

Heuristic and randomized optimization for the join ordering problem

Michael Steinbrunn^{1,*}, Guido Moerkotte^{2,**}, Alfons Kemper^{1,***}

¹ Universität Passau, Fakultät für Mathematik und Informatik, D-94030 Passau, Germany

² Universität Mannheim, Lehrstuhl für Praktische Informatik III, D-68131 Mannheim, Germany

Abstract. Recent developments in database technology, such as deductive database systems, have given rise to the demand for new, cost-effective optimization techniques for join expressions. In this paper many different algorithms that compute approximate solutions for optimizing join orders are studied since traditional dynamic programming techniques are not appropriate for complex problems. Two possible solution spaces, the space of left-deep and bushy processing trees, are evaluated from a statistical point of view. The result is that the common limitation to left-deep processing trees is only advisable for certain join graph types. Basically, optimizers from three classes are analysed: heuristic, randomized and genetic algorithms. Each one is extensively scrutinized with respect to its working principle and its fitness for the desired application. It turns out that randomized and genetic algorithms are well suited for optimizing join expressions. They generate solutions of high quality within a reasonable running time. The benefits of heuristic optimizers, namely the short running time, are often outweighed by merely moderate optimization performance.

Key words: Query optimization – Join ordering – Heuristic algorithms – Randomized algorithms – Genetic algorithms

1 Introduction

In recent years, relational database systems have become the standard in a variety of commercial and scientific applications. Because queries are stated in a non-procedural manner, the need for optimizers arises that transform the straightforward translation of a query into a cost-effective evaluation plan. Due to their high evaluation costs, joins are a primary target of query optimizers. If queries are stated interactively, there are generally only few relations involved. The optimization of these expressions can be carried out

by exhaustive search, possibly enhanced by pruning techniques that exclude unlikely candidates for good solutions. For instance, in System R [SAC⁺79], a dynamic programming algorithm is employed for the optimization of joins. This approach works well as long as only few relations are to be joined, but if the join expression consists of more than about five or six relations, dynamic programming techniques quickly become prohibitively expensive. Queries of this kind are encountered in recent developments such as deductive database systems, where join expressions may consist of a large number of relations. Another source for such queries are query-generating database system front ends and complex views. In both cases, very complex queries may be issued without the end user being aware of it. Even in object-oriented database systems [KM94], complex join expressions may be encountered; while forward traversal of object references are usually very well supported by specialized access mechanisms and would not be treated as ordinary join operations, this is not true for *backward* traversal. This would require appropriate index structures such as Access Support Relations [KM92], the processing of which, in turn, involves the handling of potentially very complex join expressions for both initial materialization and maintenance.

Hence, there is a demand for optimization techniques that can cope with such complex queries in a cost-effective manner. In this paper, we shall examine approaches for the solution of this problem and assess their advantages and disadvantages. The rest of the article is organized as follows: in Sect. 2 we give an exact definition of the problem and the terms and present several cost models that we shall be using later on in our analysis. Section 3 deals with the problem of different solution spaces for evaluation strategies. In Sect. 4 we describe common optimization strategies with the varying working principle, which are subject to a quantitative analysis in Sect. 5. Section 6 concludes the paper.

2 Problem description

The problem of determining good evaluation strategies for join expressions has been addressed from the development of the first relational database systems [WY76; YW79; SAC⁺79]. The work in this area can be divided into two

* e-mail: steinbrunn@db.fmi.uni-passau.de

** e-mail: moer@pi3.informatik.uni-mannheim.de

*** e-mail: kemper@db.fmi.uni-passau.de

Correspondence to: M. Steinbrunn

major streams: First, the development of efficient algorithms for performing the join itself, and second, algorithms that determine the nesting order in which the joins are to be performed. In this article, we shall be concentrating on the generation of low-cost join nesting orders while disregarding the specifics of join computing – Mishra and Eich [ME92] provide a good overview on this subject.

In relational database systems where queries are stated interactively, join expressions that involve more than about five or six relations are rarely encountered. Therefore, the computation of an optimal join order with lowest evaluation cost by exhaustive search is perfectly feasible – it takes but a few seconds of CPU time. But if more than about eight relations are to be joined, the generally NP-hard problem of determining the optimal order [IK84] can no longer be solved precisely; we have to rely on algorithms that compute (hopefully) good approximate solutions. Those algorithms fall into two classes: first, augmentation heuristics that build an evaluation plan step by step according to certain criteria, and second, randomized algorithms that perform some kind of “random walk” through the space of all possible solutions seeking a solution with minimal evaluation cost.

2.1 Definition of terms

The input of the optimization problem is given as the *query graph* (or *join graph*), consisting of all relations that are to be joined as its nodes and all joins specified as its edges. The edges are labelled with the *join predicate* and the *join selectivity*. The join predicate maps tuples from the cartesian product of the adjacent nodes to {false, true}, depending on whether the tuple is to be included in the result or not. The join selectivity is the ratio “number of tuples in the result/number of tuples in the cartesian product”. As a special case, the cartesian product can be considered a join operation with join predicate \equiv true and a join selectivity of 1.

The *search space* (or *solution space*) is the set of all evaluation plans that compute the same result. A *point* in the solution space is one particular plan, i.e. a solution for the problem. A solution is described by the *processing tree* for evaluating the join expression. Every point of the solution space has a *cost* associated with it; a *cost function* maps processing trees to their respective costs. The processing tree itself is a binary tree that consists of base relations as its leaves and join operations as its inner nodes; edges denote the flow of data that takes place from the leaves of the tree to the root.

The goal of the optimization is to find the point in the solution space with lowest possible cost (global minimum). As the combinatorial explosion makes exhaustive enumeration of all possible solutions infeasible and the NP-hard characteristic of the problem implies that a faster algorithm (presumably) cannot exist, we have to rely on heuristics that compute approximate results.

2.2 Cost models

Our investigations are based on the cost models discussed in this subsection. Each of these cost models measures cost

as the number of pages that have to be read from or written to secondary memory. The execution environment is not distributed. The database is assumed to be much larger than the available main memory, so all costs besides I/O can be neglected without introducing too large an error. All cost models are based on parameters listed in Table 1. The join operations themselves are equijoins. A common term for each of the cost formulae below is the cost for writing the result of the join operation to the secondary memory. This cost is

$$C_{write}(R_1 \bowtie R_2) = \frac{\sigma_{12} \cdot |R_1| \cdot |R_2|}{bs/ts_{R_{12}}}$$

2.2.1 Nested loop join

The cost for performing a nested loop join (depending on the presence of index structures) is [EN94]:

1. Without index support

$$C_{nl}(R_1 \bowtie R_2) = \underbrace{b_{R_1}}_{\text{read } R_1} + \underbrace{\left(\left\lceil \frac{b_{R_1}}{ms} \right\rceil \cdot b_{R_2} \right)}_{\text{read } R_2 \text{ and perform join}}$$

2. Primary B⁺-tree index on the join attribute of R_2

$$C_{nl}(R_1 \bowtie R_2) = \underbrace{b_{R_1}}_{\text{read } R_1} + \underbrace{|R_1| \cdot (x_{R_2} + 1)}_{\text{use index to find matching tuple in } R_2}$$

3. Secondary B⁺-tree index on the join attribute of R_2

$$C_{nl}(R_1 \bowtie R_2) = \underbrace{b_{R_1}}_{\text{read } R_1} + \underbrace{|R_1| \cdot (x_{R_2} + s_{R_2})}_{\text{use index to find matching tuple in } R_2}$$

4. Hash index on the join attribute of R_2

$$C_{nl}(R_1 \bowtie R_2) = \underbrace{b_{R_1}}_{\text{read } R_1} + \underbrace{|R_1| \cdot h}_{\text{use index to find matching tuple in } R_2}$$

h is the average number of page accesses necessary to retrieve a tuple from R_2 with a given key. We use the value $h = 1.2$ for a primary hash index, and $h = 2.2$ for a secondary hash index.

2.2.2 Sort-merge join

The cost for performing a sort-merge join operation is [EN94]:

$$C_{sm}(R_1 \bowtie R_2) = C_{R_1} + C_{R_2}$$

where C_{R_1} and C_{R_2} is computed according to the following cases:

1. The relation is sorted on the join attribute (or there is a primary B⁺-tree index on the join attribute)

$$C_{R_x} = b_{R_x},$$

i.e. there is only the cost for reading the relation.

Table 1. Cost model parameters

Parameter	Meaning
$ R $	Cardinality (number of tuples) of relation R
ts_R	Tuple size of relation R (in bytes)
bs	Size of a disk block (in bytes)
ps	Size of a tuple reference (tuple identifier, TID)
ms	Main memory size (in bs units)
σ_{12}	Join selectivity for join $R_1 \bowtie R_2$ ($\sigma_{12} = \frac{ R_1 \bowtie R_2 }{ R_1 \times R_2 }$)
b_R	Number of blocks occupied by relation R
fo	Fanout of an internal B ⁺ -tree node ($fo = \lfloor 0.69 \cdot bs / (ks + pr) \rfloor$) (ks = key size, pr = size of a page reference)
x_R	Height of a B ⁺ -tree index on the join attribute of R minus one ($x_R = \lceil \log_{fo} b_R \rceil - 1$, assuming $ps \approx ks + pr$)
s_R	Selection cardinality of R 's join attribute (average number of tuples with the same value of the join attribute)

2. There is a secondary B⁺-tree index on the join attribute

$$C_{R_x} = \left\lceil |R_x| \cdot \frac{ps}{0.69 \cdot bs} \right\rceil + b_{R_x},$$

i.e. the leaf nodes of the index tree (assuming to be 69% full) have to be scanned for pointers to the tuples of the relation, and the blocks containing the tuples themselves must be read at least once.

3. There is no sort order on the join attribute; explicit sorting is required

$$C_{R_x} = b_{R_x} \log_{ms} b_{R_x} + b_{R_x}$$

We assume the merge-sort algorithm is applied, where the number of merge passes depends on the amount of main memory available.

2.2.3 Hash join

We assume that a ‘‘hybrid hash join’’ is carried out. This algorithm performs very well over a large range of available main memory. The cost is [Sha86]:

$$b_{R_1} + b_{R_2} + 2 \cdot (b_{R_1} + b_{R_2}) \cdot (1 - q)$$

where q denotes the fraction of R_1 whose hash table fits into main memory. It is computed as:

$$q = \frac{ms - \left\lceil \frac{1.4 \cdot b_{R_1} - ms}{ms - 1} \right\rceil}{b_{R_1}}$$

The constant 1.4 accounts for the hash table’s load factor of about 71%.

2.2.4 Cost of an entire processing tree

In order to estimate the cost for evaluating an entire processing tree, the cost for each node is computed recursively (bottom-up, right-to-left) as the sum of the cost for obtaining the two son nodes and the cost for joining them in order to get the final result.

If the outer relation (R_1 in the cost formulae) is not a base relation, and the join algorithm is a nested loop join, we assume that pipelining is possible, which saves the costs

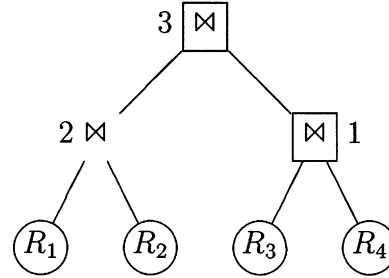


Fig. 1. Computation of processing tree costs

for writing an intermediate result to disk and to read it back into main memory.

For instance, the processing tree in Fig. 1 (where all join nodes are supposed to be nested loop joins) is evaluated as follows:

1. The join operation $R_3 \bowtie R_4$ (node 1) is performed. Because both operands are base relations, the cost for reading both R_3 and R_4 is included in the estimate. Furthermore, the result of this operation has to be written to disk as an intermediate relation.
2. The join operation $R_1 \bowtie R_2$ (node 2) is performed. Again, both operands are base relations, so the cost for scanning them has to be counted. But in contrast to node 1, no intermediate result has to be written to disk, because the tuples can be pipelined to node 3, the root of the processing tree.
3. In node 3, both intermediate results $R_1 \bowtie R_2$ and $R_3 \bowtie R_4$ are joined together in order to compute the final result. While $R_1 \bowtie R_2$ does not need to be read from disk due to the pipeline from node 2, $R_3 \bowtie R_4$ must be read back in, and the final result must be written to disk.

We note that the boxed join nodes’ results (Fig. 1) must be written to secondary memory. These considerations are valid if (and only if) the two processing nodes in question are both nested loop joins. If either node 2 or node 3 in the example tree in Fig. 1 were anything else but nested loop joins, the cost for writing the intermediate result to disk and reading it back into memory would have to be charged.

3 Solution space for the join ordering problem

Generally, the solution space is defined as the set of all processing trees that compute the result of the join expression and that contain each base relation exactly once. The leaves of the processing trees consist of the base relations, whereas the inner nodes correspond to join results of the appropriate sons. As the join operation is commutative and associative, the number of possible processing trees increases quickly with increasing number of relations involved in the join expression in question. Traditionally, a subset of the complete space, the set of so-called *left-deep processing trees*, has been of special interest to researchers [SAC⁺79; SG88; Swa89]. We shall now study the characteristics of both the complete solution space and the subset of left-deep trees as the most interesting special cases, although other tree shapes might be contemplated, e.g. right-deep trees or zig-zag trees, which are mainly of interest in distributed computing environments (see, for example [LVZ93]).

3.1 Left-deep trees

This subset consists of all processing trees where the inner relation of each join is a base relation. For a fixed number of base relations, the specification “left-deep” does not leave any degrees of freedom concerning the shape of the tree, but there are $n!$ ways to allocate n base relations to the tree’s leaves. It has been argued that good solutions are likely to exist among these trees, because such trees are capable of exploiting the cost-reducing pipelining technique on each of its join processing nodes. In case a processing tree consists solely of nested loop joins (either with or without index support), not a single intermediate result has to be materialized on secondary memory.

3.2 Bushy trees

In this solution space, we also permit join nodes where both operands are “composites” (i.e. no base relations). Thus, the solutions in this space are in no way restricted. Consequently, this solution space includes left-deep as well as other special tree shapes as (strict) subsets. Because the shape of possible processing trees can be arbitrary, the cardinality of this set is much higher than the cardinality of the left-deep space: for n base relations, there are $\binom{2(n-1)}{n-1}(n-1)!$ different solutions. However, although the degrees of freedom in constructing bushy trees are much higher, the capability of exploiting the pipelining technique is restricted to a subset of the tree’s join processing nodes. The more the shape of the tree tends toward right-deep (i.e. the join nodes’ *left* operands are base relations), the smaller is the size of this subset. For a right-deep tree, none of its join processing nodes is capable of pipelining.

In [OL90], an adaptable plan enumeration strategy for linear (chain) and star-shaped join graphs is proposed that reduces the number of plans whose costs have to be evaluated considerably. If n denotes the number of relations in the join graph, there are $(n^3 - n)/6$ (bushy tree solution space) or $(n - 1)^2$ (left-deep tree solution space) feasible joins for

linear graphs. For star graphs, there are $(n - 1) \cdot 2^{n-2}$ (bushy tree solution space) feasible joins. However, this approach requires a specially tailored “join plan enumerator” for every class of join graphs that might be encountered, and for arbitrary join graphs still the entire solution space must be considered in order to guarantee that the optimal solution cannot be missed.

4 Join ordering strategies

The problem of finding a good nesting order for n -relational joins can be tackled in several different ways:

1. *Deterministic algorithms.* Every algorithm in this class constructs a solution step by step in a deterministic manner, either by applying a heuristic or by exhaustive search.
2. *Randomized algorithms.* Algorithms in this class pursue a completely different approach: a set of *moves* is defined. These moves constitute *edges* between the different solutions of the solution space; two solutions are connected by an edge if (and only if) they can be transformed into one another by exactly one move. Each of the algorithms performs a random walk along the edges according to certain rules, terminating as soon as no more applicable moves exist or a time limit is exceeded. The best solution encountered so far is the result.
3. *Genetic algorithms.* Genetic algorithms make use of a randomized search strategy very similar to biological evolution in their search for good problem solutions. Although in this aspect genetic algorithms resemble randomized algorithms as discussed above, the approach shows enough differences to warrant a consideration of its own. The basic idea is to start with a random population and generate offspring by random crossover and mutation. The “fittest” members of the population (according to the cost function) survive the subsequent selection; the next generation is based on these. The algorithm terminates as soon as there is no further improvement or after a predetermined number of generations. The fittest member of the last population is the solution.
4. *Hybrid algorithms.* Hybrid algorithms combine the strategies of pure deterministic and pure randomized algorithms: solutions obtained by deterministic algorithms are used as starting points for randomized algorithms or as initial population members for genetic algorithms.

4.1 Deterministic algorithms

The algorithms discussed in this section either employ heuristics or a (pruned) search of the solution space in order to optimize the given join expression. We shall take a closer look at four different algorithms of this class with varying complexity and performance.

4.1.1 Dynamic programming

This is the classical algorithm that has been used for join order optimization in System-R [SAC⁺79]. It searches the

function DynProg**inputs** rels “List of relations to be joined”**outputs** pt “Processing Tree”

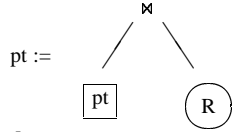
partialsolutions := {All scans for all attributes involved}

“Remove all elements from partialsolutions
with equivalent, lower-cost alternative”

```

for  $i := 2$  to |rels|
  for all pt in partialsolutions
    for all  $R$  in rels such that  $R$  not in pt

```

**end****end**“Remove all elements from partialsolutions
with equivalent, lower-cost alternative”**end****return** “Arbitrary element from partialsolutions”**Fig. 2.** Algorithm “Dynamic Programming”

solution space of left-deep processing trees. Firstly, the set of partial solutions is initialized with all possible scan nodes for all relation attributes that participate in the query. For instance, if there is an index on attribute $R.A$, then both the index scan and the ordinary file scan are considered feasible partial processing trees. In the next step, every element with a cheaper, equivalent alternative is pruned from the set of possible partial solutions, where an alternative is considered “equivalent” if it joins the same set of relations and the sort order of the partial result is the same. In the following loop, the algorithm constructs in the k th iteration a set of k -relation partial solutions from a set of $(k-1)$ -relation partial solutions. When this loop terminates, the set *partialsolutions* consists of at least one, possibly several equivalent, optimal solutions.

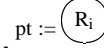
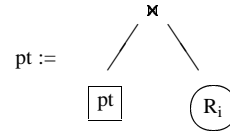
A pseudo code rendering of this algorithm is shown in Fig. 2. Apart from the removal of all equivalent alternatives but the cheapest one, the original algorithm according to the cited reference performs further pruning of the search tree: it defers the introduction of cartesian products into partial solutions as long as possible, thus removing unlikely candidates for the optimal solution. However, although this strategy reduces the computational complexity, the result is no longer guaranteed to be optimal.

A major disadvantage of this algorithm is the high memory consumption for storing partial solutions. That (and the exponential running time) makes its application for queries that involve more than about 10–15 relations prohibitively expensive.

In very recent work, Vance and Maier [VM96] devised a very efficient, so-called light-weight implementation of dynamic programming for bushy tree join optimization. Their method allows to optimize join queries with up to about 18 relations – albeit with a rather simplified cost model.

function MinSel**inputs** rels “List of relations to be joined”**outputs** pt “Processing Tree”

pt := NIL

do**if** pt = NIL **then** $R_i :=$ “Relation with smallest cardinality”**else** $R_i :=$ “Relation from rels with smallest selectivity factor
for the join with pt”**end**rels := rels \ [R_i]**while** rels $\neq []$ **return** pt;**Fig. 3.** Minimum selectivity algorithm

4.1.2 Minimum selectivity

Good solutions are generally characterized by intermediate results with small cardinality. The *minimum selectivity heuristic* builds a left-deep processing tree step by step while trying to keep intermediate relations as small as possible. In this regard, this resembles Ingres’ decomposition strategy [WY76]; however, unlike the decomposition strategy, which considers only the operands’ cardinalities, the minimum selectivity heuristic makes use of the *selectivity factor* σ of the join $R_1 \bowtie R_2$ to achieve small intermediate results. Firstly, the set of relations to be joined is divided into two subsets: the set of relations already incorporated into the intermediate result, denoted \mathcal{R}_{used} (which is initially empty), and the set of relations still to be joined with the intermediate result, denoted $\mathcal{R}_{remaining}$ (which initially consists of the set of all relations). Then, in each step of the algorithm, the relation $R_i \in \mathcal{R}_{remaining}$ with the lowest selectivity factor

$$\sigma_i := \frac{\left| R_i \bowtie \left(\bigbowtie_{R_u \in \mathcal{R}_{used}} R_u \right) \right|}{|R_i| \cdot \left| \bigbowtie_{R_u \in \mathcal{R}_{used}} R_u \right|}$$

is joined with the (so far) intermediate result and moved from $\mathcal{R}_{remaining}$ to \mathcal{R}_{used} . Figure 3 shows the complete algorithm for left-deep processing trees.

4.1.3 Krishnamurthy-Boral-Zaniolo algorithm

On the foundation of [Law78] and [MS79], Ibaraki and Kameda [IK84] showed that it is possible to compute the optimal nesting order in polynomial time, provided the query

graph forms a tree (i.e. no cycles) and the cost function is a member of a certain class. Based on this result, Krishnamurthy et al. [KBZ86] developed an algorithm (KBZ algorithm) that computes the optimal solution for a tree query in $O(n^2)$ time, where n is the number of joins.

In the first step, every relation plays, in turn, the role of the root of the query tree. For all roots, the tree is linearized by means of a *ranking function* that establishes the optimal evaluation order for that particular root. The linearized tree obeys the tree's order, in other words, a parent node is always placed before the son nodes. The evaluation order with lowest cost is the result of the algorithm.

By transforming the query tree into a rooted tree, a parent node for every node can be uniquely identified. Thus, the selectivity of a join, basically an edge attribute of the query graph, can be assigned to the nodes as well. If the cost function C can be expressed as $C(R_i \bowtie R_j) = |R_i| \cdot g(|R_j|)$, where g is an arbitrary function, the join cost can be assigned to a particular node, too. This is, in principle, possible for nested loop join algorithms, but not for merge join or hash join algorithms. The cost can be computed recursively as follows (Λ denotes the empty sequence, and l_1 and l_2 partial sequences):

$$C(\Lambda) = 0$$

$$C(R_i) = \begin{cases} |R_i| & \text{if } R_i \text{ is the root node} \\ g(|R_i|) & \text{else} \end{cases}$$

$$C(l_1 l_2) = C(l_1) + T(l_1)C(l_2)$$

The auxiliary function $T(l)$ is defined as:

$$T(l) = \begin{cases} 1 & \text{if } l = \Lambda \text{ (empty sequence)} \\ \prod_{R_k \in l} \sigma_k |R_k| & \text{else} \end{cases}$$

σ_k denotes the selectivity of the join of R_k with its parent node.

The algorithm is based on the so-called "adjacent sequence interchange property" [IK84] for cost functions that can be expressed as $C(R_i \bowtie R_j) = |R_i| \cdot g(|R_j|)$. If the join graph J is a rooted tree and A, B, U and V are sequences of J 's nodes (U and V non-null), such that the partial order defined by J is not violated by $AVUB$ and $AUVB$, then

$$C(AVUB) \leq C(AUVB) \Leftrightarrow \text{rank}(V) \leq \text{rank}(U)$$

The *rank* of a non-null sequence S is defined as

$$\text{rank}(S) := \frac{T(S) - 1}{C(S)}$$

Thus, the cost can be minimized by sorting according to the ranking function $\text{rank}(S)$, provided the partial order defined by the tree is preserved.

The algorithm for computing the minimum cost processing tree consists of the auxiliary function *linearize* and the main function *KBZ*. Firstly, the join tree is linearized according to the function *linearize* in Fig. 4, where a bottom-up merging of sequences according to the ranking function is performed. In the last step, the root node becomes the head of the sequence thus derived. However, it is possible that the root node has a higher rank than its sons, therefore a normalization of the sequence has to be carried out. That means that the first relation in the sequence (the root node) is joined

function Linearize

inputs root "Root of a (partial) tree"

outputs chain "Optimal join order for the tree-shaped join graph with root 'root'"

```

chain := []
for all succ in Sons(root)
  lin := Linearize(succ)
  "Merge lin into chain according to ranks"
end

chain := root + chain
"Normalize the root node 'root' (cf. text)"

return chain;

```

Fig. 4. Auxiliary function *linearize*

function KBZ

inputs joingraph

outputs minorder "join order"

```

tree := "Minimum spanning tree of joingraph"
mincost := ∞

forall node in tree
  lin := Linearize(node)
  "Undo normalization"

  cost := Cost(lin)
  if cost < mincost then
    minorder := lin
    mincost := cost
  end
end

return minorder;

```

Fig. 5. KBZ algorithm

with its successor. If necessary, this step has to be repeated until the order of the sequence is correct. The cost of the sequence is computed with the recursive cost function C .

In the main function *KBZ* (Fig. 5), this procedure is carried out for each relation of the join graph acting as the root node. The sequence with lowest total cost is the result of the optimization.

The algorithm can be extended to general (cyclic) join graphs in a straightforward way, namely by reducing the query graph to its minimal spanning tree using Kruskal's algorithm [Kru56]. The weight of the join graph's edges is determined by the selectivity of the appropriate join, and the minimal spanning tree is determined as the tree with the lowest product of edge weights, rather than the sum of the edges' weights, common in other applications of Kruskal's algorithm. This extension has been suggested by Krishnamurthy et al. [KBZ86]. However, if the join graph is cyclic, the result is no longer guaranteed to be optimal – it is but a heuristic approximation. When we speak of the KBZ algorithm in later sections, we refer to this extension with the computation of the minimal spanning tree of the join graph.

Due to its working principle, the KBZ algorithm requires the assignment of join algorithms to join graph edges *before* the optimization is carried out. This requirement and the restrictions concerning the cost model are the main drawbacks

function AB**inputs** joingraph**outputs** minororder “join order”

```

while number of iterations  $\leq N^2$  do
  begin
    randomize_methods;
    while number of iterations  $\leq N^2$  do
      begin
        apply_KBZ;
        change_order;
        change_methods;
      end;
    end;
    post_process;
  end
return minororder;

```

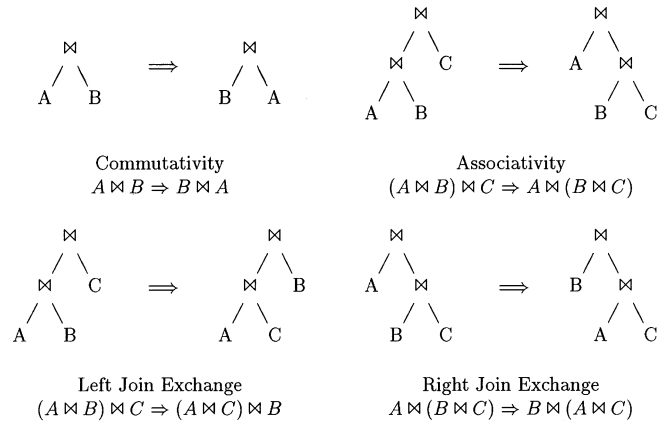
Fig. 6. AB algorithm

of the KBZ algorithm. The more sophisticated and detailed the cost model is, the more likely it is that KBZ’s optimal result based on an (almost inevitably crude) approximation is different from the real optimum. Furthermore, separating the two tasks of join order optimization and join method assignment invalidates the main advantage of the KBZ algorithm, namely to yield the optimal solution in $O(n^2)$ time. In the following section, an algorithm is discussed that tries to remedy this situation.

4.1.4 AB algorithm

The AB algorithm has been developed by Swami and Iyer [SI93]. It is based on the KBZ algorithm with various enhancements trying to remove the restrictions that are imposed on the join method placement. The algorithm permits the use of two different join methods, namely nested loop and sort-merge. The sort-merge cost model has been simplified by Swami and Iyer such that it conforms to the requirements of the KBZ algorithm [$C(R_1 \bowtie R_2) = |R_1| \cdot g(|R_2|)$] for some function g ; cf. Sect. 4.1.3]. The algorithm runs as follows (cf. Fig. 6):

1. In *randomize_methods*, each join in the join graph is assigned a randomly selected join method. If the join graph is cyclic, a random spanning tree is selected first.
2. The resulting tree query is optimized by the KBZ algorithm (*apply_KBZ*).
3. *change_order* attempts to further reduce the cost by swapping relations such that “interesting orders” can be exploited.
4. This step comprises a single scan through the join order achieved so far. For each join, an attempt is made to reduce the total cost by changing the join method employed (*change_method*).
5. Steps 2 to 4 are iterated until no further improvement is possible or N^2 iterations are performed (N = number of joins in the join graph).
6. Steps 1 to 5 are repeated as long as the total number of iterations of the inner loop does not exceed N^2 .

**Fig. 7.** Moves for bushy tree solution space traversal

7. In a post-processing step (*post_process*), once more the order of the relations is changed in an attempt to reduce the cost.

The AB algorithm comprises elements of heuristic and randomized optimizers. The inner loop searches heuristically for a local minimum, whereas in the outer loop several random starting points are generated in the manner of the iterative improvement algorithm (cf. Sect. 4.2.1). However, without ignoring the contribution of the KBZ algorithm, even with the AB extension it is hardly possible to make use of a sophisticated cost model.

4.2 Randomized algorithms

Randomized algorithms view solutions as *points* in a solution space and connect these points by edges that are defined by a set of *moves*. The algorithms discussed below perform some kind of random walk through the solution space along the edges defined by the moves. The kind of moves that are considered depend on the solution space: if left-deep processing trees are desired, each solution can be represented uniquely by an ordered list of relations participating in the join. Two different moves are proposed in [SG88] and [Swa89] for modifying these solutions: Swap and 3Cycle. Swap exchanges the positions of two arbitrary relations in the list, while 3Cycle performs a cyclic rotation of three arbitrary relations in the list. For instance, if $R_1R_2R_3R_4R_5$ was a point in the solution space, the application of Swap might lead to $R_1R_4R_3R_2R_5$, whereas 3Cycle could yield $R_5R_2R_1R_4R_3$.

If the complete solution space with arbitrarily shaped (bushy) processing trees is considered, the moves depicted in Fig. 7 (introduced by Ioannidis and Kang [IK90]) are used for traversal of the solution space.

4.2.1 Iterative improvement

If the solution space of the join optimization problem did contain but one global cost minimum without any local minima, we could use a simple hill-climbing algorithm for finding this minimum. However, because the solution space

function IterativeImprovement**outputs** minstate “Optimized processing tree”

```

mincost := ∞
do
  state := “Random starting point”
  cost := Cost(state)

  do
    newstate := “state after random move”
    newcost := Cost(newstate)
    if newcost < cost then
      state := newstate
      cost := newcost
    end
  while “Local minimum not reached”

  if cost < mincost then
    minstate := state
    mincost := cost
  end

  while “Time limit not exceeded”

```

return minstate;

Fig. 8. Iterative improvement algorithm

does contain local minima, hill-climbing would almost certainly yield one of them. The *iterative improvement algorithm* [SG88; Swa89; IK90] tries to overcome this problem in the following way (Fig. 8): after selecting a random starting point, the algorithm seeks a minimum cost point using a strategy similar to hill-climbing. Beginning at the starting point, a random neighbour (i.e. a point that can be reached by exactly one move) is selected. If the cost associated with the neighbouring point is lower than the cost of the current point, the move is carried out and a new neighbour with the lower cost is sought. This strategy is insofar different from genuine hill-climbing, as no attempt is made to determine the neighbour with the lowest cost. The reason for this behaviour is the generally very high number of neighbours that would have to be checked. The same holds for the check of whether a given point is a local minimum or not. Instead of systematically enumerating all possible neighbours and checking each one individually, a point is assumed to be a local minimum if no lower-cost neighbour can be found in a certain number of tries.

This procedure is repeated until a predetermined number of starting points is processed or a time limit is exceeded. The lowest local minimum encountered is the result.

4.2.2 Simulated annealing

Iterative Improvement suffers from a major drawback: because moves are accepted only if they improve the result obtained so far, it is possible that even with a high number of starting points the final result is still unacceptable. This is the case especially when the solution space contains a large number of high-cost local minima. In this case, the algorithm gets easily “trapped” in one of the high-cost local minima.

Simulated annealing (Fig. 9) is a variant on iterative improvement that removes this restriction [IW87; SG88]. In

function SimulatedAnnealing**inputs** state “Random starting point”**outputs** minstate “Optimized processing tree”

```

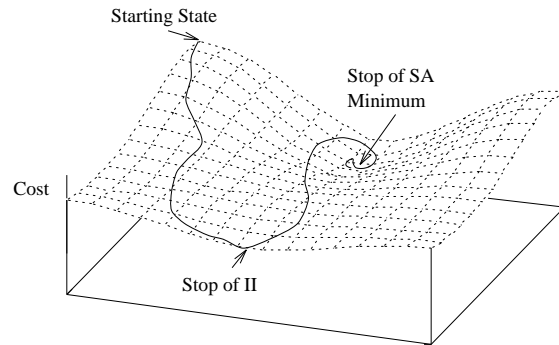
minstate := state; cost := Cost(state); mincost := cost
temp := “Starting temperature”
do
  do
    newstate := “state after random move”
    newcost := Cost(newstate)
    if newcost ≤ cost then
      state := newstate
      cost := newcost
    else “With probability  $e^{\frac{\text{newcost} - \text{cost}}{\text{temp}}}$ ”
      state := newstate
      cost := newcost
    end
  end

  if cost < mincost then
    minstate := state
    mincost := cost
  end
  end

  while “Equilibrium not reached”
    “Reduce Temperature”
  while “Not frozen”

```

return minstate;

Fig. 9. Simulated annealing algorithm**Fig. 10.** Iterative improvement versus simulated annealing

simulated annealing, a move may be carried out even if the neighbouring point is of higher cost. Therefore, the algorithm does not get trapped in local minima as easily as iterative improvement. As the name of simulated annealing suggests, the algorithm tries to simulate the annealing process of crystals. In this natural process, the system eventually reaches a state of minimum energy. The slower the temperature reduction is carried out, the lower the energy of the final state (one large crystal is of lower energy than several smaller ones combined). Figure 10 illustrates this behaviour: one iteration of iterative improvement stops in the first local minimum, whereas simulated annealing overcomes the high-cost barrier that separates it from the global minimum, because the simulated annealing algorithm always accepts moves that lead to a lower cost state, but also accepts moves that increase costs with a probability that depends on the temperature and the difference between the actual and the new state’s cost.

Of course, the exact behaviour is determined by parameters like starting temperature, temperature reduction and

stopping condition. Several variants have been proposed in the literature – we shall present the detailed parameters in the next section where we analyse and compare those simulated annealing variants.

4.2.3 Two-phase optimization

The basic idea for this variant is the combination of iterative improvement and simulated annealing in order to combine the advantages of both [IK90]. Iterative improvement, if applied repeatedly, is capable of covering a large part of the solution space and descends rapidly into a local minimum, whereas simulated annealing is very well suited for thoroughly covering the neighbourhood of a given point in the solution space. Thus, two-phase optimization works as follows:

1. For a number of randomly selected starting points, local minima are sought by way of iterative improvement (Fig. 8).
2. From the lowest of these local minima, the simulated annealing algorithm (Fig. 9) is started in order to search the neighbourhood for better solutions.

Because only the close proximity of the local minimum needs to be covered, the initial temperature for the simulated annealing pass is set lower than it would be for the simulated annealing algorithm run by itself.

4.2.4 Toured simulated annealing

An approach similar to two-phase optimization has been proposed by Lanzelotte et al. [LVZ93] in the context of a distributed computing environment. In toured simulated annealing, several simulated annealing “tours” with different starting points are performed. Each starting point is derived from a deterministic algorithm that greedily builds processing trees using some augmentation heuristic. For instance, the minimum selectivity heuristic (discussed in Sect. 4.1.2) could be used to provide these starting points.

Similarly to two-phase optimization, the main benefit of toured simulated annealing is the reduced running time. The starting temperature for the different tours is set much lower (0.1 times the initial plan’s cost) than for simulated annealing with a random starting point, so the annealing process does not spend much time accepting moves that do not improve the current solution.

4.2.5 Random sampling

In [GLPK94], a radically different idea is pursued. All randomized algorithms discussed so far are based on transformations that attempt to reduce a given solution’s evaluation cost according to a set of rules until no further improvement can be achieved. However, an analysis of the cost distribution in the solution space reveals that a significant fraction of solutions is rather close to the optimum. An algorithm that draws a truly random sample of solutions should therefore contain the same fraction of good solutions as the entire space; however, designing such an algorithm that selects

each processing tree with equal probability is not trivial. In the above-mentioned work, such an algorithm (designed for acyclic join graphs) is presented; its application is most appropriate when a reasonably good (evaluation cost of less than two times the minimum cost) evaluation plan has to be identified quickly, as the experimental results in [GLPK94] indicate.

4.3 Genetic algorithms

Genetic algorithms are designed to simulate the natural evolution process. As in nature, where the fittest members of a population are most likely to survive and propagate their features to their offspring, genetic algorithms propagate solutions for a given problem from generation to generation, combining them to achieve further improvement. We provide a brief overview of the terminology and the working principles of genetic algorithms. For a comprehensive introduction, the reader is referred to [Gol89].

4.3.1 Terminology

Because genetic algorithms are designed to simulate biological evolution, much of the terminology used to describe them is borrowed from biology. One of the most important characteristics of genetic algorithms is that they do not work on a single solution, but on a set of solutions, the *population*. A single solution is sometimes called a *phenotype*. Solutions are always represented as *strings* (*chromosomes*), composed of *characters* (*genes*) that can take one of several different *values* (*alleles*). The *locus* of a gene corresponds to the *position* of a character in a string. Each problem that is to be solved by genetic algorithms must have its solutions represented as character strings by an appropriate *encoding*.

The “fitness” of a solution is measured according to an objective function that has to be maximized or minimized. Generally, in a well-designed genetic algorithm, both the average fitness and the fitness of the best solution increases with every new generation.

4.3.2 Basic algorithm

The working principle of the genetic algorithm that we use to optimize join expressions is the same as of the generic algorithm described below.

Firstly, a population of random character strings is generated. This is the “zero” generation of solutions. Then, each next generation is determined as follows:

1. A certain fraction of the fittest members of the population is propagated into the next generation (*selection*).
2. A certain fraction of the fittest members of the population is *combined* yielding offspring (*crossover*).
3. A certain fraction of the population (not necessarily the fittest) is altered randomly (*mutation*).

This loop is iterated until the best solution in the population has reached the desired quality, a certain, predetermined number of generations has been produced or no improvement

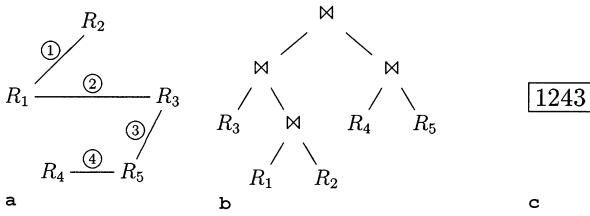


Fig. 11. Encoding of bushy processing trees: a join graph; b processing tree; c encoded tree

has been observed for a certain number of generations. In the next section, we shall examine how this generic algorithm can be adapted to the problem of optimizing join expressions.

4.3.3 Genetic algorithm for optimizing join expressions

Because genetic algorithms were not almost studied as intensively for join order optimization as other randomized algorithms, we shall discuss the questions associated with the employment of genetic algorithms for optimizing join expressions in more detail. In particular, we will not merely provide the techniques that we finally implemented, but some of the alternatives we considered (and tested) as well. Even if the basic algorithm remains unmodified, many variations for solution encoding, selection, crossover and mutation may be contemplated.

4.3.3.1 Encoding. Before a genetic algorithm can be applied to solve a problem, an appropriate encoding for the solution and an objective function has to be chosen. For join optimization, the solutions are processing trees, either left-deep or bushy, and the objective function is the evaluation cost of the processing tree that is to be minimized. For encoding processing trees, we considered two different schemes:

1. Ordered list

a) Left-deep trees. Solutions are represented as an ordered list of leaves. For instance, the processing tree $((((R_1 \bowtie R_4) \bowtie R_3) \bowtie R_2) \bowtie R_5)$ is encoded as “14325”.

b) Bushy trees. Bushy trees *without* cartesian products are encoded as an ordered list of join graph edges. This scheme has been proposed in [BFI91]. As an example of this encoding scheme, we represent the processing tree depicted in Fig. 11b as a character string. In a preliminary step, every edge of the join graph is labelled by an arbitrary number, such as in Fig. 11a. Then, the processing tree is encoded bottom-up and left-to-right, just the way as it would be evaluated. So, the first join of the tree joins relations R_1 and R_2 , i.e. edge 1 of the join graph. In the next steps, R_{12} and R_3 are joined, then R_4 and R_5 , and finally R_{123} and R_{45} , contributing edges 2, 4 and 3, respectively. Thus, the final encoding for our sample processing tree is “1243” (Fig. 11c).

2. Ordinal number encoding

a) Left-deep trees. A chromosome consists of a sequence of ordinal numbers of the processing tree’s

list of leaves. For instance, the processing tree $((((R_1 \bowtie R_4) \bowtie R_3) \bowtie R_2) \bowtie R_5)$ is encoded as follows:

- An ordered list L of all participating relations is made (for instance, based on their indices), such as $L = [R_1, R_2, R_3, R_4, R_5]$.
- The first relation in the processing tree, R_1 , is also the first relation in our list L , so its index “1” is the first gene of the chromosome. R_1 is then removed from the list L , so $L := [R_2, R_3, R_4, R_5]$.
- The second relation in the processing tree, R_4 , is the third relation in the list L , so “3” becomes the second gene of the chromosome. After removal of R_4 , L becomes $[R_2, R_3, R_5]$.
- This process is repeated until the list L is exhausted. In our example, the final encoding for the processing tree is “13211”.

b) Bushy trees. For bushy trees, the ordinal numbers in the chromosome denote join nodes similar to the ordered list of join edges described above. But instead of specifying the join node by the corresponding join graph edge, the join’s operands are used for that purpose. For instance, the processing tree in Fig. 11b is encoded as follows:

- An ordered list of all participating relations is made exactly as for left-deep tree encoding: $L := [R_1, R_2, R_3, R_4, R_5]$.
- The first join node in the processing tree is $R_1 \bowtie R_2$, which involves R_1 and R_2 with index “1” and “2”, respectively, so “12” becomes the first gene of the chromosome. R_1 and R_2 are replaced by R_{12} , so $L := [R_{12}, R_3, R_4, R_5]$.
- The next node in the processing tree joins relation R_3 with the result $R_1 \bowtie R_2$ (index 2 and 1), yielding gene “21” and $L := [R_{123}, R_4, R_5]$.
- Repeating this process finally leads to the complete chromosome “12 21 23 12”.

In the actual implementation, the chromosome’s genes carry additional information, namely operand order [encoding (1b)] and join algorithm (all encoding schemes).

4.3.3.2 Selection. The selection operator is used to separate good and bad solutions in the population. The motivation is to remove bad solutions and to increase the share of good solutions. Mimicking nature, selection is realized as shown in Fig. 12. The sample population consists of four solutions, and the objective function, cost, has to be minimized. The cost value for each of the solutions is listed in the table in Fig. 12. Each solution is assigned a sector of size inverse proportional to its cost value on a biased roulette wheel. Four spins of the wheel might yield the result in the second table, where Solution 4 has not been selected – it “became extinct due to lack of adaptation.”

This selection scheme is based on the fitness ratio of the members of the population: the better a member satisfies the objective function, the more it dominates the wheel, so one (relative) “super” population member may cause the premature convergence to a mediocre solution, because of the

Population			
Nr.	String	Cost	Fraction of the Wheel
1	31542	31784	31.1%
2	45321	46924	30.1%
3	51234	174937	21.2%
4	14325	227776	17.6%
Sum			100.0%

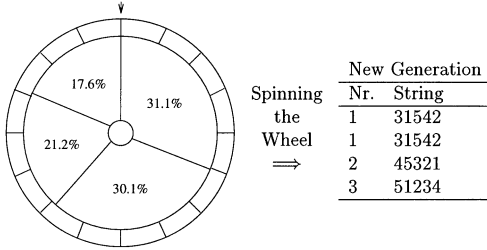


Fig. 12. Selection

disappearance of other members' features. Those features may be valuable, even if the solution as a whole is not of high quality. To avoid this, we use *ranking-based selection*. This means that it is not the value of the objective function itself but only its rank is used for biasing the selection wheel. In Fig. 12, for instance, it is not the cost values that would determine the fraction of the wheel a solution is assigned to, but just its rank value, i.e. 4 for solution 1, 3 for solution 2, 2 for solution 3 and 1 for solution 4.

General experience shows that ranking-based selection usually makes the evolution process advance more slowly, but the risk of untimely losing important information contained in weaker solutions is much lower.

Another variant is to keep the best solution in any case. This strategy (sometimes referred to as "elitist") helps speed up the convergence to the (near) optimal solution, because the risk of losing an already very good solution is eliminated.

4.3.3.3 Crossover. The crossover operator is a means of combining partially good solutions in order to obtain a superior result. The realization of a crossover operator depends heavily on the chosen encoding. For instance, the crossover operator has to make sure that the characteristics of the particular encoding are not violated. Such a characteristic is the uniqueness of each character in the string for the ordered list encoding scheme. The crossover operator and the encoding scheme are tightly coupled, because often the implementation of a particular crossover operator is facilitated (or even made possible at all) if a particular encoding scheme is used. Basically, we considered two different crossover operators, namely *subsequence exchange* and *subset exchange*. They work as follows:

1. Subsequence exchange (ordered list encoding). An example of the application of this operator is shown in Fig. 13. It assumes the ordered list encoding scheme. In each of the two offspring chromosomes, a random subsequence is permuted according to the genes' order of appearance in the other parent. For instance, in Fig. 13, the subsequence "532" is selected from the string "45321". The first gene of its offspring remains the same as in the parent (4). The second gene is taken from the first gene

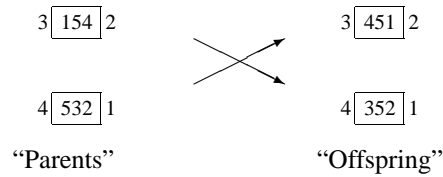


Fig. 13. Crossover 1 – subsequence exchange for ordered list encoding

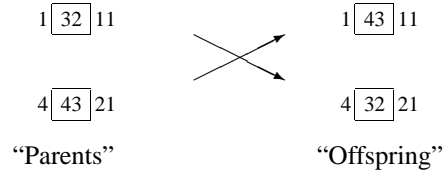


Fig. 14. Crossover 2 – subsequence exchange for ordinal number encoding

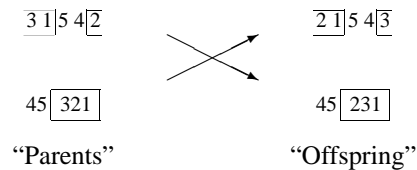


Fig. 15. Crossover 3 – subset exchange

of the other parent (3). The second gene of the other parent (1) cannot be used, because it is already present, so the third gene of the offspring is taken from the third gene of the other parent. Continuing this process yields at last the offspring chromosome "43521". Determining the second offspring is carried out similarly.

2. Subsequence exchange (ordinal number encoding). This operator is a slight variation of the above. It is intended for use in conjunction with the ordinal number encoding. In contrast to the first version of the sequence exchange operator, the two subsequences that are selected in the two parents must be of equal length. These subsequences are then simply swapped. This is only feasible with the ordinal number encoding because we do not have to worry about duplicated characters. Fig. 14 shows a sample application of this operator.
3. Subset exchange (ordered list encoding). The basic idea for this operator is to avoid any potential problems with duplicated characters by simply selecting two random subsequences with equal length in both parents that consist of the same set of characters. These two sequences are then simply swapped between the two parents in order to create two offspring. Figure 15 depicts an example of the application of this crossover operator.

4.3.3.4 Mutation. The mutation operator is needed for introducing features that are not present in any member of the population. Mutation is carried out by random alteration of a randomly selected chromosome. If the operator must not introduce duplicate characters, as in ordered list chromosomes, two random genes are simply swapped in order to carry out the mutation; with ordinal number encoding, a random gene of the chromosome is assigned a new, random value.

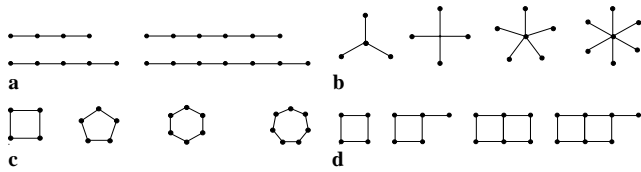


Fig. 16. Join graphs: a chain; b star; c cycle; d grid

Mutation, the “spice” of the evolution process, must not be applied too liberally lest the process be severely disrupted. Usually, only a few mutations are performed in one generation.

If the “elitist” variant of the selection operator is used, one might also consider excepting the best solution in the population from being mutated. The reasons for doing so are explained in the paragraph above describing the selection operator.

5 Quantitative analysis

5.1 Preliminaries

The generation of queries for the benchmarks permits independent setting of the following parameters:

1. Class of the join graph
2. Distribution of relation cardinalities
3. Attribute domains

The shape of the join graph can be chosen from the following four classes: chain, star, cycle and grid (Fig. 16).

Relation cardinalities and domain sizes fall into four categories, S, M, L and XL, as specified in Table 2; for instance, 35% of all relations comprise between 1000 and 10 000 tuples. These figures were chosen such that join results and cost values are neither too small (because of the error that would be introduced due to the page granularity of the cost model) nor too large (loss of accuracy due to limited floating point arithmetic resolution).

The query itself is specified such that all relations from a particular join graph are to be joined; the selectivities that are associated with the graph’s edges are computed according to the estimate used in System-R [SAC⁺79], i.e. $\sigma = 1 / \min(\text{dom}(\text{attribute}_1), \text{dom}(\text{attribute}_2))$. Index structures (either hash tables or B^+ -trees) facilitate read access on 20% of all relation attributes. While constructing a join graph, relation cardinalities and attribute domain sizes are drawn independently; however, various “sanity checks” ensure that, for instance, a relation’s cardinality cannot exceed the product of its attribute domain sizes.

Each point in the following diagrams represents the average of at least thirty optimized queries, which proved to be a good compromise between the conflicting goals “avoidance of spurious results” (due to atypical behaviour of single runs) and “running time,” as preliminary tests showed.

From the optimization strategies discussed in Sect. 4, we implemented the following algorithms: the System-R algorithm, the minimum selectivity heuristic and the KBZ algorithm from the class of deterministic optimizers, and simulated annealing, iterative improvement and genetic (all three

in several variants) from the class of randomized/genetic optimizers. All deterministic algorithms yield solutions in the subspace of left-deep processing trees, whereas some of the randomized/genetic algorithms operate in the entire solution space (bushy trees).

All cost figures are scaled with respect to the best solution available apart from System-R (because the System-R algorithm could not be run for all parameter settings due to its high running time). For instance, a solution with a scaled cost of two is twice as expensive to evaluate as the best plan computed by any algorithm for that particular query. However, a curve for the algorithm that actually did compute the best solution is not necessarily shown in every plot. In other words, the set of algorithms that compete for the best solution is always the same, regardless of the subset that is depicted in a particular plot.

5.2 Solution spaces

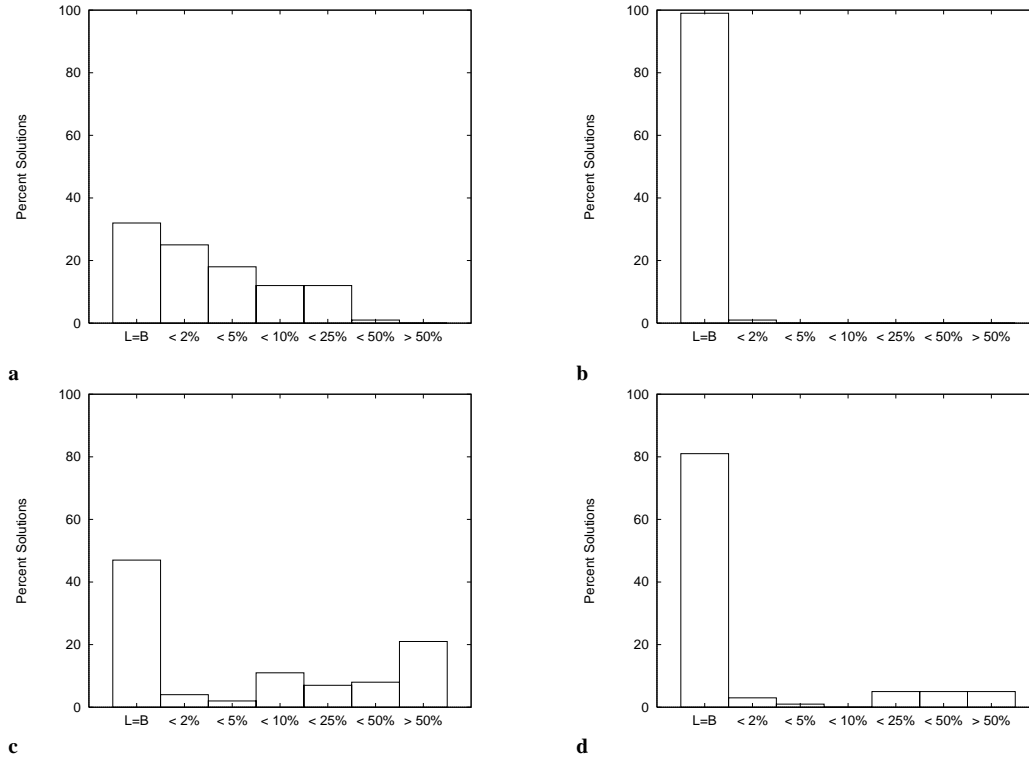
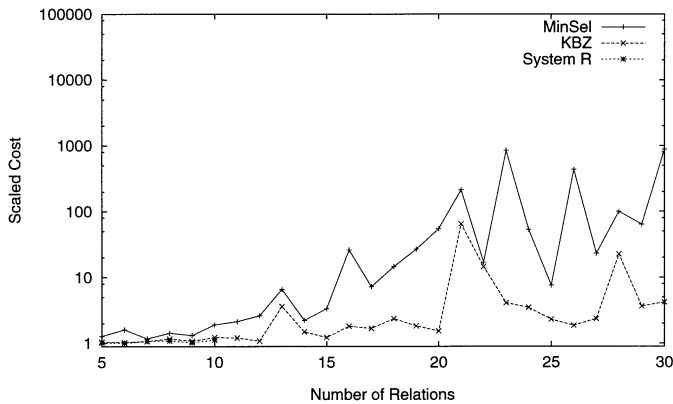
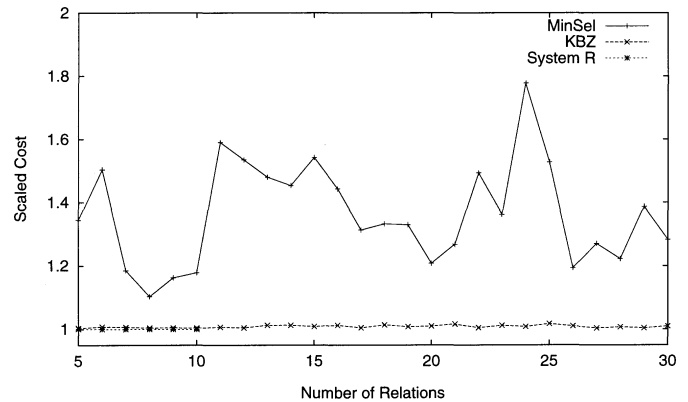
Before presenting the benchmark results, we will take a closer look at the two solution spaces. The left-deep tree space is a subset of the bushy tree space, so we can expect lower running times of optimizers that operate in the left-deep space. On the other hand, there is the danger of missing good solutions that are not left-deep trees.

In order to get some insight into the advantages of using one solution space instead of the other, we determined both the “left-deep optimal” and “bushy optimal” solutions for one hundred randomly selected queries with six participating relations. The histograms for the four different join graph types in Fig. 17a–d show the percentage of cases where the left-deep tree optimum and bushy tree optimum is of equal cost (i.e. the optimal solution is in fact a left-deep tree; labelled L=B). Following from left to right: the percentages of cases where the bushy tree optimum has less than 2%, between 2% and 5%, etc., lower cost than the left-deep tree optimum.

Considering these histograms, it becomes apparent that the shape of the join graph makes a big difference: for chain and cycle, we can find in half of all cases a better solution in the bushy tree solution space; for cycle, about one fifth even more than 50% cheaper than the best left-deep tree solution. Consequently, the investment in searching the bushy tree solution space should be profitable. On the other hand, for star join graphs in most of the cases the optima are left-deep trees anyway, because other tree shapes necessarily comprise cartesian products. Finally, for the grid join graph, the situation is not as clear as for the other three: 80% of the optima are left-deep trees, but a non-neglectable fraction of the bushy tree optima are far cheaper than their left-deep counterparts. A choice in favour of the bushy tree solution space would depend heavily on the optimization algorithms’ capability to locate these solutions (cf. [IK91]). In the remainder of this section, we will investigate whether the bushy tree optimizers can exploit the potential of good solutions in the bushy tree solution space.

Table 2. Relation cardinalities and domain sizes

Class	Relation cardinality	Percentage (%)	Class	Domain size	Percentage (%)
S	10–100	15	S	2–10	5
M	100–1000	30	M	10–100	50
L	1000–10000	35	L	100–500	30
XL	10000–100000	20	XL	500–1000	15
100%			100%		

**Fig. 17.** Left-deep versus bushy processing trees: **a** chain; **b** star; **c** cycle; **d** grid**Fig. 18.** Deterministic algorithms; chain join graph**Fig. 19.** Deterministic algorithms; star join graph

5.3 Benchmark results for deterministic algorithms

In the first series of benchmarks, we shall examine deterministic algorithms. Figures 18–21 show the results for the System-R algorithm (Sect. 4.1.1), the minimum selectivity heuristic (Sect. 4.1.2), and the KBZ algorithm (Sect. 4.1.3). Because none of the cost formulae in Sect. 2.2 fulfils the KBZ algorithm's requirement, we used a simple approximation that counts the processed tuples for a nested loop join

(without considering index structures) in order to be able to run the algorithm. The cost of a complete evaluation plan thus derived, however, was computed according to the exact formulae.

On each of the diagrams, the scaled cost (cost of the optimized evaluation plan divided by the cost of the best plan derived by *any* of the optimization algorithms discussed in this section except System-R) is plotted against the number of relations participating in the query; the join graph type

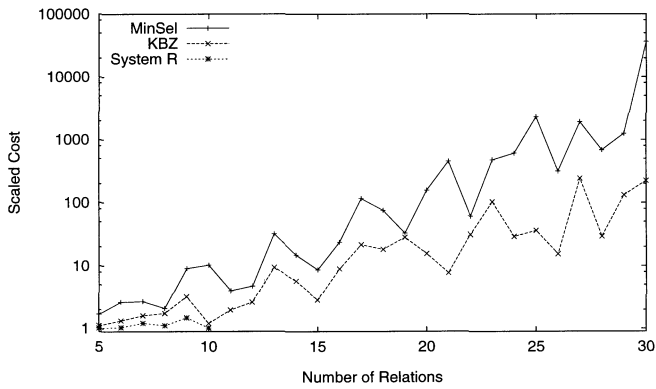


Fig. 20. Deterministic algorithms; cycle join graph

is noted in the respective caption. Note the smaller scale in the y -axis for the star join graph. All deterministic algorithms yield left-deep processing trees; in addition, the best join method is determined locally for each join node, i.e. proceeding bottom-up and selecting the least costly join method for each node. The results for the System-R optimization are plotted for five to ten participating relations in order to provide some “absolute” basis for comparison purposes.

Despite the simple cost approximation for running the KBZ algorithm, this optimizer turns out to be the best of the two (heuristic) deterministic optimizers for the chain, star and cycle join graphs. Especially, the solutions for the star join graph can hardly be improved by any of the other algorithms we tested, be it deterministic, randomized or genetic. For cycle and grid, the results are not quite as competitive, because these join graphs are cyclic and a spanning tree must be selected prior to the application of the KBZ algorithm. This effect becomes especially apparent for the grid join graph, where the minimum selectivity heuristic performs best. The System-R algorithm, which computes the optimal left-deep processing tree without cartesian products, achieves on the average a cost factor slightly above one, even though some solutions have cost factors below one, i.e. better than any of the approximate solutions. The reason why this cost factor is often higher than one is due to the limitation of the System-R algorithm to left-deep processing trees without cartesian products.

To summarize the results for heuristic optimizers, we can note the following points: firstly, both of the discussed optimizers have a very short running time – the KBZ algorithm managed to compute the results for the 30-relation star queries in less than 2s of CPU time each, the minimum selectivity heuristic in less than 0.1s for each one of the same queries. Secondly, the performance in terms of quality is – except for star join graphs – competitive only for small queries, where KBZ performs best for the join graphs with low connectivity and minimum selectivity for the grid join graph. However, for small queries, one would probably not want to rely on heuristic algorithms, but compute the optimal solution using some kind of search strategy.

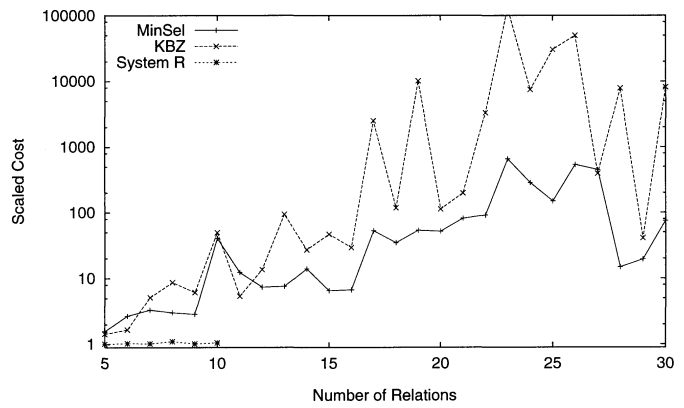


Fig. 21. Deterministic algorithms; grid join graph

5.4 Benchmark results for randomized and genetic algorithms

The next set of benchmarks is carried out with randomized algorithms (cf. Sect. 4.2) and genetic algorithms (cf. Sect. 4.3). We will compare three variants of iterative improvement (called IIJ, IIH [SG88] and IIIO [IK90]), and of simulated annealing (called SAJ, SAH [SG88] and SAIO [IK90]) and two variants of genetic algorithms (Genetic, BushyGenetic). Furthermore, the results of the System-R optimization are shown for five to ten participating relations. The parameters for each algorithm are derived from the cited references (II, SA) or they were determined in preliminary tests (Genetic, BushyGenetic). In addition, for all algorithms generating left-deep trees, a search proceeding from the leaves of the tree to the root is performed for all trial solutions in order to determine the most appropriate join method on each join node.

Exactly as in the first set of benchmarks with the heuristic algorithms, the scaled cost is plotted against the number of relations participating in the join – please note the different scale in Fig. 22. The parameters of the algorithms mentioned above are as follows:

1. SAJ

- A move is either a Swap or a 3Cycle, i.e. only left-deep processing trees are considered.
- The starting temperature is chosen such that at least 40% of all moves are accepted.
- The number of iterations of the inner loop is the same as the number of joins in the query.
- After every iteration of the inner loop, the temperature is reduced to 97.5% of its old value.
- The system is considered frozen when the best solution encountered so far cannot be improved in five subsequent outer loop iterations (i.e. temperature reductions) and less than 2% of the generated moves are accepted.

2. SAH

- A move is either a Swap or a 3Cycle, i.e. only left-deep processing trees are considered.
- The starting temperature is determined as follows: the standard deviation σ for the cost is estimated from a set of sample solutions and multiplied by a constant value (20).

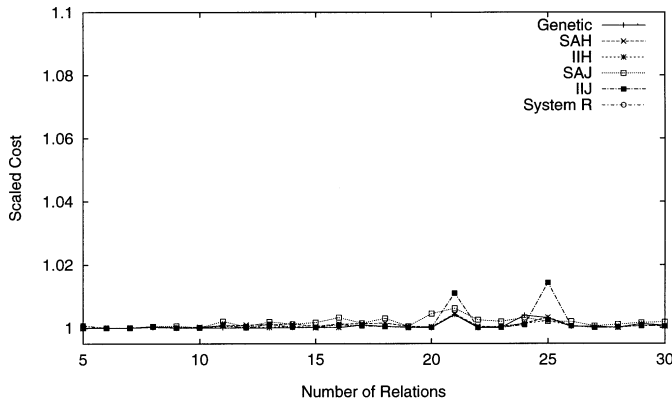


Fig. 22. Randomized algorithms, left-deep tree solution space; chain join graph

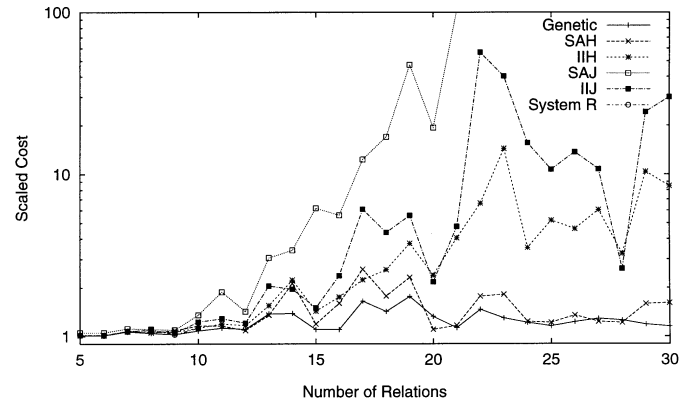


Fig. 23. Randomized algorithms, left-deep tree solution space; star join graph

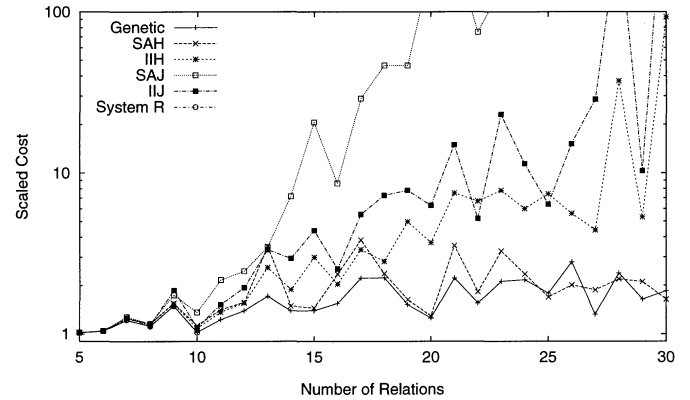


Fig. 24. Randomized algorithms, left-deep tree solution space; cycle join graph

- The inner loop is performed until the cost distribution of the generated solutions is sufficiently stable (for details see [SG88]).
- After every iteration of the inner loop, the temperature is multiplied by $\max(0.5, e^{-\frac{\lambda t}{\sigma}})$ ($\lambda = 0.7$; σ see above).
- The system is considered frozen when the difference between the minimum and maximum costs among the accepted states at the current temperature equals the maximum change in cost in any accepted move at the current temperature.

3. SAIO

- Moves are chosen from join method change, commutativity, associativity, left join exchange and right join exchange. The entire solution space (bushy processing trees) is considered.
- The starting temperature is twice the cost of the (randomly selected) starting state.
- The number of iterations of the inner loop is 16 times the number of joins in the query.
- After every iteration of the inner loop, the temperature is reduced to 95% of its old value.
- The system is considered frozen when the best solution encountered so far cannot be improved in four subsequent outer loop iterations (i.e. temperature reductions) and the temperature falls below one.

4. Iterative improvement (IIH, IIJ, IIIO)

- All starting points are chosen randomly.
- For an algorithm II_x , moves are chosen from the same set as the corresponding SA_x algorithm.
- Local minima are determined according to [SG88] (IIH, IIJ) and [IK90] (III O), i.e. a solution is considered a local minimum if k randomly selected neighbours fail to improve the result. k is the number of join graph edges for IIH and IIJ; for III O, k is the number of neighbouring states.
- In order to perform a “fair” comparison between II and SA, the total number of solutions considered is approximately the same for both the corresponding II_x and SA_x algorithms.

5. Two-phase optimization (III O+SAIO)

- Ten random starting points for the II phase.

- SA phase starts with the minimum from the II phase and the starting temperature is 0.1 times its cost.
- ### 6. Genetic algorithms (Genetic/BushyGenetic)
- Solution space: left-deep processing trees/
bushy processing trees
 - Encoding: ordered list of leaves/
ordinal number encoding
 - Ranking-based selection operator
 - Sequence exchange crossover operator
 - Population: 128
 - Crossover rate 65% (65% of all members of the population participate in crossover)
 - Mutation rate 5% (5% of all solutions are subject to random mutation)
 - Termination condition: 30 generations without improvement / 50 generations without improvement

In Figs.22–25, the results for the left-deep tree optimizers are depicted. Although the parameter setting for SAH/SAJ and IIH/IIJ is similar, we note that the J variants perform poorly for all but one of the join graph types. SAH and IIH perform much better, where, in turn, SAH is superior to IIH. In all cases, SAH and the genetic algorithm computed the best evaluation plans among the left-deep tree optimizers, with a slight superiority of the genetic algorithm. Apparently, the sophisticated equilibrium/freezing condition for SAH is the main reason for its good results. A closer look at the benchmark data revealed that indeed SAJ visited

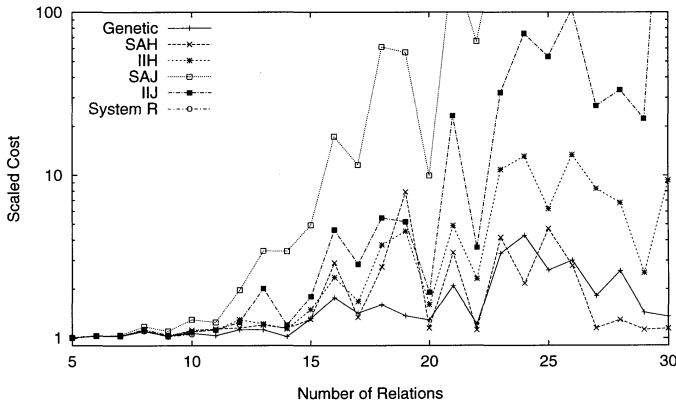


Fig. 25. Randomized algorithms, left-deep tree solution space; grid join graph

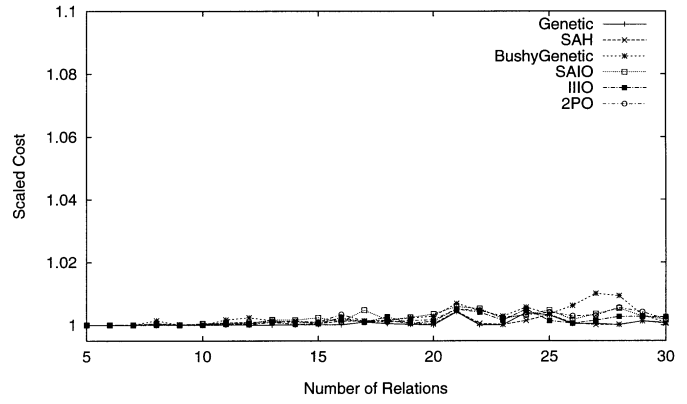


Fig. 27. Randomized algorithms, bushy tree solution space; star join graph

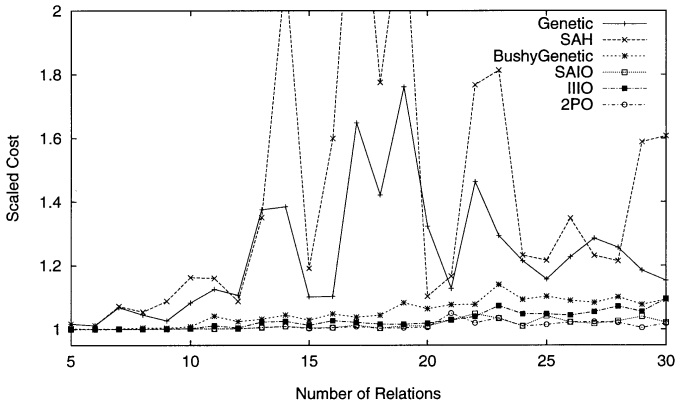


Fig. 26. Randomized algorithms, bushy tree solution space; chain join graph

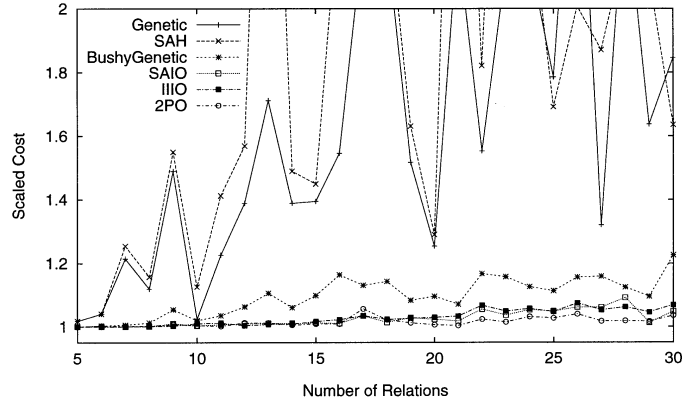


Fig. 28. Randomized algorithms, bushy tree solution space; cycle join graph

much less solution alternatives than SAH. The II variants that were designed to consider about as many different solutions as the respective SA algorithms reflect this fact: IIH achieves better results than IIJ. Apart from the quality of the derived results, another important criterion for selecting an optimizer is its running time, which we will investigate later. In the meantime, we will look at the performance of those optimizers that operate in the bushy tree solution space.

These optimizers, namely SAIO, IIIO, two-phase optimization (2PO) and a genetic algorithm (BushyGenetic), are compared in Figs. 26–29. In addition, the best two left-deep optimizers’ curves (SAH and Genetic) are included in order to facilitate direct comparison. It turns out that, in terms of quality, none of the implemented algorithms performed better than the 2PO algorithm, regardless of the join graph type, although the gap between SA, II and 2PO is quite narrow. In the left-deep case the genetic algorithm showed a slight superiority over the SA results, but this is not so in the bushy tree solution space. Although the genetic algorithm does not perform particularly poor, it cannot quite equal the quality of SA or 2PO.

Only for star queries, all algorithms exhibit very similar behaviour (divergence just about 1%), so the algorithms’ running time would be the decisive factor in this case. For all other join graphs, every bushy tree optimizer easily outperforms even the best implemented left-deep tree optimizer, which confirms that these algorithms are indeed capable of

locating the superior solutions of the bushy tree solution space.

Let us now look at the running times for the different optimizers. Although the quality of the generated solutions is a very important characteristic, the running time of an algorithm has a considerable impact on the final choice. The intended application area determines how much time can be spent on the optimization: queries that are stated interactively and run only once do not warrant the same amount of optimization as compiled queries that are repeated hundreds or thousands of times. In Fig. 30, the average running times (CPU times) for Genetic, SAH, BushyGenetic, SAIO, IIIO, 2PO and System-R are plotted against the number of relations participating in the queries (chain join graph). The running times for the various algorithms were determined on a SPARCstation 20/612MP.

From the six randomized/genetic algorithms, SAIO has the longest running times with up to 2800s of CPU time for 30-relation queries. Although 2PO yields slightly better solutions, it requires a running time of about 1300s, only half of SAIO’s time. As expected, left-deep tree optimizers (SAH, Genetic) run faster than bushy tree optimizers, but gain of speed must be paid by loss of quality. Surprisingly, the BushyGenetic algorithm runs even faster than both the left-deep optimizers, even though it yields solutions that are at least as good as theirs: it can handle the 30-relation queries on average in about 500s. In contrast to the random-

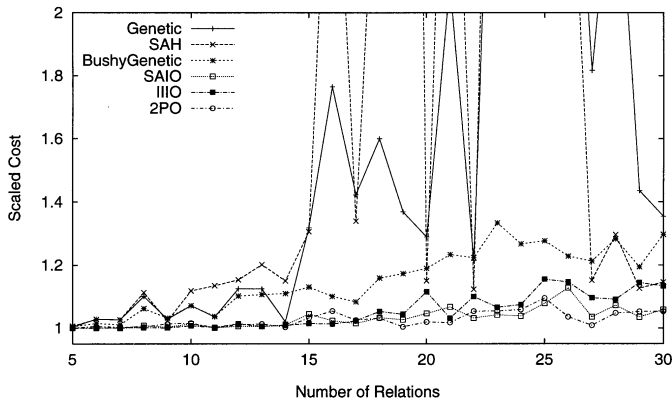


Fig. 29. Randomized algorithms, bushy tree solution space; grid join graph

ized/genetic algorithms, the computation of the optimal left-deep tree without cartesian products is in our implementation feasible only for queries involving up to 12 twelve relations; for instance, 13-relation queries already require on average an optimization time of more than 1 h.

In Fig. 31, for the same six algorithms, the approach to their respective final solutions is shown. Each time the currently best solution is improved, the gain in absolute cost units is noted (y -axis), together with the time of its occurrence (x -axis). Because only a single optimization run for a 20-relation query (chain join graph) is plotted, we cannot draw any far reaching conclusions, but nevertheless the curves reflect the algorithms' typical behaviour quite well. Both simulated annealing algorithms, SAIO as well as SAH, spend a good deal of the total running time investigating high-cost processing trees; SAIO required more than 150s to reach a cost level of less than ten times the cost of the final solution. SAH ran faster, but it still took a very long time for the approach to its final solution. On the other hand, both 2PO and iterative improvement (IIIO) achieved a very good result in less than 1s and less than 3s running time, respectively. For the genetic algorithms, it is especially BushyGenetic that can reach acceptable solutions very quickly. Even the initial population consisted of at least one member with an evaluation cost that is as low as SAIO's after running more than 100 times as long. Although the drawing of the initial population is not guaranteed to be unbiased in our implementation, we can note that genetic algorithms nicely supplement the approach in [GLPK94] (Sect. 4.2.5): in a first step, a random sample could be drawn using the algorithm presented in [GLPK94], which can be used in a second step as the initial population for the genetic algorithm.

5.5 Summary

Comparing the performance of the various optimization algorithms, we can draw the following conclusions.

Algorithms that perform an exhaustive or near exhaustive enumeration of the solution space, such as dynamic programming, can compute the optimal result, but the high running time makes their application only feasible for queries that are not too complex (i.e. less than about 10–15 relations for left-deep processing trees). For the same reason, searching

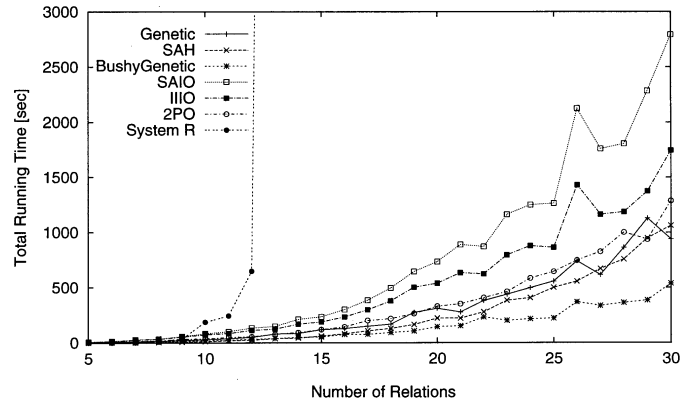


Fig. 30. Total running times

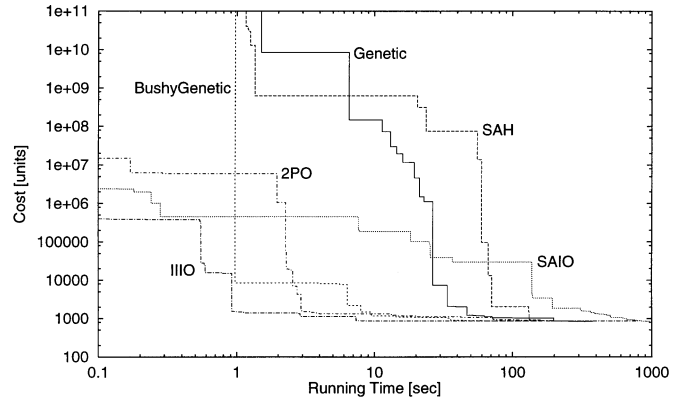


Fig. 31. Approach to the final solution

the bushy tree solution space can be carried out only for very simple queries (in our experiments, about six to seven relations), so the advantages of this solution space can hardly be exploited.

Heuristic optimizers avoid the high time complexity of exhaustive enumeration, but the results are, especially for complex queries with many participating relations, rarely acceptable. The KBZ algorithm, although yielding the optimal left-deep solution under certain circumstances, is difficult to apply in practice: the need for cost model approximations and problems concerning join method assignment limits its usefulness. We found that only for star queries is the KBZ algorithm competitive; its short running time compared to alternative algorithms (especially randomized/genetic) makes it the solution of choice.

Finally, randomized and genetic algorithms operating in the bushy tree solution space are the most appropriate optimizers in the general case, provided the problems are too complex to be tackled by exhaustive enumeration. The question of which of the discussed algorithms is the most adequate depends on the particular application area, namely whether short running time or best optimization performance is the primary goal. If good solutions are of highest importance, 2PO, the algorithm that performed best in our experiments, is a very good choice; other SA variants, for instance toured simulated annealing (TSA; [LVZ93]), that we did not implement are likely to achieve quite similar results. The "pure" SA algorithm has a much higher running time without yielding significantly better solutions. If short running

time is more important, iterative improvement (IIIO), the genetic algorithm (BushyGenetic), and, to a lesser extent, 2PO are feasible alternatives; especially the first two degrade gracefully if they are preempted: in the example run in Fig. 31, they achieved acceptable results in less than 1s. Moreover, as mentioned above, genetic algorithms can be combined very well with the transformationless approach in [GLPK94].

6 Conclusion

We have studied several algorithms for the optimization of join expressions. Because of new database applications, the complexity of the optimization task has increased; more relations participate in join expressions than in traditional relational database queries. Enumeration of all possible evaluation plans is no longer feasible. Algorithms that compute approximate solutions, namely heuristic, randomized and genetic algorithms, show different capabilities for solving the optimization task. Heuristic algorithms compute solutions very quickly, but the evaluation plans are in many cases far from the optimum. Randomized and genetic algorithms are much better suited for join optimizations; although they require a longer running time, the results are far better.

For the question of the adequate solution space, we have found that, with the exception of the star join graph, the bushy tree solution space is preferable in spite of the fact that “pipelining” (avoiding writing intermediate results to secondary memory) can be carried out mainly by left-deep processing trees.

Another consideration is the extensibility of randomized and genetic algorithms: both can be designed to optimize not merely pure join expressions, but complete relational queries. In addition, some of them (namely the iterative improvement and genetic algorithms) can be easily modified to make use of parallel computer architectures.

Acknowledgements. We should like to thank the referees for their valuable comments and suggestions on an earlier draft of this paper, and our colleagues for countless CPU-hours on their workstations. This work was partially supported by grant DFG Ke 401/6-1 and SFB 346.

References

- [BF191] Bennet K, Ferris MC, Ioannidis YE (1991) A genetic algorithm for database query optimization. In: Proc 4th Int Conf Genetic Algorithms, San Diego, Calif, pp 400–407
- [EN94] Elmasri E, Navathe SB (1994) Fundamentals of database systems, 2nd edn. Benjamin/Cummings, Redwood City, Calif
- [GLPK94] Galindo-Legaria C, Pellenkoff A, Kersten M (1994) Fast, randomized join-order selectoin – why use transformations? In: Proc Conf Very Large Data Bases (VLDB), Santiago, Chile, September, pp 85–95
- [Gol89] Goldberg DE (1989) Genetic algorithms in search, optimization & machine learning. Addison-Wesley, Reading, Mass
- [IK84] Ibaraki T, Kameda T (1984) Optimal nesting for computing N -relational joins. ACM Trans Database Syst 9:482–502
- [IK90] Ioannidis YE, Kang YC (1990) Randomized algorithms for optimizing large join queries. In: Proc ACM SIGMOD Conf Management of Data, Atlantic City, NJ, April, pp 312–321
- [IK91] Ioannidis YE, Kang YC (1991) Left-deep vs. bushy trees: an analysis of strategy spaces and its implications for query optimization. In: Proc ACM SIGMOD Conf Management of Data, Denver, Colo, May, pp 168–177
- [IW87] Ioannidis YE, Wong E (1987) Query optimization by simulated annealing. In: Proc ACM SIGMOD Conf Management of Data, San Francisco, Calif, pp 9–22
- [KBZ86] Krishnamurthy R, Boral H, Zaniolo C (1986) Optimization of non-recursive queries. In: Proc Cof Very Large Data Bases (VLDB), Kyoto, Japan, pp 128–137
- [KM92] Kemper A, Moerkotte G (1992) Access support relations: an indexing method for object bases. Inf Syst 17:117–146
- [KM94] Kemper A, Moerkotte G (1994) Object-oriented database management: applications in engineering and computer science. Prentice-Hall, Englewood Cliffs, NJ
- [Kru56] Kruskal JB (1956) On the shortest spanning subtree of a graph and the travelling salesman problem. Proc Am Math Soc 7:48–50
- [Law78] Lawler E (1978) Sequencing jobs to minimize total weighted completion time subject to precedence constraints. Ann Discrete Math 2:75–90
- [LV93] Lanzelotte R, Valduriez P, Zait M (1993) On the effectiveness of optimization search strategies for parallel execution spaces. In: Proc Conf Very Large Data Bases (VLDB), Dublin, Ireland, pp 493–504
- [ME92] Mishra P, Eich MH (1992) Join processing in relational databases. ACM Comput Surv 24:63–113
- [MS79] Monma C, Sidney J (1979) Sequencing with series-parallel precedence constraints. Math Oper Res 4:215–224
- [OL90] Ono K, Lohman GM (1990) Measuring the complexity of join enumeration in query optimization. In: Proc Conf Very Large Data Base (VLDB), Brisbane, Australia, pp 314–325
- [SAC⁺79] Selinger PG, Astrahan MM, Chamberli DD, Lorie RA, Price TG (1979) Access path selection in a relational database management system. In: Proc ACM SIGMOD Conf Management of Data, Boston, USA, May, pp 23–24
- [SG88] Swami A, Gupta A (1988) Optimization of large join queries. In: Proc ACM SIGMOD Conf Management of Data, Chicago, Ill, May, pp 8–17
- [Sha86] Shapiro LD (1986) Join processing in database systems with large main memories. ACM Trans Database Syst 11:239–264
- [SI93] Swami A, Iyer B (1993) A polynomial time algorithm for optimizing join queries. In: Proc IEEE Conf Data Eng, Vienna, Austria, April, pp 345–354
- [Swa89] Swami A (1989) Optimization of large join queries: combining heuristics and combinatorial techniques. In: Proc ACM SIGMOD Conf Management of Data, Portland, Ore, May, pp 367–376
- [VM96] Vance B, Maier D (1996) Rapid bushy join-order optimization with cartesian product. In: Proc ACM SIGMOD Conf Management of Data, Montreal, June, pp 35–46
- [WY76] Wong E, Youssefi K (1976) Decomposition – a strategy for query processing. ACM Trans Database Syst 1:223–241
- [YW79] Youssefi K, Wong E (1979) Query processing in a relational database management system. In: Proc Conf Very Large Data Bases (VLDB), Rio de Janeiro, Brazil, pp 409–417