages*, R. H. R. M. D. Stemple (editor), Morgan Kaufmann, Gleneden Beach, Oregon, June 1989, 103.
- [246] G. Shaw and S. Zdonik, "A Object-Oriented Query Algebra", *IEEE Database Eng.* 12, 3 (September 1989), 29.
- [247] G. M. Shaw and S. B. Zdonik, "A Query Algebra for Object-Oriented Databases", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1990, 154.
- [248] E. J. Shekita and M. J. Carey, "A Performance Evaluation of Pointer-Based Joins", *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 300.
- [249] S. W. Sherman and R. S. Brice, "Performance of a Database Manager in a Virtual Memory System", *ACM Trans. on Database Systems* 1, 4 (December 1976), 317.
- [250] A. P. Sheth and J. A. Larson, "Federated Database Systems for Managing Distributed, Heterogeneous, and Autonomous Databases", *ACM Computing Surveys* 22, 3 (September 1990), 183.
- [251] A. Sikeler, "VAR-PAGE-LRU: A Buffer Replacement Algorithm Supporting Different Page Sizes", *Lecture Notes in Computer Science* 303 (April 1988), 336, Springer Verlag.
- [252] A. Silberschatz, M. Stonebraker and J. Ullman, "Database Systems: Achievements and Opportunities", *Communications of the ACM, Special Section on Next-Generation Database Systems* 34, 10 (October 1991), 110.
- [253] H. W. Six and P. Widmayer, "Spatial Searching in Geometric Databases", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1988, 496.
- [254] J. M. Smith and P. Y. T. Chang, "Optimizing the Performance of a Relational Algebra Database Interface", *Communications of the ACM* 18, 10 (October 1975), 568.
- [255] R. Snodgrass and K. Shannon, "Semantic Clustering", *Fourth Int'l Workshop on Persistent Object Systems*, Martha's Vineyard, MA, September 1990, 361.
- [256] R. Snodgrass, "Temporal Databases: Status and Research Directions", *ACM SIGMOD Record, Special Issue on Directions for Future Database Research and Development* 19, 4 (December 1990), 83.
- [257] J. A. Solworth and C. U. Orji, "Write-Only Disk Caches", *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 123.
- [258] V. Srinivasan and M. J. Carey, "Performance of B-Tree Concurrency Control Algorithms", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 416.
- [259] V. Srinivasan and M. Carey, "On-Line Index Construction Algorithms", *Computer Sciences Technical Report 1008*, March 1991.
- [260] V. Srinivasan and M. J. Carey, "Performance of On-Line Index Construction Algorithms", *Proc. Int'l. Conf. on Extending Database Technology*, Vienna, Austria, March 1992.
- [261] J. W. Stamos and H. C. Young, "A Symmetric Fragment and Replicate Algorithm for Distributed Joins", *Technical Report RJ7188* (December 5, 1989), IBM Almaden Research Lab.
- [262] M. Stonebraker, "Implementation of Integrity Constraints and Views by Query Modification", *Proc. ACM SIGMOD Conf.*, San Jose, CA, June 1975.
- [263] M. Stonebraker, "Operating system support for database management", *Communications of the ACM* 24, 7 (July 1981).
- [264] M. Stonebraker, "The Design and Implementation of Distributed INGRES", in *The INGRES Papers*, M. Stonebraker (editor), Addison-Wesley, Reading, MA, 1986, 187.
- [265] M. Stonebraker, "The Case for Shared-Nothing", *IEEE Database Eng.* 9, 1 (March 1986).
- [266] M. Stonebraker, "The Design of the POSTGRES Storage System", *Proc. Int'l. Conf. on Very Large Data Bases*, Brighton, England, August 1987, 289. Reprinted in M. Stonebraker, *Readings in Database Systems*, Morgan-Kaufman, San Mateo, CA, 1988.
- [267] M. Stonebraker, R. Katz, D. Patterson and J. Ousterhout, "The Design of XPRS", *Proc. Int'l. Conf. on Very Large Data Bases*, Long Beach, CA, August 1988, 318.

- [268] M. Stonebraker, P. Aoki and M. Seltzer, "Parallelism in XPRS", *UCB/Electronics Research Lab. Memorandum M89/16*, Berkeley, February 1989.
- [269] M. Stonebraker and G. A. Schloss, "Distributed RAID — A New Multiple Copy Algorithm", *Proc. IEEE Conf. on Data Eng.*, Los Angeles, CA, February 1990, 430.
- [270] M. Stonebraker, L. A. Rowe and M. Hirohama, "The Implementation of Postgres", *IEEE Trans. on Knowledge and Data Eng.* 2, 1 (March 1990), 125.
- [271] M. Stonebraker, A. Jhingran, J. Goh and S. Potamianos, "On Rules, Procedures, Caching and Views in Data Base Systems", *Proc. ACM SIGMOD Conf.*, Atlantic City, NJ, May 1990, 281.
- [272] M. Stonebraker, "Managing Persistent Objects in a Multi-level Store", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 2.
- [273] D. D. Straube and M. T. Ozsü, "Query transformation rules for an object algebra", Univ. of Alberta, Dept. of Computing Sciences Tech. Rep. 89-23, August 1989.
- [274] S. Y. W. Su, *Database Computers: Principles, Architectures and Techniques*, McGraw-Hill, New York, NY, 1988.
- [275] S. Su, "An Association Algebra for Processing Object-Oriented Databases", *Proc. IEEE Conf. on Data Eng.*, Kobe, Japan, April 1991.
- [276] Teradata, *DBC/1012 Data Base Computer, Concepts and Facilities*, Teradata Corporation, Los Angeles, CA, 1983.
- [277] F. W. Tompa and J. A. Blakeley, "Maintaining materialized views without accessing base data", *Information Systems* 13, 4 (1988), 393.
- [278] I. L. Traiger, "Virtual Memory Management for Data Base Systems", *ACM Operating Systems Review* 16, 4 (October 1982), 26.
- [279] I. L. Traiger, J. Gray, C. A. Galtieri and B. G. Lindsay, "Transactions and Consistency in Distributed Database Systems", *ACM Trans. on Database Systems* 7, 3 (September 1982), 323.
- [280] J. W. Tukey, *Exploratory Data Analysis*, Addison-Wesley, Reading, MA, 1977.
- [281] A. Unnikrishnan, P. Shankar and Y. V. Venkatesh, "Threaded Linear Hierarchical Quadrees for Computation of Geometric Properties of Binary Images", *IEEE Trans. on Software Eng.* 14, 5 (May 1988), 659.
- [282] P. Valduriez, "Join Indices", *ACM Trans. on Database Systems* 12, 2 (June 1987), 218.
- [283] S. L. Vandenberg and D. J. DeWitt, "Algebraic Support for Complex Objects with Arrays, Identity, and Inheritance", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 158.
- [284] C. B. Walton, "Investigating Skew and Scalability in Parallel Joins", *Department of Computer Sciences Technical Report Tech. Rep.-89-39* (December 1989), University of Texas.
- [285] C. Walton, A. Dale and R. Jenevein, "A Taxonomy and Performance Model of Data Skew in Parallel Joins", *Proc. Int'l. Conf. on Very Large Data Bases*, Barcelona, Spain, 1991.
- [286] G. Weikum, P. Zabback and P. Scheuermann, "Dynamic File Allocation in Disk Arrays", *Proc. ACM SIGMOD Conf.*, Denver, CO, May 1991, 406.
- [287] P. Williams, D. Daniels, L. Haas, G. Lapis, B. Lindsay, P. Ng, R. Obermarck, P. Selinger, A. Walker, P. Wilms and R. Yost, "R*: An Overview of the Architecture", in *Readings in Database Systems*, M. Stonebraker (editor), Morgan-Kaufman, San Mateo, CA, 1988.
- [288] J. L. Wolf, D. M. Dias and P. S. Yu, "An Effective Algorithm for Parallelizing Sort Merge in the Presence of Data Skew", *Proc. Int'l. Symp. on Databases in Parallel and Distributed Systems*, Dublin, Ireland, July 1990, 103.
- [289] J. Wolf, "An Effective Algorithm for Parallelizing Hash Joins in the Presence of Data Skew", *Proc. IEEE Conf. on Data Eng.*, Kobe, Japan, April 1991.
- [290] R. Wolniewicz and G. Graefe, "Automatic Optimization and Parallelization in Scientific Databases", in *preparation*, 1992.
- [291] E. Wong and K. Youssefi, "Decomposition - A Strategy for Query Processing", *ACM Trans. on Database Systems* 1, 3 (September 1976), 223.
- [292] E. Wong and R. H. Katz, "Distributing a Database for Parallelism", *Proc. ACM SIGMOD Conf.*, San Jose, CA, May 1983, 23.
- [293] H. Yang and P. A. Larson, "Query Transformation for PSJ-queries", *Proc. Int'l. Conf. on Very Large Data Bases*, Brighton, England, August 1987, 245.
- [294] K. Youssefi and E. Wong, "Query Processing in a Relational Database Management System", *Proc. Int'l. Conf. on Very Large Data Bases*, Rio de Janeiro, October 1979, 409.
- [295] L. Yu, "An Evaluation Framework for Algebraic Object-Oriented Query Models", *Proc. IEEE Conf. on Data Eng.*, Kobe, Japan, April 1991.

- [296] C. Zaniolo, "Design of Relational Views Over Network Schemas", *Proc. ACM SIGMOD Conf.*, Boston, MA, May-June 1979, 179.
- [297] H. Zeller, "Parallel Query Execution in NonStop SQL", *Digest of Papers, 35th CompCon Conf.*, San Francisco, CA, Feb-Mar 1990, 484.
- [298] H. Zeller and J. Gray, "An Adaptive Hash Join Algorithm for Multiuser Environments", *Proc. Int'l. Conf. on Very Large Data Bases*, Brisbane, Australia, 1990, 186.

Table of Contents

Abstract	1
Index Terms	1
Introduction	1
1. Architecture of Query Execution Engines	4
2. Sorting and Hashing	10
2.1. Sorting	10
2.2. Hashing	15
3. Disk Access	20
3.1. File Scans	20
3.2. Associative Access using Indices	20
3.3. Faster Storage Devices	23
3.4. Buffer Management	24
3.5. Physical Database Design	25
4. Aggregation and Duplicate Removal	26
4.1. Aggregation Algorithms Based on Sorting	27
4.2. Aggregation Algorithms Based on Hashing	28
4.3. A Rough Performance Comparison	29
4.4. Additional Remarks on Aggregation	30
5. Binary Matching Operations	31
5.1. Nested Loops Join Algorithms	32
5.2. Merge-Join Algorithms	33
5.3. Hash Join Algorithms	34
5.4. Pointer-Based Joins	36
5.5. A Rough Performance Comparison	37
6. Universal Quantification	38
7. Duality of Sorting and Hashing	42
8. Execution of Complex Query Plans	49
9. Mechanisms for Parallel Query Execution	54
9.1. Parallel vs. Distributed Database Systems	54
9.2. Forms of Parallelism	55
9.3. Implementation Strategies	56
9.4. Load Balancing and Skew	59
9.5. Tuning a Parallel System	60
9.6. Architectures and Architecture-Independence	62
10. Parallel Algorithms	64
10.1. Parallel Selections and Updates	65
10.2. Parallel Sorting	65
10.3. Parallel Aggregation and Duplicate Removal	68
10.4. Parallel Joins and Other Binary Matching Operations	68
10.5. Parallel Universal Quantification	71
11. Non-Standard Query Processing Algorithms	71
11.1. Nested Relations	72
11.2. Temporal and Scientific Database Management	73
11.3. Object-Oriented Database Systems	74
11.4. More Meta-Operators	75
12. Additional Techniques for Performance Improvement	77

12.1. Precomputation and Derived Data	77
12.2. Data Compression	78
12.3. Surrogate Processing	79
12.4. Bit Vector Filtering	80
12.5. Specialized Hardware	82
Summary and Outlook	82
Acknowledgements	83
References	83