

Printed at the Mathematical Centre, 49, 2e Boerhaavestraat, Amsterdam.

The Mathematical Centre, founded the 11-th of February 1946, is a non-profit institution aiming at the promotion of pure mathematics and its applications. It is sponsored by the Netherlands Government through the Netherlands Organization for the Advancement of Pure Research (Z.W.O).

MATHEMATICAL CENTRE TRACTS 81

**FOUNDATIONS OF
COMPUTER SCIENCE II**

K.R. APT (ed.)

J.W. DE BAKKER (ed.)

PART 1

MATHEMATISCH CENTRUM

AMSTERDAM 1979

AMS(MOS) subject classification scheme (1970): 68A10

ACM - Computing Reviews - categories: 5.25, 5.32

ISBN 90 6196 140 8

First printing 1976

Second printing 1979

Preface		i
E.L. LAWLER:	<i>Graphical algorithms and their complexity</i>	3
J. VAN LEEUWEN:	<i>The complexity of data organization</i>	37

PREFACE

The Second Advanced Course on the Foundations of Computer Science was organized by the Mathematical Centre as part of an international cooperation under the auspices of the European Communities, and took place at the University of Amsterdam, May 31 - June 11, 1976. The Course consisted of six series of lectures, and the notes of five of these have been collected in the present Tracts, jointly edited by K.R. Apt and the undersigned.

I am very grateful to the Netherlands Organization for the Advancement of Pure Research (Z.W.O.) and to the European Communities for generously supplying the money to organize the Course, to the lecturers for their excellent contributions, to Mrs. S.J. Kuipers-Hoekstra for her invaluable help with the organization, and to the staff of the publication services of the Mathematical Centre for their undaunted efforts in the typing and printing of these books.

J.W. de Bakker

Director of the Course

1. Introduction	3
2. Topological ordering	4
3. Recognition of series parallel digraphs.	5
4. Isomorphism of transitive series parallel digraphs	16
5. Subgraph isomorphism.	19
6. Finding minimum spanning trees	20
7. Generating all maximal independent sets.	25
8. Finding a maximum independent set.	28
9. Computing the chromatic number of a graph.	29
References.	31

GRAPHICAL ALGORITHMS AND THEIR COMPLEXITY

E.L. LAWLER

University of California, Berkeley, USA

1. INTRODUCTION

Graphs have important applications in all branches of science and technology. It follows that algorithms which solve problems on graphs are important for the problems they solve. But there are additional reasons why graphs and graphical algorithms have particular fascination for computer scientists. Graphical algorithms pose a challenging array of problems in data structures, algorithmic analysis, and complexity theory.

There is now a great wealth of literature on the subject of graphical algorithms. The present notes are not intended as a survey. The author has made a rather arbitrary and personal selection of problems to discuss, with the expectation that these problems may serve as an introduction to the methodology of the subject area.

As a first, very simple, problem, we find a "topological ordering" of the nodes of an acyclic digraph. This can be done in $O(n+m)$ time, provided the digraph is specified by an "arc-list" data structure. (Hereafter, we consistently let n denote the number of nodes and m the number of arcs.)

We then consider a "recognition" problem; determining if an acyclic digraph is "series parallel". By careful tailoring of the algorithm, we are able to obtain a bound of $O(n^2)$ on the running time - a second example of an algorithm with polynomial-bounded running time.

We then take a brief excursion into the problem of graphic isomorphism. There is no known polynomial-bounded algorithm for determining whether or not two graphs are isomorphic, nor is there convincing evidence that such an algorithm does not exist. Yet we can show that for the special case of transitive series parallel digraphs there is such an algorithm.

A very basic and important problem is that of finding a minimum cost spanning tree in an undirected graph. Various algorithms have been proposed to solve this problem in $O(n^2)$ or $O(m \log n)$ time.

Recently, several algorithms have been proposed which have $O(m \log \log n)$ running time. We give a brief description of the simplest of these new algorithms. The reduced running time is obtained through the clever use of data structures for the implementation of priority queues.

Given a graph, one sometimes wants to generate all subgraphs of a particular type, e.g. trees, cycles. As typical of a problem of this type, we consider the generation of all maximal independent sets. A newly developed algorithm is described, and shown to have running time which is linear in the number of maximal independent sets.

The problem of finding a maximum (cardinality) independent set is "NP-complete", and hence it is extremely unlikely that a polynomial-bounded algorithm exists. One obvious way to solve the problem is to generate all maximal independent sets and to pick out one with a maximum number of elements. However, it is possible to do significantly better than this, and we indicate techniques one can apply to reduce the exponential running time of the search procedure.

Finally, we consider another very difficult NP-complete problem; computing the chromatic number of a graph. We describe the application of dynamic programming to this problem, and make use of the previous results on generating maximal independent sets, to obtain an algorithm whose running time is bounded by a relatively low order exponential function.

2. TOPOLOGICAL ORDERING

It is a well known result that a digraph $G = (N, A)$ is acyclic (contains no directed cycles) if and only if it is possible to "topologically" number its nodes, so that $(i, j) \in A$ only if $i < j$. (That is, all arcs are directed from a lower-numbered node to a higher-numbered node.) How efficiently can we test an arbitrary digraph G to determine if it is acyclic, and if it is acyclic, to determine a topological numbering of its nodes?

The "classical" procedure for solving this problem is as follows. Find a node with in-degree zero, give that node the number 1, and delete it from G . Then find another node with in-degree zero, give it number 2 and delete it. Continue in this way until either all nodes have been num-

bered or until there is no node with in-degree zero in the remaining sub-graph (in which case G is not acyclic).

The efficiency with which this procedure can be implemented depends critically on the way in which the digraph is specified to the computer. If G is represented by an adjacency matrix, there is no way to implement the procedure in less than $O(n^2)$ time. In fact, it has been proved that no algorithm can solve this problem (or any problem in a much wider class) without "looking at" at least $O(n^2)$ entries in the adjacency matrix [12].

We can, however, do much better if G is specified by means of "arc lists". That is, for each node x , there is a simple linear list of all the nodes to which there are arcs directed from x . Then we can proceed as follows: First compute the in-degree of each node. This is done quite easily in $O(n+m)$ time by simply scanning the arc lists. (When y is found in the list for x , the in-degree of y is incremented by one.) Then scan through the array of in-degrees, picking out those nodes which have in-degree zero, and make a list of these nodes. This requires $O(n)$ time. Then choose a node with in-degree zero and give it the number 1. The arc list for this node is used to decrement the in-degrees of the nodes to which it has arcs. If any of these in-degrees are reduced to zero, the node in question is placed on the list of nodes with in-degree zero. Then choose another node with in-degree zero and repeat. It is apparent that the decrementation operations require only $O(m)$ time overall, because there are exactly m entries in the arc lists. Hence the algorithm is seen to require $O(n+m)$ time overall.

This very simple example serves to illustrate the importance of the data structure which is used to specify a graph or digraph.

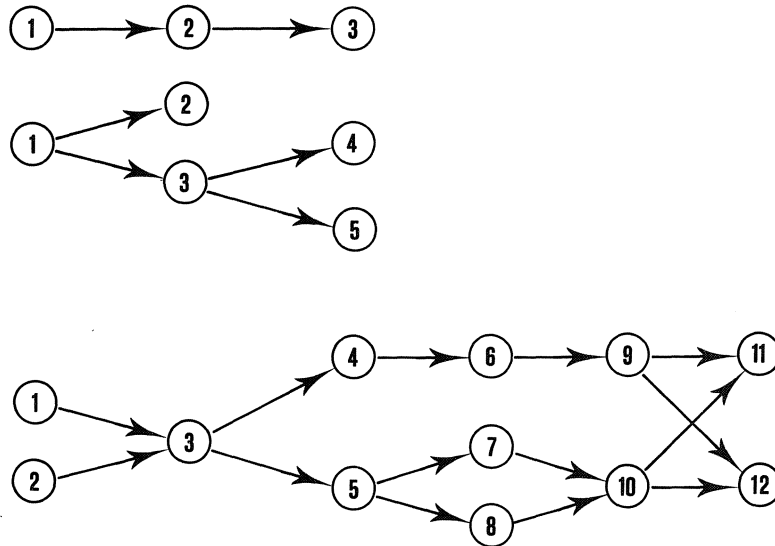
3. RECOGNITION OF SERIES PARALLEL DIGRAPHS

"Series parallel" digraphs have a number of very useful properties. In particular, a number of combinatorial problems which are NP-complete for general digraphs are computationally tractible for the special case of series parallel digraphs [5]. In the present section, we shall concern ourselves with the problem of recognizing series parallel digraphs, i.e. testing an arbitrary digraph to determine if it is series parallel. In later sections, we shall discuss some interesting properties of these digraphs.

First let us give a recursive definition of this class of digraphs.

A digraph on a single node is *transitive series parallel*. If $G_1 = (N_1, A_1)$ and $G_2 = (N_2, A_2)$, $N_1 \cap N_2 = \phi$, are transitive series parallel digraphs, then $G_1 \cup G_2 = (N_1 \cup N_2, A_1 \cup A_2)$, the *parallel composition* of G_1, G_2 , is transitive series parallel. Also, $G_1 \rightarrow G_2 = (N_1 \cup N_2, A_1 \cup A_2 \cup (N_1 \times N_2))$, the *series composition* of G_1, G_2 , is transitive series parallel. Only those digraphs which can be obtained by a finite number of series and parallel compositions of single-node digraphs are transitive series parallel. Any digraph whose transitive closure is transitive series parallel is *series parallel*.

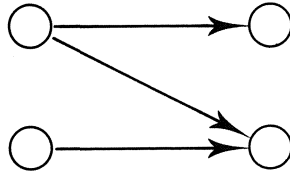
Some examples of series parallel digraphs are given in Figure 1.



Examples of Series Parallel Digraphs

Figure 1

Clearly every series parallel digraph is acyclic. Every acyclic digraph with three nodes or less is series parallel. The only 4-node acyclic digraph that is not series parallel is the one shown in Figure 2. It is not hard to show that if a digraph contains this 4-node digraph as an induced subgraph, then the digraph is not series parallel. The recognition procedure will provide a proof of the converse for transitive digraphs. Thus, a digraph is series parallel if and only if its transitive closure does not contain the digraph of Figure 2 as a "forbidden subgraph".

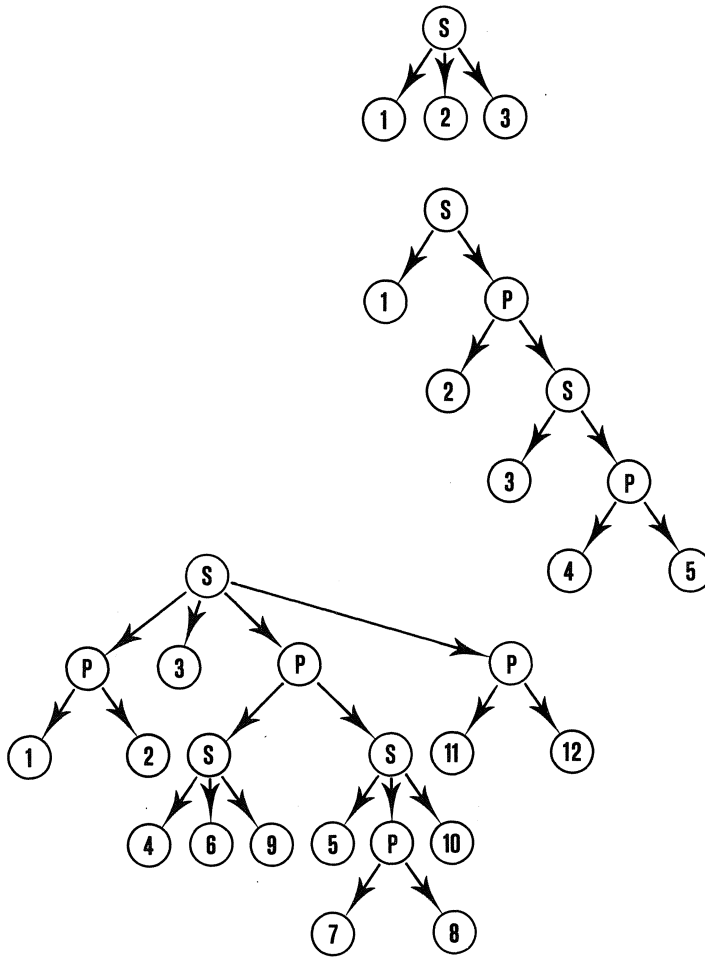


Forbidden Subgraph

Figure 2.

The structure of a series parallel digraph is represented by a "decomposition tree". Decomposition trees for the digraphs in Figure 1 are shown in Figure 3. Each leaf of the decomposition tree of G is identified with a node of G . Each internal node is designated "S" or "P", indicating serial and parallel composition of the subgraphs identified with its sons. The ordering of the sons of a P-node is unimportant, whereas the ordering of the sons of S-node is significant, indicating the order of series composition.

Let us adopt the convention that no son of an S-node is an S-node and no son of a P-node is a P-node. We assert that, subject to this assumption, the decomposition tree of a series parallel digraph is uniquely determined, up to a permutation of the sons of P-nodes.



Decomposition Trees for Digraphs in Figure 1

Figure 3

Our recognition procedure attempts to construct the decomposition tree for a given digraph G . If, at any step, it is not possible to carry the construction further, then we shall be able to exhibit the forbidden subgraph as an induced subgraph of the transitive closure of G , proving that G is not series parallel.

Let us suppose that the nodes of the given digraph $G = (N, A)$ are topologically ordered. As we have seen, such a numbering can be obtained

for an acyclic digraph in $O(n+m)$ time.

Our strategy will be to obtain the decomposition tree T_{j+1} for the subgraph induced on nodes $1, 2, \dots, j+1$ from the tree T_j obtained for nodes $1, 2, \dots, j$. It is, of course, a simple matter to obtain T_2 , to start the procedure.

To facilitate our task, we shall compute the values of two Boolean variables, SER and PAR, for each node of T_j . At each leaf node i ,

$$\left. \begin{array}{l} \text{SER}(i) = 1 \\ \text{PAR}(i) = 0 \end{array} \right\} \text{ if } (i, j+1) \in A$$

and

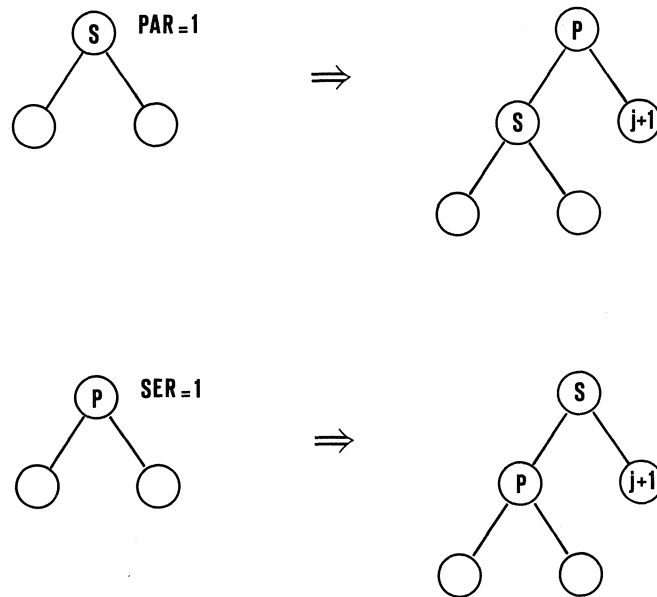
$$\left. \begin{array}{l} \text{SER}(i) = 0 \\ \text{PAR}(i) = 1 \end{array} \right\} \text{ if } (i, j+1) \notin A$$

For each internal node x of T_j , $\text{PAR}(x) = 1$ if and only if $\text{PAR}(y) = 1$ for each son y of x . For each internal P-node x , $\text{SER}(x) = 1$ if and only if $\text{SER}(y) = 1$ for each son y of x . For each internal S-node x , $\text{SER}(x) = 1$ if and only if $\text{SER}(y) = 1$ for the rightmost son y of x .

It should be clear that the values of SER and PAR can be computed for all nodes of T_j in $O(n)$ time, proceeding from the leaf nodes toward the root node.

In words, $\text{PAR}(x) = 1$, if and only if node $j+1$ is "in parallel" with the subgraphs induced by the nodes which are descendants of x in T_j . (Note that this property depends upon topological ordering.) Node $j+1$ is "in series" with the subgraph induced by node x of T_j if and only if either $\text{SER}(x) = 1$ or else x is a P-node and $\text{PAR}(x') = 0$, where x' is a brother to the right of x . Clearly, it is not possible for both $\text{SER}(x) = 1$ and $\text{PAR}(x) = 1$.

Now let us try to insert node $j+1$ into T_j , to obtain the tree T_{j+1} . If either $\text{SER}(r) = 1$ or $\text{PAR}(r) = 1$ at the root node r , then $j+1$ is in series or in parallel, respectively, with the subgraph induced on nodes $1, 2, \dots, j$. In particular, if r is a P-node and $\text{SER}(r) = 1$ or if r is an S-node and $\text{PAR}(r) = 1$, we must create a new root node, as indicated



Creation of New Root Node

Figure 4

in Figure 4. These are special cases. Otherwise, follow the procedure for S-nodes or P-nodes, as appropriate, outlined below.

S-Node Procedure

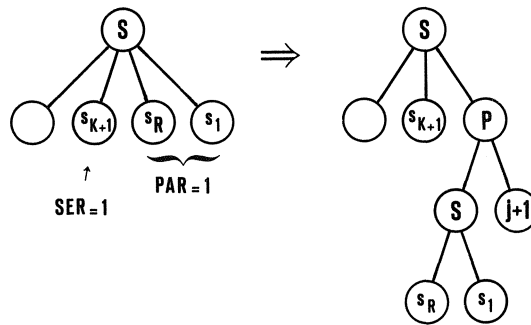
It is assumed that $PAR(r) = 0$ for the S-node r .

Case 1. If $SER(r) = 1$, then insert $j+1$ into T_j as the new rightmost son of

r. This completes the construction of T_{j+1} .

Case 2. Suppose Case 1 does not apply. Let the sons of r , indexed from right to left, be s_1, s_2, \dots, s_p ; $p \geq 2$. If $\text{PAR}(s_1) = 0$ or if $\text{SER}(s_2) = 1$, then node $j+1$ is in series with the subgraph induced on the descendants of s_2, s_3, \dots, s_p . Consequently, if $j+1$ is to be inserted anywhere into T_j , $j+1$ must be inserted somewhere into the subtree rooted at s_1 . If s_1 is itself a leaf node (in which case $\text{PAR}(s_1) = 1$), replace s_1 by a new P-node with s_1 and $j+1$ as sons, thereby completing the construction of T_{j+1} . If s_1 is a P-node, then apply the P-node procedure with r set to s_1 . (Note that in this case $\text{SER}(s_1) = 0$.)

Case 3. Suppose Cases 1 and 2 do not apply. (In this case, $\text{PAR}(s_1) = 1$ and $\text{SER}(s_2) = 0$.) If $\text{PAR}(s_1) = \text{PAR}(s_2) = \text{PAR}(s_k) = 1$, for $k \geq 2$, and $\text{SER}(s_{k+1}) = 1$, then node $j+1$ is in parallel with the subgraph induced on the descendants of s_1, s_2, \dots, s_k and in series with the subgraph induced on the descendants of $s_{k+1}, s_{k+2}, \dots, s_p$. Consequently, $j+1$ can be inserted into T_j by creating a new S-node and P-node, as shown in Figure 5. This completes the construction of T_{j+1} .

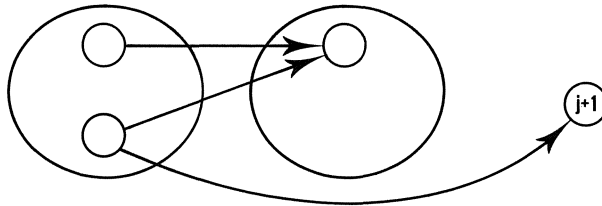


Insertion of Node $j+1$ in Case 3 of S-node Procedure

Figure 5

Case 4. Suppose Cases 1, 2, and 3 do not apply. That is, $\text{PAR}(s_1) = \text{PAR}(s_2) = \dots = \text{PAR}(s_k) = 1$, $k \geq 1$, and $\text{PAR}(s_{k+1}) = 0$, $\text{SER}(s_{k+1}) = 0$. In this case it is possible to identify the forbidden subgraph in the transitive closure of the subgraph induced on nodes $1, 2, \dots, j+1$, as shown in Fig. 6.

$$\begin{array}{l} \text{SER}(s_{k+1}) = 0 \\ \text{PAR}(s_{k+1}) = 0 \end{array} \quad \text{PAR}(s_1) = 1$$



Forbidden Subgraph
Identified in Case 4 of S-Node Procedure
Figure 6.

P-Node Procedure

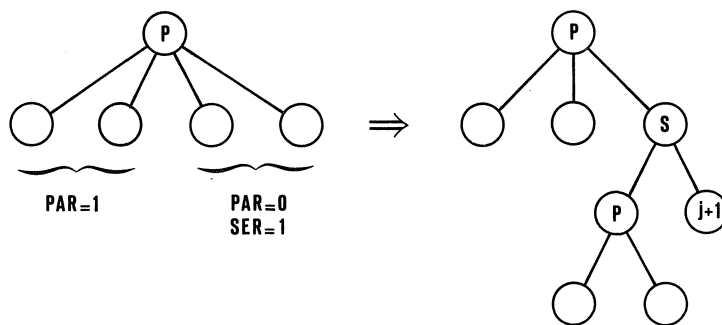
It is assumed that $\text{SER}(r) = 0$ for the P-node r .

Case 1. If $\text{PAR}(r) = 1$, then insert $j+1$ into T_j as a new son of r . This completes the construction of T_{j+1} .

Case 2. Suppose Case 1 does not apply. Let the sons of r , indexed arbitrarily, be s_1, s_2, \dots, s_p ; $p \geq 2$. If $\text{PAR}(s_i) = 1$ for all $i \neq k$, and $\text{PAR}(s_k) = 0$, then node $j+1$ is in parallel with the subgraph induced on the descendants of s_i , $i \neq k$. Consequently, if $j+1$ is to be inserted anywhere into T_j , $j+1$ must be inserted somewhere into the subtree rooted to s_k . If s_k is itself a leaf node (in which case $\text{SER}(s_k) = 1$), replace s_k by a new S-node with s_k and $j+1$ as sons, thereby completing the construction of T_{j+1} . If s_k is an S-node, then apply the S-node procedure with r set to s_k . (Note that in this case $\text{PAR}(s_k) = 0$.)

Case 3. Suppose Cases 1 and 2 do not apply. (In this case, there are at least two sons s_i such that $PAR(s_i) = 0$.) If $SER(s_i) = 1$ for each son s_i such that $PAR(s_i) = 0$, then node $j+1$ is in parallel with the subgraph induced on the descendants of all sons s_i for which $PAR(s_i) = 1$ and in series with the subgraph induced on the descendants of all sons s_i for which $PAR(s_i) = 0$. Consequently, $j+1$ can be inserted into T_j by creating a new S-node and P-node, as shown in Figure 7. This completes the construction of T_{j+1} .

Case 4. Suppose Cases 1, 2 and 3 do not apply. That is, there are sons, say s_1, s_2 , such that $PAR(s_1) = PAR(s_2) = 0$, and $SER(s_2) = 0$. In this case it is possible to identify the forbidden subgraph in the transitive closure of the subgraph induced on nodes $1, 2, \dots, j+1$, as shown in Figure 8.

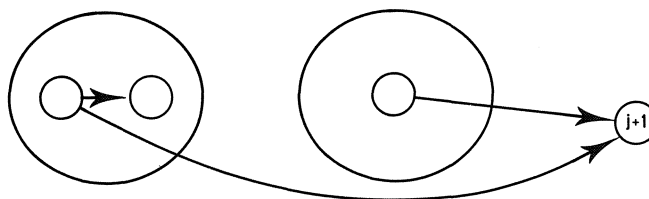


Insertion of node $j+1$ in Case 3 of P-node Procedure.

Figure 7.

$SER(s_2) = 0$
 $PAR(s_2) = 0$

$PAR(s_1) = 0$



Forbidden Subgraph Identified in Case 4 of P-node Procedure

Figure 8.

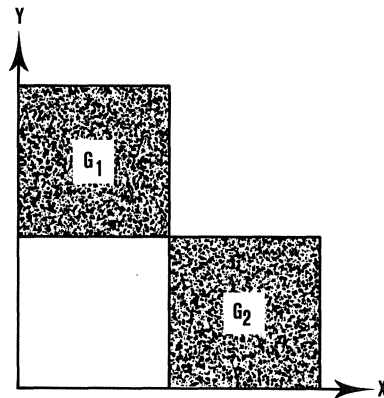
It should be apparent that by starting at the root node and proceeding toward the leaves, one either obtains T_{j+1} from T_j in $O(n)$ running time, or shows that the digraph G is nonseries parallel. Thus, the overall running time required for recognition of series parallel digraphs is $O(n^2)$.

"Two Dimensionality" of Transitive Series Parallel Digraphs

Every partial order on n elements can be represented as the intersection of total orders. The minimum number of total orders which are necessary to so represent a partial order is called its "dimension".

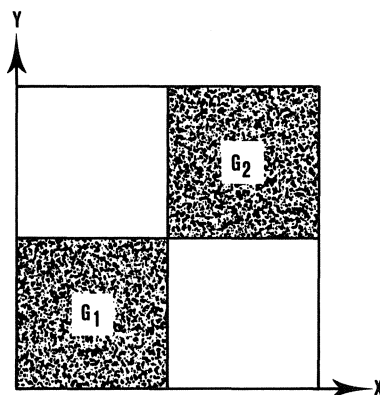
Every transitive series parallel digraph represents a partial order of dimension 2, but not conversely. To capture the notion of two-dimensionality geometrically, imagine that the n elements of the partial order " \leq " are identified with points in the plane. Let element i have coordinates (x_i, y_i) and element j have coordinates (x_j, y_j) . Then $i \leq j$, if and only if $x_i \leq x_j$ and $y_i \leq y_j$. (Note that antisymmetry requires that either $x_i \neq x_j$ or $y_i \neq y_j$.)

Thus an element is " \leq " a second element if and only if it is located below and to the left of the second element. Note that parallel composition can be effected by locating G_1, G_2 as shown in Figure 9 and series composition as shown in Figure 10.



Parallel Composition

Figure 9



Series Composition

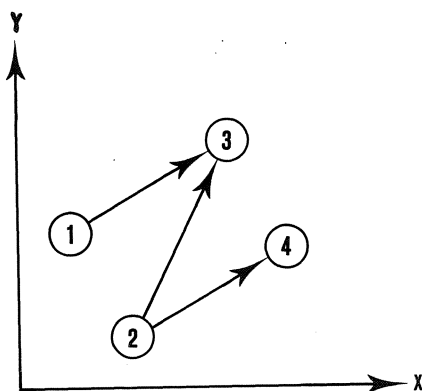
Figure 10

One can obtain two total orders whose intersection represents the transitive closure of a given series parallel digraph by simple operations performed on the decomposition tree. For each leaf node j , the two total orders are simply $\alpha = \beta = (j)$. For an S-node whose sons are s_1, s_2, \dots, s_k , from right to left, with total orders $\alpha_1, \beta_1; \alpha_2, \beta_2; \dots; \alpha_k, \beta_k$, the total orders are $\alpha = \alpha_k \alpha_{k-1} \dots \alpha_1$, $\beta = \beta_k \beta_{k-1} \dots \beta_1$. For a P-node, the total orders are $\alpha = \alpha_k \alpha_{k-1}, \dots, \alpha_1$ and $\beta = \beta_1 \beta_2 \dots \beta_k$.

It should be evident that these total orders can be computed in $O(n)$ time, given the decomposition tree of the series parallel digraph.

Note that not every two-dimensional partial order is transitive series parallel. This can be confirmed from the fact that the forbidden subgraph can be so represented, as shown in Figure 11. Please note that one representation of the forbidden subgraph is $\alpha = (1, 2, 3, 4)$, $\beta = (2, 4, 1, 3)$. One can always renumber the nodes so that one of the partial orders is $\alpha = (1, 2, \dots, n)$. Then all information about the partial order can be obtained from β . In particular, one can read off the complete adjacency matrix in $O(n^2)$ time.

The preceding observations show that one can compute the transitive closure of a series



Forbidden Subgraph is Two-Dimensional

Figure 11

parallel digraph in time $O(n^2)$. This contrasts with a time of $O(n^{2.81})$ for the best available transitive closure algorithm for general digraphs.

4. ISOMORPHISM OF TRANSITIVE SERIES PARALLEL DIGRAPHS

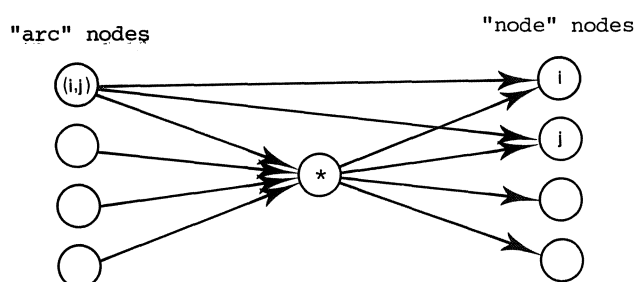
The status of the graphic isomorphism problem is unclear. Specifically, we know that the problem is in NP, but we do not know if it is either polynomial-bounded or NP-complete. We do know that a number of problems are polynomial equivalent to the graphic isomorphism problem: isomorphism of semigroups, digraphs, bipartite graphs, chordal graphs, and series parallel digraphs.

We present a simple reduction of the (undirected) graph isomorphism problem to the problem of isomorphism for series parallel digraphs.

For a given graph $G = (N, A)$, construct a digraph $G' = (N', A')$ with $m+n+1$ nodes, as follows. $N' = N \cup A \cup \{*\}$, where "*" is a special node. There is an arc $((i, j), *) \in A'$ from each "arc" node of G' to $*$, and an arc $(*, j)$ from $*$ to each "node" node of G' . There are also arcs $((i, j), i)$, $((i, j), j)$ from each "arc" node (i, j) to the "node" nodes i and j . The resulting

digraph G' , indicated schematically in Figure 12, is clearly series parallel.

We also assert, without formal proof, that two graphs G_1, G_2 are isomorphic if and only if G'_1, G'_2 , constructed as above, are isomorphic. It therefore follows



Series Parallel Digraph in Isomorphism Reduction

Figure 12

that we cannot expect to solve the isomorphism problem for series parallel digraphs without solving the general graphic isomorphism problem as well.

What we shall do is to solve the isomorphism problem for transitive series parallel digraphs. Or to put it another way, we shall be able to test two series parallel digraphs G_1, G_2 , for isomorphism of their transitive closures.

It is apparent that the decomposition tree for a series parallel digraph is also a decomposition tree for its transitive closure. We also noted, but shall not provide a formal proof, that the decomposition tree of a transitive series parallel digraph is uniquely determined, up to a reordering of the sons of the P -nodes. It follows that the isomorphism problem for transitive series parallel digraphs reduces to an isomorphism problem for decomposition trees.

Suppose T_1 and T_2 are the decomposition trees for the transitive series parallel digraphs G_1, G_2 . These graphs can be isomorphic only if both T_1, T_2 have an S -node or a P -node as roots. It follows that at each level, T_1, T_2 each

contain only S-nodes and leaves or only P-nodes and leaves.

We shall test for isomorphism of T_1, T_2 by working from level to level, from the leaves toward the roots. We shall speak of a node of T_1 as being "isomorphic" to a node of T_2 if the subgraphs induced on their descendants are isomorphic. Clearly each leaf of T_1 is isomorphic to each leaf of T_2 , since all single-node digraphs are isomorphic.

Now suppose we have determined isomorphism of all pairs of nodes at a P-level and we wish to determine isomorphism for nodes at the S-level immediately above. Suppose node s_1 of T^1 has sons p_1, p_2, \dots, p_k , from right to left and node s_2 of T^2 has sons p'_1, p'_2, \dots, p'_k . Clearly s_1 can be isomorphic to s_2 only if $k = k'$. Assuming this is the case, s_1 is isomorphic to s_2 if and only if $p_1, p'_1; p_2, p'_2; \dots; p_k, p'_k$ are pairwise isomorphic. It follows that it is easy to determine isomorphism of S-nodes from isomorphism of their sons.

Now suppose we have determined isomorphism of all pairs of nodes at an S-level and we wish to determine isomorphism for nodes at the P-level immediately above. Let node p_1 of T^1 have sons s_1, s_2, \dots, s_k and node p_2 of T^2 have sons s'_1, s'_2, \dots, s'_k . Again, clearly p_1 can be isomorphic to p_2 only if $k = k'$. Assuming this is true, p_1 is isomorphic to p_2 if and only if it is possible to establish a one-to-one correspondence (a complete matching) between isomorphic sons of p_1 and p_2 . Note that isomorphism is an equivalence relation. Hence if $s_1 = s'_1$ and $s'_1 = s_2$ then $s_1 = s_2$. A necessary and sufficient condition for a complete matching to exist is that for each s_i , s_i is isomorphic to q sons s'_j , $q \geq 1$ and that each of these sons be isomorphic to exactly q sons s_i . It follows that if a complete matching exists, it can be found by a simple algorithm, as follows:

- Step 0 Set $i = 0$. No nodes are matched.
- Step 1 If $i = k$, then p_1 and p_2 are isomorphic. Otherwise, set $i = i+1$.
- Step 2 Try to find an unmatched node s'_j , isomorphic to s_i . If no such node exists, stop; p_1 and p_2 are not isomorphic. If such a node s'_j can be found, match s_i with s'_j , and return to step 1.

The two transitive series parallel digraphs are isomorphic if and only if their root nodes are isomorphic. It is evident that this can be determined

in $O(n^2)$ time.

5. SUBGRAPH ISOMORPHISM

The *subgraph isomorphism problem* is to determine if a given graph G_1 is isomorphic to a subgraph of a second graph G_2 . It is well-known that the subgraph isomorphism problem is NP-complete. This result follows immediately from the fact that the maximum independent set problem is NP-complete, since this problem is equivalent to determining if a given complete graph is isomorphic to a subgraph of a second graph.

The subgraph isomorphism problem can, however, be solved efficiently for transitive series parallel digraphs.

For simplicity, let us suppose that the decomposition trees for G_1 , G_2 have the same type of root node (S or P). If this is not the case, then the problem becomes that of determining if the root for G_1 is isomorphic to a son of the root node for G_2 . (Note that we shall continue speaking of "isomorphism" of nodes, when what we mean is "isomorphism of the subgraph induced by the descendants of the node in T_1 to a subgraph of the subgraph induced by the descendants of the node in T_2 ".)

We must now make some modifications of the rules stated in the previous section.

As before, suppose we have determined isomorphism for all pairs of nodes at a P-level and we wish to determine isomorphism for nodes at the S-level immediately above. Suppose node s_1 of T_1 has sons p_1, p_2, \dots, p_k , from right to left, and node s_2 of T_2 has sons $p'_1, p'_2, \dots, p'_{k'}$, where $k' \geq k$. We find the rightmost son, $p'_{j(1)}$ to which p_1 is isomorphic, then the rightmost son, $p'_{j(2)}$ to the left of $p'_{j(1)}$ to which p_2 is isomorphic. The nodes s_1, s_2 are isomorphic if and only if all the nodes p_1, p_2, \dots, p_k can be matched up in this way.

Now suppose we have determined isomorphism of all pairs of nodes at an S-level and we wish to determine isomorphism for nodes at the P-level immediately above.

Suppose node p_1 of T_1 has sons s_1, s_2, \dots, s_k and node p_2 of T_2 has sons $s'_1, s'_2, \dots, s'_{k'}$, where $k' \geq k$. Node s_1 is isomorphic to s_2 if and only if the nodes s_1, s_2, \dots, s_k can be matched to distinct nodes from among $s'_1, s'_2, \dots, s'_{k'}$.

This is a bipartite matching problem which can be solved in $O(k')^{5/2}$ time [2]. Let k_j denote the number of sons of the j th P-node at level i of the T_2 . Then

$$\sum_j k_j^{5/2} \leq \left(\sum_j k_j \right)^{5/2} = n_{i+1}^{5/2},$$

where n_{i+1} is the number of nodes at level $i+1$. But

$$\begin{aligned} \sum_i n_i &< 2n, \\ \sum_i n_i^{5/2} &\leq \left(\sum_i n_i \right)^{5/2} < (2n)^{5/2}. \end{aligned}$$

It follows that the overall running time required for solving all matching problems is $O(n^{5/2})$. All other operations are dominated by matching, hence subgraph isomorphism can be tested in $O(n^{5/2})$ time.

Essentially the same approach was used by Edmonds and by Matula [8] for the problem of determining whether one tree is isomorphic to a subtree of another. The present procedure can be viewed as an extension of their techniques to a broader class of isomorphism problems.

6. FINDING MINIMUM SPANNING TREES

Let $G = (N, A)$ be a connected undirected graph, with a value $c(v, w)$ specified for each arc $(v, w) \in A$. We wish to find a spanning tree $T = (N, A')$, such that

$$\sum_{(v, w) \in A'} c(v, w).$$

This problem has been solved in time $O(n^2)$ or $O(m \log n)$. Recently, a variety of algorithms have been suggested, all with $O(m \log \log n)$ running time [2, 14]. In the present section we describe what appears to be the simplest of all known algorithms with $O(m \log \log n)$ running time. The presentation is adapted from Cheriton and Tarjan [2].

LEMMA 1. Let $X \subseteq N$. Let $(v, w) \in A$ be such that $v \in X$, $w \in N - X$, and $c(v, w) = \min\{c(x, y) \mid (x, y) \in A, x \in X, y \in N - X\}$. Then some minimum spanning tree contains (v, w) .

This lemma justifies the following very general procedure:

First step: Pick some node. Choose the smallest arc incident to the node. This is the first arc of the spanning tree.

General step: The arcs so far selected define a forest (a graph consisting of a set of trees) which is a subgraph of G . (Isolated nodes, not incident to any arcs selected so far, are regarded as trees with one node.) Pick a tree in the forest. Select the smallest unselected arc (v,w) incident to a node in T and to a node in a *different* tree T' . Delete all arcs smaller than (v,w) which are incident to T (these arcs form cycles in T). Add (v,w) to the minimum spanning tree (updating the forest by combining T and T'). Repeat the general step until all nodes are connected.

This general method requires a mechanism for keeping track of the subtrees in the forest. To keep track of the nodes in each tree of the forest, we propose to use the well-known "disjoint set union" algorithm [13]. Given a collection of disjoint sets (in this case sets of the nodes in each tree) the set union algorithm implements three operations:

FIND (x) : returns the name of the set containing node x .
 UNION (i,j): combines sets i and j , naming the new set i .
 INIT (i,L) : initializes set i to contain all nodes in list L .

The worst-case time required for $O(m)$ FINDs and $O(n)$ UNIONs is $O(m\alpha(m,n))$, where α is a functional inverse of Ackermann's function [13]. Initialization requires $O(n)$ time for n trees, each consisting of a single node.

To keep track of the arcs incident to each tree of the forest, we use mechanisms of priority queues. We maintain queues of all arcs incident to each tree of the forest (except, of course, those arcs which have been deleted). We need three operations on these queues:

QUNION (i,j): combines queues i and j , naming the new queue i .
 MIN (i) : returns the smallest arc in queue i which has an endpoint outside tree i . MIN (i) also deletes this arc and all smaller arcs from queue i . (Note: MIN (i) must use the FIND operation to test whether a given arc has an endpoint outside tree i .)
 QINIT (i,L) : initializes queue i to contain all arcs in list L .

An implementation of the general algorithm using these operations as primitives appears below. We assume the graph has node set $N = \{1,2,\dots,n\}$ and

that for each node v , $I(v) = \{(v,w) \mid (v,w) \in A\}$ is a list of the arcs incident to V . (Note that (v,w) is considered an ordered pair $(v,w) \in I(v)$, $(w,v) \in I(w)$.) Each tree i of the forest is represented by a set i of its nodes and by a queue i of its incident arcs. Initially each node v is represented by the set $\{v\}$ and by a queue consisting of the arcs in $I(v)$. At all times during execution of the algorithm, every arc (v,w) in queue k satisfies $\text{FIND}(v) = k$.

General Algorithm

```

For i = 1 until n do
  begin
    INIT(i,{i});
    QINIT(i,I(i));
  end
While more than one tree do
  begin
    pick some tree K;
    (i,j) ← MIN(K);
    add arc (i,j) to spanning tree;
    x ← FIND(j);
    UNION(x,k);
    QUNION(x,k)
  end;

```

We shall consider two methods for implementing priority queues:

(a) unordered lists and (b) ordered lists of size k .

(a) Unordered Lists

QINIT (i,L) requires $O(m(i))$ time where $m(i)$ is the size of queue i (and of L).

MIN (i,L) requires $O(m(i))$ time plus time for $m(i)$ FINDs.

QUNION (i,j) requires $O(1)$ time.

(b) Ordered Lists of Size k

Another possibility is to represent a queue as a list of sets, each set initially of size k , plus possibly one set of size less than k . Each set is in sorted order, but there is no ordering relationship among the sets. We

carry out QINIT (i,L) by dividing L into $\lfloor |L|/k \rfloor$ sets of size k plus at most one set of size less than k. We then sort these sets. This process requires $O(|L| \log k)$ time. MIN (i) is carried out by inspecting the arcs in each set of queue i in order, discarding those that don't connect two different trees. We then compare the smallest arc from each set and select the overall minimum arc. MIN (i) requires $O(s(i) + d(i))$ time, plus time for $s(i) + d(i) - 1$ FINDs, where $s(i)$ is the number of sorted sets in queue i and $d(i)$ is the number of arcs deleted from queue i by MIN (i).

(Note: $d(i) \geq 1$, always, because the arc returned by MIN (i) is deleted.)

UNION (i,j) is accomplished by simply merging the list of sets in queue i and the list of sets in queue j. This requires $O(1)$ time.

We shall use the "uniform" selection rule to determine the order in which trees are processed. This works as follows. Initially all of the n trees (each a single node) arc placed on a queue. At the general step, the first tree T at the front of the queue is examined. When this tree T is connected to another tree T', both T and T' are deleted from the queue and the new combined tree is placed at the rear of the queue. It is easy to implement this selection rule so that the total overhead for selection is $O(n)$.

Without loss of generality, assume that the nodes of the graph are numbered, so that the i^{th} execution of MIN is MIN (i). Let $m(i)$ be the number arcs in queue i then MIN (i) is executed. For each i, $1 \leq i \leq n-1$, let T_i be the tree in F containing node i when MIN (i) is executed. T_n is the minimum spanning tree constructed.

We define a *stage number*, $st(T)$ for T, by $st(T) = 0$, if T contains a single node, and $st(T) = \min\{st(T'), st(T'')\} + 1$, if T is formed by connecting trees T' and T'' by an arc. It is easy to prove by induction that if $st(T) = j$, then T has at least 2^j vertices, so there are at most $\lfloor \log n \rfloor$ stages at which trees are connected.

LEMMA 2. *Two different trees T_i and T_j with the same stage number are node disjoint.*

Proof. Suppose T_i and T_j share a node. Then either $T_i \subseteq T_j$ or $T_j \subseteq T_i$. Without loss of generality, assume $T_i \subseteq T_j$. The trees initially on the queue have nondecreasing stage numbers (all zero), and the property is preserved as the algorithm proceeds. Just before MIN (i) is executed, T_i has a stage number as small as any tree on the queue. Thus, if $T_i \neq T_j$, all

trees T on the queue which are subtrees of T_j have $\text{st}(T) \geq \text{st}(T_i)$, which means $\text{st}(T_j) \geq \text{st}(T_i) + 1$. \square

COROLLARY 3. For any stage s , $0 \leq s < \log n$,

$$\sum_{\text{st}(T_i)=s} m(i) \leq 2m.$$

Moreover,

$$\sum_{i=1}^{n-1} m(i) \leq 2m \log n.$$

Proof. Each arc is represented at most twice in the queues, and $0 \leq \text{st}(T_i) \leq \log n - 1$, if $1 \leq i \leq n-1$, by the proof of Lemma 2 and the fact that $|T_i| \geq 2^{\text{st}(T_i)}$. \square

We now propose to implement the algorithm as follows:

Phase 1: Initialize all queues as unordered lists, and execute the algorithm until stage $\log \log n$.

Phase 2: Re-initialize all queues as ordered lists of size $k = \log n$, and execute the algorithm until completion.

The time required for the operations QINIT and MIN in the table below. Note that the total number of QUNIONS performed is exactly $n - 1$, with $O(1)$ time for each QUNION. The time required for FINDS and UNION is dominated by $O(m \log \log n)$, which is seen to be the running time for the overall algorithm.

	QINIT	MIN
Phase 1	$O(m)$	$O(m \log \log n)$
Phase 2	$O(m \log \log n)$	$O(m)$

Note that if unordered lists were used for the entire computation, the time for MIN would be $O(m \log n)$. If ordered lists of size $\log n$ were used for the entire computation something like $\frac{m}{\log n} + n$ sets could be created at initialization, and the time for MIN would then contain a term of the form $n \log n$. Thus two phases are necessary to attain the desired running time of $O(m \log \log n)$.

7. GENERATING ALL MAXIMAL INDEPENDENT SETS

An *independent set* of nodes of a graph $G = (N, A)$ is a subset $I \subseteq N$ such that no two nodes in I are adjacent. There are situations, e.g. in graph coloring, in which one may want to generate all the *maximal* independent sets of a given graph. This can be accomplished very efficiently by an algorithm due to Tsukiyama, et al [16].

Note that an independent set in G can be identified with a complete subgraph in the complement of G . Also, the complement of an independent set is a covering of arcs by nodes. Thus any algorithm for generating maximal independent sets can be used to generate maximal cliques or minimal covers as well.

Before describing the algorithm of the four Japanese, let us review an older algorithm due to Paull and Unger [11].

Let I_j denote the family of all independent sets which are maximal in the subgraph induced on nodes $1, 2, \dots, j$. Clearly $I_1 = \{\{1\}\}$. We then propose to find I_{j+1} from I_j , finally obtaining I_n , the family of maximal independent sets of G itself.

Let $I \in I_j$, and let $A(j+1)$ denote the set of nodes adjacent to node $j+1$. If $I \cap A(j+1) = \emptyset$, it is clear that $I \cup \{j+1\} \in I_{j+1}$. And if $I \cap A(j+1) \neq \emptyset$, then $I \in I_{j+1}$. We thus deduce that

$$|I_1| \leq \dots \leq |I_j| \leq |I_{j+1}| \leq \dots \leq |I_n| = K.$$

Now consider a typical set $I' \in I_{j+1}$. If $j+1 \notin I'$, then $I' \in I_j$. If $j+1 \in I'$, then $I' - \{j+1\} \subseteq I$, for some $I \in I_j$ and $I' = (I - A(j+1)) \cup \{j+1\}$. It thus follows that if we form the multiset

$$I'_{j+1} = \{I' \mid I' = (I - A(j+1)) \cup \{j+1\}, I \in I_j\},$$

then

$$I_{j+1} \subseteq I_j \cup I'_{j+1}.$$

It is very important to note that we refer to I'_{j+1} as a *multiset*. That is, I'_{j+1} may contain duplicate entries. A given set $I' \in I'_{j+1}$ may result from two distinct sets $I_1, I_2 \in I_j$. That is,

$$(6.1) \quad I' - \{j+1\} = I_1 - A(j+1) = I_2 - A(j+1),$$

even though $I_1 \neq I_2$.

In order to obtain I_{j+1} from $I_j \cup I'_{j+1}$, it is necessary to purge from $I_j \cup I'_{j+1}$ those sets which are either nonmaximal or duplicates of other entries. Since $|I_j| = |I'_{j+1}| \leq k$, this can be done with at most $O(k^2)$ pairwise comparisons of sets.

Now suppose $A(j+1)$ is given by a list and each $I \in I_j$ is specified by an n -vector. Then $O(|A(j+1)|K)$ time suffices to form the multiset I'_{j+1} . Each set comparison requires $O(n)$ time, so purging $I_j \cup I'_{j+1}$ of nonmaximal and duplicate entries can be done in $O(nK^2)$ time. It is clear that the overall running time of the procedure is bounded by $O(n^2K^2)$ and storage $O(m+nK)$.

What the four Japanese have done is to reduce the running time to $O(mnK)$ and space requirements to $O(m+n)$. Their key contribution is to show how to form I'_{j+1} without introducing either nonmaximal or duplicate sets.

The avoidance of nonmaximal sets is very simple. Let $I \in I_j$. Then

$$I' = (I - A(j+1)) \cup \{j+1\}$$

is a maximal independent set on the subgraph induced on $1, 2, \dots, j+1$, if and only if, for all $k < j+1$, $k \notin I$,

$$(6.2) \quad A(k) \cap I' \neq \phi.$$

This test can be performed for all $I' \in I'_{j+1}$ in $O((n+m)K)$ time.

Now let us consider the matter of duplication. Suppose a situation were to exist, as in (6.1). The rule we shall use is the following. For a given set $I \in I_j$, we shall introduce $I' = (I - A(j+1)) \cup \{j+1\}$ into I'_{j+1} only if $I \cap A(j+1)$ is *lexicographically smallest*, with respect to all sets $I \in I_j$ for which $I - A(j+1)$ is the same.

Recall that lexicography induces a total order. Let $I_1 \neq I_2$, where

$$I_1 = \{e_{j(1)}, e_{j(2)}, \dots, e_{j(p)}\},$$

$$I_2 = \{e_{j'(1)}, e_{j'(2)}, \dots, e_{j'(q)}\},$$

$j(1) < j(2) < \dots < j(p)$; $j'(1) < j'(2) < \dots < j'(q)$, and $p \leq q$. Suppose there is an $k \leq p$, such that $j(i) = j'(i)$, for $i = 1, 2, \dots, k-1$, and $j'(k) < j(k)$. Then we say that I_2 is lexicographically smaller than I_1 ; otherwise I_2 is lexicographically larger.

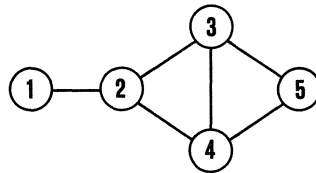
In words, a set $I \in I_j$ satisfies the lexicography condition if and only if there does not exist a node $k < j+1$, $k \notin I$ which is nonadjacent to both $I - A(j+1)$ and to all lower-numbered nodes in $I \cap A(j+1)$. In other words, if and only if, for all $k < j+1$, $k \notin I$,

$$A(k) \cap (I - (A(j+1) \cap \{k+1, k+2, \dots, j\})) \neq \emptyset.$$

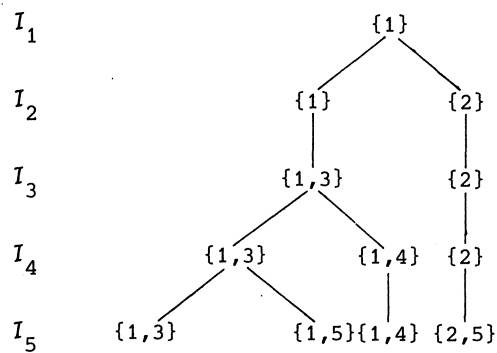
This test can also be performed for all $I \in I_j$ in $O((n+m)K)$ time.

It is now apparent that I_{j+1} can be formed from I_j in $O((n+m)K)$ running time. Hence I_n can be obtained in $O(n(n+m)K)$ time overall. Assume $m \geq n$, and let us call this time $O(mnK)$. (If $m < n-1$, then there is at least one node with degree zero, and any such node is contained in every maximal independent set.)

Let us now see how storage requirements can be reduced. The process of generating the maximal independent sets can be visualized by means of a tree, as shown in Figure 14, for the graph in Figure 13.



Graph for Example
Figure 13



Generation of Maximal Independent Sets
for Example

Figure 14

The nodes at level j of the tree are identified with sets in I_j , with the tree rooted to $\{1\}$, the unique set in I_1 . Each node has at least one son. (Recall that for each $I \in I_j$, either I or $I \cup \{j+1\}$ is in I_{j+1}). The left son, if it exists, is identified with $(I - A(j+1)) \cup \{j+1\}$. This son is provided only if the maximality and lexicography conditions are satisfied. The right son, if it exists, is identified with I , but this son exists only if $I \cap A(j+1) \neq \emptyset$. (otherwise I is not maximal in I_{j+1} , being a subset of $I \cup \{j+1\} \in I_{j+1}$.)

We now carry out our computation in the form of a depth-first search of the tree, always traveling downward to the left-most son of a node that we have not yet visited. When we reach a node in I_n , we simply output the maximal set in question (rather than maintaining it in storage). We then backtrack. The backtracking procedure is facilitated by maintaining a stack. We leave it as an exercise for the reader to show that the stack can be implemented in $O(n)$ space.

8. FINDING A MAXIMUM INDEPENDENT SET

Let us now consider the problem of finding a *maximum* (i.e. maximum cardinality) independent set in a graph. As we know, this is an NP-complete problem, so we can hardly hope to find such a set in running time which is polynomial in m and n .

It is a well-known result that there are at most $3^{n/3}$ maximal independent sets in a graph [10]. It thus follows that the algorithm of the previous section can be applied to solve the problem in $O(mn(3^{1/3})^n)$ time. Note: $3^{1/3} \approx 1.445$.

Tarjan and Trojanowski [14] have shown that the problem can be solved in $O(2^{n/3})$ time. We shall not go into the complete elaboration of their procedure, but only indicate the central idea behind their procedure.

First note that the problem is trivial if the maximum degree of any node in G does not exceed 2. In this case, each component consists of a chain or a cycle, and it is easy to pick out a maximum number of nodes from each component to be part of an independent set.

Hence let us assume that there is at least one node j of degree three or more. Obviously a maximum independent set either (a) contains j or (b) it does not. If we make the former assumption, then none of the nodes in $A(j)$ can be contained in the solution. It then suffices to find an inde-

pendent set in the graph from which j and the nodes in $A(j)$ are deleted. In the latter case, we are left with the problem of finding the best solution we can in the subgraph obtained by deleting node j only. The first subgraph has at most $n-4$ nodes, and the latter $n-1$.

Let $c(n)$ = the running time required to solve the problem for an n -node graph, by our procedure. Then we have

$$c(n) \leq c(n-1) + c(n-4) + p(n),$$

where $p(n)$ is a polynomial function required for housekeeping, comparing solutions for the two subproblems, etc.

Assuming $c(n) = a^n$, for some base a , an upper bound on the order of $c(n)$ is given by solving for a such that

$$a^n = a^{n-1} + a^{n-4}$$

or

$$a^4 = a^3 + 1.$$

This indicates $a \approx 1.38$, which, of course, is better than $3^{1/3} \approx 1.445$.

Tarjan and Trojanowski succeeded in reducing the base a to less than $2^{1/3}$ by a much more involved case analysis.

9. COMPUTING THE CHROMATIC NUMBER OF A GRAPH

A k -coloring of a graph $G = (N, A)$ is a partition of the node set N into k independent sets. If there exists such a partition, G is said to be k -colorable. The *chromatic number* of G is the least value of k such that G is k -colorable.

It is well known that the chromatic number problem is NP-complete. In fact, even the problem of determining the chromatic number of a graph to within a given factor $r < 2$ has been shown to be NP-complete [3].

A number of algorithms for this problem have been proposed [9], yet there appears to have been very little discussion in the literature concerning nontrivial upper bounds on the complexity of the problem. We propose to make use of the results obtained above concerning maximal in-

dependent sets to obtain an algorithm with worst case running time of $O(mn(1+\sqrt[3]{3})^n)$. Note: $1 + \sqrt[3]{3} \approx 2.445$.

It is easy to see that if a graph is k -colorable, at least one of the k colors may be assumed to be a maximal independent set. For a given graph $G = (N, A)$, let $K(N')$ denote the chromatic number of the subgraph induced on $N' \subseteq N$. It follows that if N' is nonempty there exists a maximal independent set I of the subgraph induced on N' , such that

$$K(N') = 1 + K(N'-I).$$

There are a finite number of maximal independent sets I of the induced subgraph. By minimizing over them we obtain the "dynamic programming" equations

$$K(N') = 1 + \min_{I \subseteq N'} \{K(N'-I)\}, \quad N' \neq \emptyset$$

$$K(\emptyset) = 0.$$

We now estimate the running time required to solve these equations for all $N' \subseteq N$, in the worst case.

Suppose, for fixed N' , we have already found $K(N'')$, for all proper subsets $N'' \subset N'$. The time required to compute $K(N')$ is then proportional to the number of maximal independent sets of the subgraph induced on N' , plus the time required to generate them. But, as we have seen, this time is no worse than $O(mr3^{r/3})$, where $|N'| = r$. Summing over all possible $N' \subseteq N$, and invoking the binomial theorem, we obtain

$$\sum_{r=0}^n \binom{n}{r} mr3^{r/3} < mn \sum_{r=0}^n \binom{n}{r} 3^{r/3} = mn(1+\sqrt[3]{3})^n.$$

This yields the desired result.

REFERENCES

- [1] A. AHO, J.E. HOPCROFT & J. ULLMAN, *The design and analysis of computer algorithms*, Addison-Wesley, Reading, Mass., 1974.
- [2] D. CHERITON & R.E. TARJAN, *Finding minimum spanning trees*, to appear in *Siam J. Comput.*
- [3] M.R. GAREY & D.S. JOHNSON, *The complexity of near-optimal graph coloring*, *Journal ACM* 23 (1976) 43-49.
- [4] J.E. HOPCROFT & R.M. KARP, *A $n^{5/2}$ algorithm for maximum matchings in bipartite graphs*, *SIAM J. Comput.* 2 (1973) 225-231.
- [5] E.L. LAWLER, *Sequencing jobs to minimize total weighted completion time subject to precedence constraints*, to appear.
- [6] E.L. LAWLER, *A note on the complexity of the chromatic number problem*, *Information Processing Letters*, 5 (1976) 66-67.
- [7] E.L. LAWLER, R.E. TARJAN & J. VALDEZ, *Analysis and isomorphism of series parallel digraphs*, to appear.
- [8] D.W. MATULA, *Subtree isomorphism in $O(n^{5/2})$* , to appear.
- [9] D.W. MATULA, G. MARBLE & J.D. ISAACSON, *Graph coloring algorithms*, *Graph Theory and Computing*, R.C. Read, ed., Academic Press, New York, 1972, 109-122.
- [10] J.W. MOON & L. MOSER, *On cliques in graphs*, *Israel J. Math.* (1965) 23-28.
- [11] M.C. PAUL & S.H. UNGER, *Minimizing the number of states in incompletely specified sequential functions*, *IRE Trans. Electr. Computers* EC-8 (1959) 356-357.
- [12] R.L. RIVEST & J. VUILLEMIN, *A generalization and proof of the Aanderaa-Rosenberg conjecture*, *Proc. Seventh Annual ACM Symposium on Theory of Computing*, May 1975.
- [13] R.E. TARJAN, *Efficiency of a good but not linear set union algorithm*, *J. ACM* 22 (1975) 215-225.
- [14] R.E. TARJAN & R. TROJANOWSKI, *Finding a maximum independent set*, *Technical Report CS-76-550*, Computer Science Dept., Stanford University.

- [15] A.C. YAO, *An $O(|E| \log \log |V|)$ algorithm for finding minimum spanning trees*, Information Processing Letters, 4 (1975) 21-23.
- [16] S. TSUKIYAMA, M. IDE, H. ARUJOSHI & H. OZAKI, *A new algorithm for generating all the maximal independent sets*, to appear.

Note added in proof.

It appears that it is possible to recognize series parallel digraphs in $O(m)$ time. This newer procedure will be communicated in [7].

0. Introduction	37
1. Efficiency versus data-representation.	39
2. File-merging	51
3. Tables and balanced trees.	67
4. Path compression	84
5. Associative search structures.	102
6. Pattern matching	126
References.	142

THE COMPLEXITY OF DATA ORGANIZATION

J. VAN LEEUWEN *

State University of New York, Buffalo, USA

0. INTRODUCTION

0.1. The way *data* is presented to the computer and represented in the system can be a dominating factor in the performance of the computer as a data-processing tool.

By its very nature the area of data-organisation extends from the study of combinatorial algorithms to the design of information management systems, and it is too broad to be conveyed in a series of six lectures only.

We decided to emphasize a few topics in concrete complexity theory which have recently received much attention and which have lead to a number of fundamental techniques and new results which may be applicable to many problems which computer-programmers encounter.

0.2. The present lectures neither emphasize (say) a specific domain in computational complexity nor give you the design-philosophy of relational database, but instead we will concentrate on useful *techniques* (tricks?) in data-organisation which can bring improvements in many computer-programs performing data-manipulations of some sort.

We have specifically avoided to present results which are already adequately treated in such excellent texts as Knuth [43],[44] or Aho, Hopcroft, and Ullman [3].

The reader is assumed to have some familiarity with computer-programming and the fundamental information-structures as discussed in Knuth [43] or Wirth [66]. Roberts [54] gives a useful survey of the basic file-organisation techniques.

*) Present address: Dept. of Computer Science, Whitmore Laboratory, the Pennsylvania State University, University Park, Pennsylvania 16802.

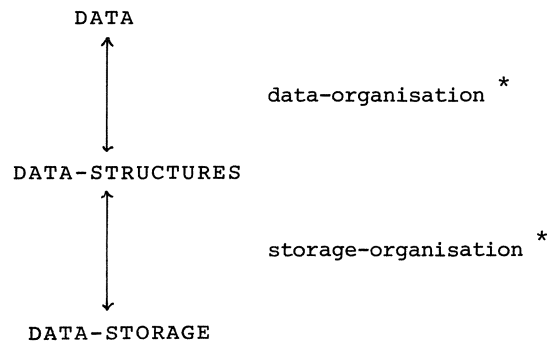


fig.1. Some terminology

* Sometimes called "file-organisation" and "data-organisation" respectively.

1. EFFICIENCY VERSUS DATA-REPRESENTATION

1.1. Let us start from an intuitive concept of *data*. The user is largely responsible for *data-collection*, and must make sure that he gathers the information needed for a successful data-processing system.

1.2. To facilitate *retrieval* and efficient *manipulation* information should not be stored randomly, but is preferably organised in structures which allow for easy access and (say) deletion and addition of data at all times.

1.3. Computer storage is divided into directly addressable *core-storage* ("main memory") and "supplementary" or *secondary storage* ("auxiliary memory"). Data is usually stored in auxiliary memory on mass-storage devices like magnetic tapes, magnetic drums, (magnetic) cards, data-cells, or disc-packs, and only small parts will be held for processing in main memory at a time. In time-shared systems the data will necessarily be segmented, with due restrictions on size per segment.

1.4. The user is not likely to be concerned with the "hardware" of a data management system, but is using an intermediate language instead. This can be a "host"-language (i.e. some general purpose programming language) providing helpful primitive data-structures and a flexible mode-definition mechanism for creating more complex structures, or a special *data-language*.

Data-structures are merely the model of data-storage as observed by the user. The "system" presumably interpretes his model in "real" storage.

1.5. Information is usually provided in small packages of logically connected data called *records*. The *fields* of a record must be specified, and may consist of other records. The information in a record can be accessed by using the *field names* as *selectors*.

1.6. Some fields may be used (or added) for uniquely identifying a record. The contents of such fields together form the *primary key* of a record.

Later we shall simply identify a record and its primary key (as it is the only part our algorithms will be using), and define it to be an *ordered k-tuple* (b_1, b_2, \dots, b_k) of values taken from some ordered set R.

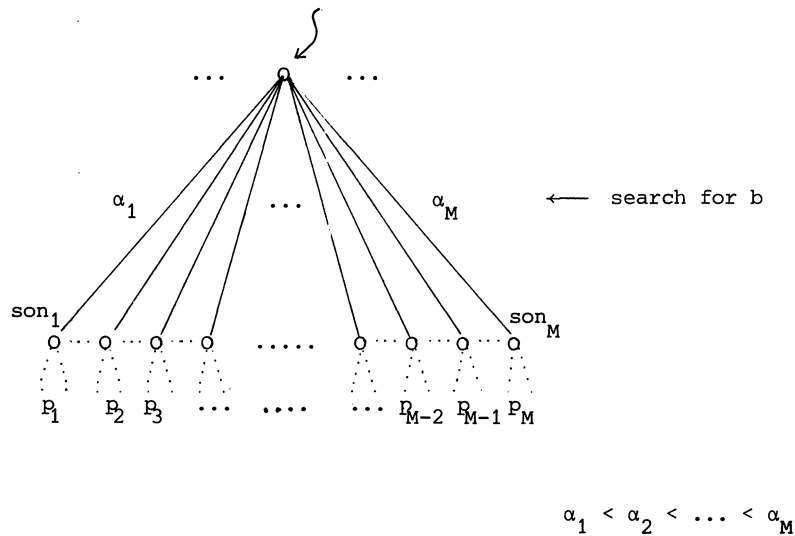
1.7. Consider an arbitrary, unstructured collection of (perhaps) randomly stored records

$$V = \{v_1, v_2, \dots, v_N\} \subseteq R^k.$$

Numerous *address-calculation schemes* exist for "finding" a record when its primary key is presented (Knuth [44]).

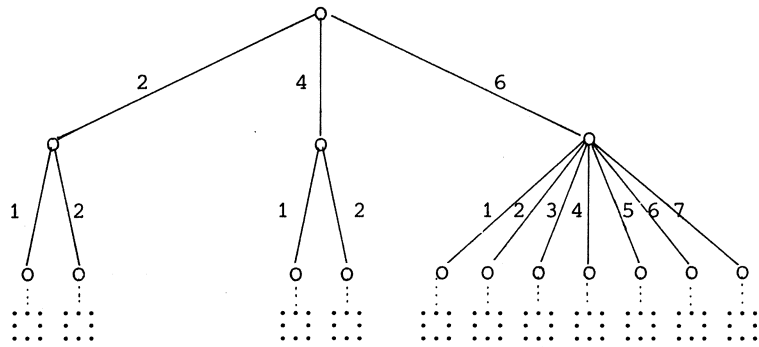
1.8. The simplest method would be to interpret keys as *indices*, and to store the addresses in a *table* (fig. 2). *Sequential search* is easy to program, but will require an average of $N/2$ probes. If keys were sorted in *lexicographic order*, then presumably some kind of binary search over "outgoing" edges for each successive component of the key would reduce the search-time to $\sim \log N$ (at worst) *per component*. Fredman [31] indicated how an information-theoretic argument can help to improve it further. (The proof we develop here seems new.)

1.9. Consider an arbitrary node reached in progressing down the tree



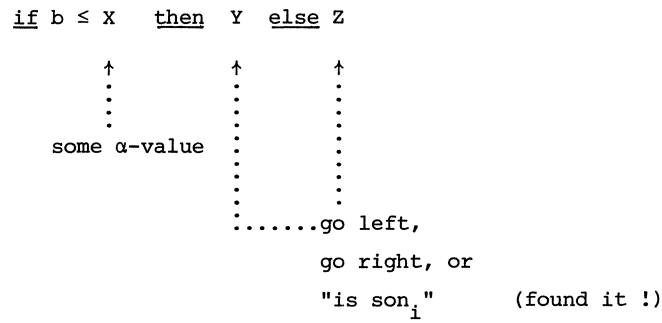
2	1	⋮
		⋮
2	2	⋮
		⋮
4	1	⋮
		⋮
4	2	⋮
		⋮
6	1	⋮
		⋮
6	2	⋮
		⋮
6	3	⋮
		⋮
6	4	⋮
		⋮
6	5	⋮
		⋮
6	6	⋮
		⋮
6	7	⋮
		⋮

fig.2. A table and a (lexicographic) tree



Suppose "son_i" has p_i descendant leaves, p_i ∈ N. In searching for an i such that α_i = b (sending us off to son_i), it is useful in probing to emphasize more frequently occurring α's.

The proper formulation makes use of a binary *search-tree*: each node in such a tree carries a query

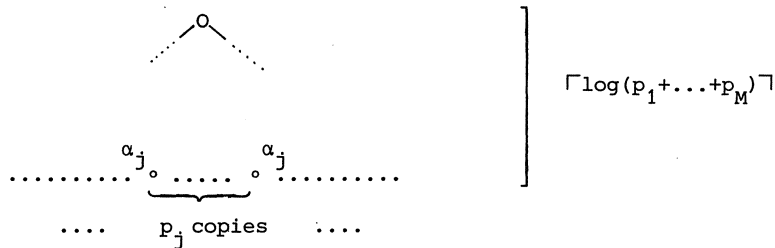


LEMMA. There exists a search-tree for finding b in

$$\leq \lceil \log(p_1 + \dots + p_M) \rceil - \lfloor \log(p_i + 1) \rfloor + 1 \text{ steps,}$$

where i is such that b = α_i.

Proof. Make p_j copies of α_j (for each j = 1, ..., M), and build a binary tree of minimum path-length



Call a node *j-critical* if and only if all its descendant leaves carry α_j and there is no node less deep in the tree with that property. For each j these are one or two *j-critical* nodes.

Make it a search-tree by assigning appropriate queries *topdown* such that the search for some b (presumably = α_i) is always directed to an

i -critical node. All nodes that are irrelevant now (thus, not on a search-path to some critical node) may be *purged*, after which the tree may have to be condensed back to binary form (thus reducing the search-length of some paths even further).

The identity of b is uncovered at the father of a critical node (and it obviously is the fastest way of identifying it).

The height of the "father" is bounded by the smallest t such that $d_t = 0$ in the following recurrence

$$d_0 = p_i$$

$$d_s = \left\lceil \frac{d_{s-1} - 1}{2} \right\rceil, \quad \text{for } s = 1, 2, \dots$$

which is $t = \lceil \log(p_i + 1) \rceil$. Thus a decision is found after $\lceil \log(p_1 + \dots + p_M) \rceil - \lceil \log(p_i + 1) \rceil + 1$ steps. \square

1.10. The resulting search-tree (after all unnecessary information has been purged) takes only little more space than was needed at the node of 1.9 anyway, and should be substituted. The original lexicographic tree becomes a cascade of local search-trees, and it enables one to find the address of a record (b_1, b_2, \dots, b_k) quickly.

THEOREM. *A record can be identified within $\sim \log N + 2k$ queries.*

Proof. Follow the trail of identifying (b_1, \dots, b_k) (fig. 3).

By lemma 1.9 the total number of queries needed is bounded by

$$\sum_0^{k-1} (\lceil \log N_i \rceil - \lceil \log(N_{i+1} + 1) \rceil + 1) =$$

$$= \lceil \log N_0 \rceil + \sum_1^{k-1} (\lceil \log N_i \rceil - \lceil \log(N_i + 1) \rceil) - 1 + k$$

$$\cong \log N + 2k. \quad \square$$

Recall that the original, unbalanced search-strategy did cost us about $\sim \log N$ queries *per component*.

1.11. Large collections of logically related records are hardly ever stored

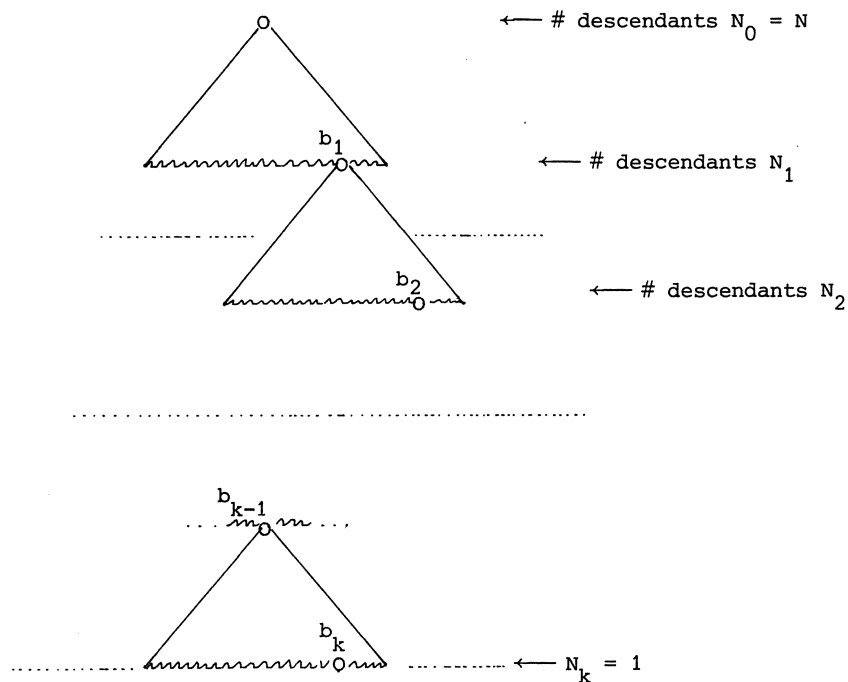


fig.3. Cascade

randomly, and are normally arranged in one or more perhaps specifically structured units called *files*. The reason for *file-design* (or *file-organisation*) is to store data appropriately on available storage-media to allow for efficient handling of the operations which the user requests.

1.12. In modern terminology a collection of stored operational files used by the application systems of some particular enterprise is called a *database* (Engles [25]).

1.13. Files can be distinguished by the supporting retrieval mechanism. The basic types are

(i) *sequential files*

- A record can be accessed only by scanning the entire file for the point where it is stored. It is common to *batch* successive record-

requests, and to avoid the need for rewinding the file for each individual request.

(ii) *random access files*

- The address of a record is determined from its primary key with a simple (perhaps tabulated) function, and access to its location follows (almost) immediately.

(iii) *indexed (sequential) files*

- A derived key is used as index in a directory structure to yield the *bucket* where a record must be, and ordinary retrieval from the bucket follows.

1.14. An important task in programming consists of finding the best data-organisation for an application in order that the data can be processed most efficiently.

However, the data may not be available in most desirable form or may be of much larger volume than can fit in core, and it can be equally important to modify algorithms for an application to work better and better under *given* constraints on the data-organisation.

These two directions of research together form the domain of the *complexity-theory of data organisation*.

1.15. An interesting example to show how a given data organisation can affect the choice of an algorithm is *Warren's transitive closure algorithm* for 0-1 matrices (Warren [63]).

We shall first consider the more traditional algorithm of Warshall (Warshall [64]) for this task.

1.16. An $n \times n$ 0-1 matrix $M = (m_{ij})$ can always be interpreted as the adjacency-matrix of a directed graph G_M on $\{1, \dots, n\}$ where

$$m_{ij} = 1 \iff \begin{array}{c} \leftarrow \longrightarrow \\ i \circ \quad \nearrow \circ j \end{array}$$

The 0-1 matrix $M^* = (m_{ij}^*)$ is called the *transitive closure* of M if and only if its coefficients satisfy

$$m_{ij}^* = 1 \iff \text{there is a path } i \rightarrow \dots \rightarrow j \text{ (of length } \geq 0 \text{) in } G_M.$$

1.17. Most transitive closure algorithms can be distinguished by the stepwise manner in which all paths $i \rightarrow \dots \rightarrow j$ are built up. One only needs to consider paths free of repeating nodes.

1.18. In Warshall's algorithm one builds M^* in stages $M^*(s)$ where

$$m_{ij}^*(s) = 1 \iff \text{there is a path}$$

$$\underbrace{i \rightarrow \dots \rightarrow j}_{\substack{\vdots \\ \in \{1, \dots, s\}}}$$

(of length ≥ 0) in G_M .

1.19. Obviously $M^*(0) = M$ and $M^*(n) = M^*$.

1.20. For $s = 1, \dots, n$ one can construct $M^*(s)$ from the previous stage by observing that

$$m_{ij}^*(s) = m_{ij}^*(s-1) \vee m_{is}^*(s-1) \cdot m_{sj}^*(s-1)$$

$$\quad \quad \quad \Delta$$

fixed factor for row i .

Denoting the i -th row of $M^*(s)$ by $M_i^*(s)$ we obtain the rule

$$M_i^*(s) := \underline{\text{if}} \ m_{is}^*(s-1) \ \underline{\text{then}} \ M_i^*(s-1) \vee M_s^*(s-1) \ \underline{\text{else}} \ M_i^*(s-1).$$

The algorithm can now be formulated as

```

"let  $M^*$  be  $M$ ";
for s := 1 to n do
  for i := 1 to n do
     $M_i^* := \underline{\text{if}} \ m_{is}^* \ \underline{\text{then}} \ M_i^* \vee M_s^*$ 
  od
od;

```

1.21. Warshall's algorithm works very well if each row of the matrix can

be *packed* in a single word and words can be "or"-ed directly.

As a contrast it is of interest to evaluate the algorithm for (very) large matrices which are stored externally in a sequential file. (It also serves as an adequate model for a paging environment).

Let consecutive records correspond to consecutive rows of the matrix, and assume that one can store ~ 2 records (at least) in memory at a time.

We shall use the number of records "paged in" by the algorithm as a measure for its complexity.

1.22. The algorithm must now be formulated as

```

rewind;
for s := 1 to n do
  scan for record s;
  Ms* := get record;
  rewind;
for i := 1 to n do
  Mi* := get record
  if mis* then reset record;
                Mi* := Mi* v Ms*;
                put Mi*
fi
od;
rewind
od;

```

and it follows that $\sim n^2 + n$ records must be read into memory.

1.23. With the given (sequential) organisation of the data Warshall's algorithm immediately becomes less attractive because the need for the coefficient m_{is}^* forces one to read in every record even if there is no subsequent action on the record.

In an algorithm recently proposed by Warren ([63]) this has been eliminated, and records are paged in only when necessary. (Our proof seems to be new.)

1.24. In Warren's algorithm M^* is built up in two passes.

Pass I yields an intermediate matrix M^{\otimes} in stages $M^{\otimes}(s)$ where

$$m_{ij}^{\otimes}(s) = 1 \iff \text{there is a path}$$

$$\underbrace{i \rightarrow \dots \rightarrow j}_{\substack{\vdots \\ \in \{1, \dots, i-1\}}}$$

for all $i \leq s$

and pass II subsequently yields M^* in stages $M^*(s)$ where

$$m_{ij}^*(s) = 1 \iff \text{there is a path}$$

$$\underbrace{i^{\otimes} \rightarrow \dots \rightarrow j^{\otimes}}_{\substack{\vdots \\ \in \{i+1, \dots, n\}}}$$

for all $i \leq s$.

1.25. Obviously $M^{\otimes}(1) = M$, $M^*(0) = M^{\otimes}(n)$, and $M^*(n-1) = M^*$.

For $s = 2, \dots, n$ one can construct $M^{\otimes}(s)$ from the previous stage by updating row s . One can similarly construct $M^*(s)$ from the previous stage for $s = 1, \dots, n-1$.

1.26. To compute $m_{sj}^{\otimes}(s)$ it is helpful to conceive of intermediate coefficients $m_{sj}^{\otimes}(i, s)$ where

$$m_{sj}^{\otimes}(i, s) = 1 \iff \text{there is a path}$$

$$\underbrace{s \rightarrow \dots \rightarrow j}_{\substack{\vdots \\ \in \{1, \dots, i\}}}$$

($0 \leq i < s$).

Obviously $m_{sj}^{\otimes}(0, s) = m_{sj}^{\otimes}(s-1)$, and the values of $m_{sj}^{\otimes}(i, s)$ for $i = 1, \dots, s-1$ can be computed using the rule

$$m_{sj}^{\otimes}(i, s) = m_{sj}^{\otimes}(i-1, s) \vee m_{si}^{\otimes}(i-1, s) \cdot m_{ij}^{\otimes}(s-1)$$

fixed factor for the row.

Denoting the s -th row of $M^{\otimes}(s)$ by $M_s^{\otimes}(s)$ we obtain the iteration


```

for i := 1 to s-1 do
  if  $m_{si}^{\otimes}(s)$  then  $M_s^{\otimes}(s) := M_s^{\otimes}(s) \vee M_i^{\otimes}(s-1)$ 
  fi
od

```

(where $M_i^{\otimes}(s-1) \equiv M_i^{\otimes}(s)$).

1.27. To compute $m_{sj}^*(s)$ it is similarly helpful to conceive of intermediate coefficients $m_{sj}^*(i,s)$, where

$m_{sj}^*(i,s) = 1 \iff$ there is a path

$$\underbrace{s \xrightarrow{\otimes} \dots \xrightarrow{\otimes} j}_{\substack{\vdots \\ \in \{s+1, \dots, i\}}}$$

($s \leq i \leq n$).

Obviously $m_{sj}^*(s,s) = m_{sj}^{\otimes}(n)$, and the values of $m_{sj}^*(i,s)$ for $i = s+1, \dots, n$ can be computed using the rule

$$m_{sj}^*(i,s) = m_{sj}^*(i-1,s) \vee m_{si}^*(i-1,s) \cdot m_{ij}^*(s-1)$$

$\begin{matrix} \Delta \\ \vdots \\ \vdots \end{matrix}$
 fixed factor for the row.

We can therefore obtain $m_{sj}^*(s)$ from the iteration

```

for i := s+1 to n do
  if  $m_{si}^*(s) = 1$  then  $M_s^*(s) := M_s^*(s) \vee M_i^*(s-1)$ 
  fi
od

```

(where $M_i^*(s-1) \equiv M_i^*(s) \equiv M_i^*(0)$).

1.28. The entire algorithm can now be formulated as

```

rewind;
scan for record 2;
for s := 2 to n do
   $M_s^* :=$  get record;
  rewind;
  for i := 1 to s-1 do

```

```

if m*si then M*i := get record;
           M*s := M*s v M*i
           else scan for next record
fi
od;
put M*s
od;
rewind;
for s := 1 to n-1 do
M*s := get record
for i := s+1 to n do
if m*si then M*i := get record;
           M*s := M*s v M*i
           else scan for next record
fi
od;
rewind;
scan for record s;
put M*s
od;

```

1.29. The number of records "paged in" by Warren's algorithm is bounded by

$$\begin{aligned}
 & 2(n-1) + \text{\#off-diagonal ones in } M^* \\
 & \sim n^2 + n - \text{\#zero's in } M^*
 \end{aligned}$$

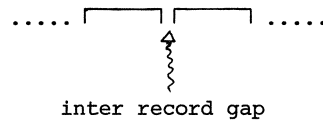
which shows that it behaves better than Marshall's the more sparse M^* is.

1.30. It follows that the design of good algorithms which use a pre-determined data-organisation can lead to non-trivial questions which the theory must consider just as well, and in the complexity-theory of data organisation one has to pay attention to the trade-offs between new data-structures and new algorithms.

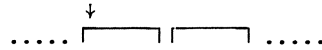
2. FILE-MERGING

2.1. Sequential files are *tape-like*.

We shall assume that all records have equal size and that records are stored with small, detectable intermissions



2.2. There is a *file-pointer*, normally positioned at the beginning of a record

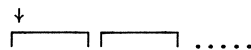


2.3. We shall later consider essentially random access files which allow an arbitrary (but per application bounded) number of pointers.

2.4. Any general purpose programming language provides some or all of the following operations on a sequential file:

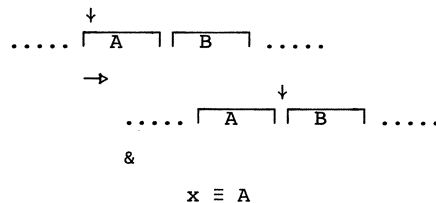
rewind

position the file-pointer back to the beginning of the tape



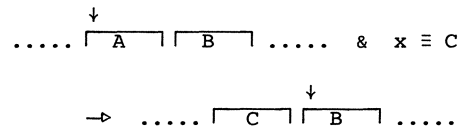
$x := \text{get record}$

the contents of the record currently pointed at is read into record-variable x and the pointer is advanced to the beginning of the next record



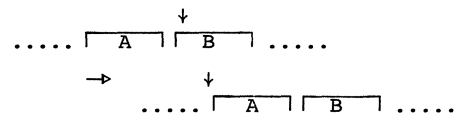
put x

the contents of record-variable x is stored into the current record-position on the tape. *The operation is usually allowed only if the pointer is located at the end of the file*, otherwise we assume it simply overwrites the previous record contents



reset record

positions the pointer at the beginning of the immediately preceding record (and void at the beginning of the tape)



and occasionally one may want to use derived operations like

scan for b

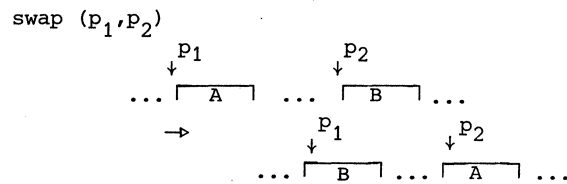
```

repeat
  x := get record
until eof  $\vee$  x.key = b;

```

(or an equivalent interpretation of it).

2.5. On files with more than one pointer we shall in fact allow primitive operations such as



and perhaps a boolean-valued instruction

ordered (p_1, p_2)
 yielding true iff p_1 is before p_2 and $p_1.\text{key} \leq p_2.\text{key}$.

EXAMPLE. A random access file



with n records can be inverted in linear time by

```

while  $p_2 > p_1$  do swap ( $p_1, p_2$ );
     $p_1 := p_1 + s$ ;
     $p_2 := p_2 - s$ 
od;

```

(where s is the record size).

2.6. A (sequential) file is said to be *ordered* if and only if the keys of consecutive records form a nondecreasing chain.

2.7. The process of combining two ordered (sequential) files into a single ordered (sequential) file is called *merging*.

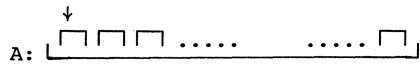
A merging procedure is said to be *stable* if and only if it leaves the relative order of records with the same key in the original files unchanged. (See Knuth [44].)

2.8. It is now reasonably well-understood how to merge two files in case some auxiliary tapes are available (Wirth [66]). It is less obvious how to proceed with minimal storage-requirements, and in particular how to do a linear time merge without the use of an unbounded number of (perhaps hidden) pointers or links.

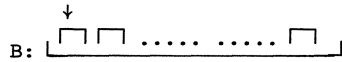
This problem was first cited in Knuth ([44], p.388) and solved in 1974 by Horvath [40] and by Luis Trabb Pardo [58].

To acquire a better appreciation for the specific constraints of the problem we shall first consider a linear time two tape merge-procedure of Floyd & Smith [28].

2.9. Suppose we are given two ordered sequential files of records of a certain size s



A: N records, but perhaps
only $n \leq N$ with a distinct key.



B: M records, but perhaps
only $m \leq M$ with a distinct key.

We shall assume very little about s , but only require that it provides enough bits for the distinct keys:

$$s \geq \log(n+m).$$

2.10. The idea is to copy file A (at the end of B) and file B (at the end of A) and insert the proper sequence numbers in intermediate records first, showing in what order the file-members must be assembled.

An immediate problem shows up if we want to perform this in linear time: the sequence-numbers for file A (and similarly for file B) will occupy

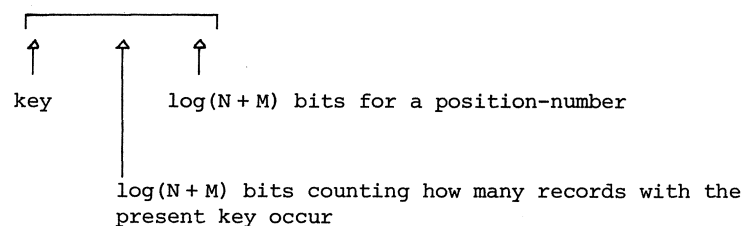
$$\sim N \cdot \log(N+M) \text{ in tape-length}$$

and this may not at all be bounded in terms of the original file-length

$$\sim N \cdot s$$

which could be as little as $N \cdot \log(n+m)$.

2.11. Instead of complete sequence-numbers we shall therefore insert intermediate records once for each distinct key only (a common trick in *data-reduction*), but augment the information per record to



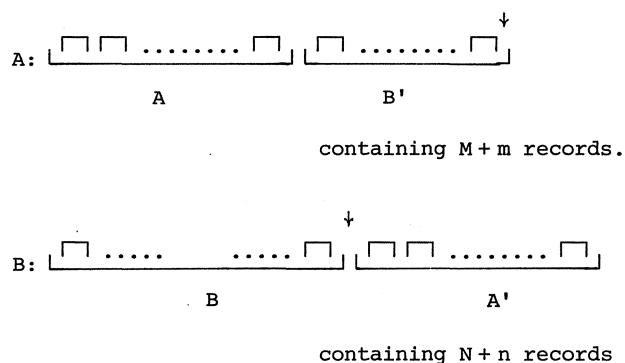
This will require an extra space of only

$$\begin{aligned}
 &\sim n (s + 2 \log(N+M)) \\
 &\leq N \cdot s + 2n \log(N+M) \\
 &\leq N \cdot s + \gamma' \cdot (N+M) \log n \\
 &\leq \gamma \cdot (Ns + Ms)
 \end{aligned}$$

for tape A (and similarly for tape B), which is now bounded by a constant factor in the length of the original files.

2.12. The first step of the algorithm will copy A and insert occurrence-numbers for all distinct keys.

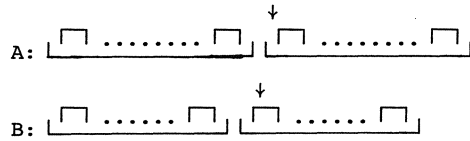
After a similar procedure is applied to copy tape B into B' at the end of tape A we have



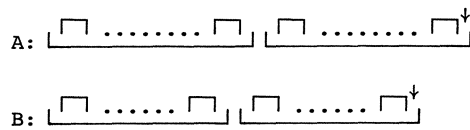
and the time required so far is

$$\begin{aligned}
 &\sim Ns + Ms + \underbrace{Ms}_{\text{rewind}} + \underbrace{3Ns}_{\text{scan B}} + \underbrace{Ns}_{\text{scan + copy}} + \underbrace{Ns + \gamma(Ns + Ms)}_{\text{copy onto A'}} + Ns + \\
 &+ \gamma(Ns + Ms) + \underbrace{Ms}_{\text{rewind over}} + \underbrace{3Ms}_{\text{scan + copy}} + \underbrace{Ms + \gamma(Ns + Ms)}_{\text{copy onto B'}} \\
 &\quad \text{A'} \text{ and B} \quad \quad \quad \text{B} \\
 &\sim (6+3\gamma) Ns + (7+3\gamma) Ms.
 \end{aligned}$$

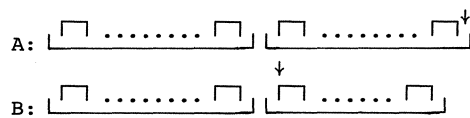
2.13. Rewind over B'



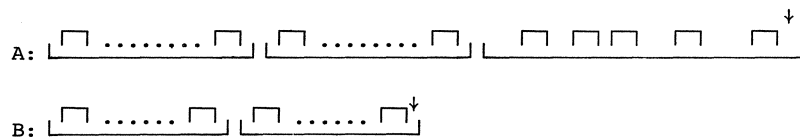
and perform a "dummy" merge of A and B by looking at the intermediate, accounting records only and inserting the appropriate position-number for the first record of each distinct key in the proper field.



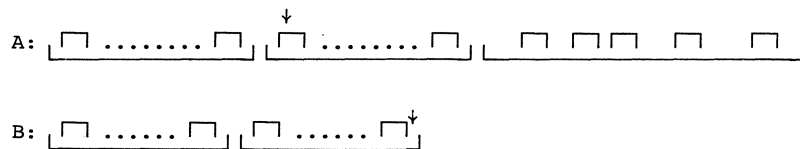
Rewind over A'



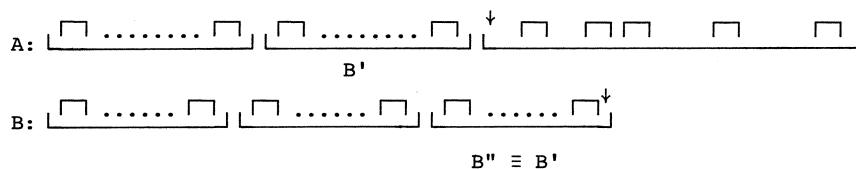
and copy the records of A (taken from A') at the end of tape A in the proper relative final position (which we can determine from the information in the accounting records which should *not* be copied over)



Rewind tape A to the beginning of B'



and copy all of B' onto tape B:



Finally, rewind over B'' and copy the records of B (taken from B'') into their final positions on tape A (again using the information from the accounting records)



The merged version of the original files appears as the last block on tape A.

2.14. We have now obtained the essential result of Floyd & Smith.

THEOREM. *The two-tape merge-procedure is stable and requires time linear in the length of the merged files.*

Proof. Stability follows from the construction. The time required for all successive stages adds up to

$$\sim (18 + 13\gamma)(N_s + M_s). \quad \square$$

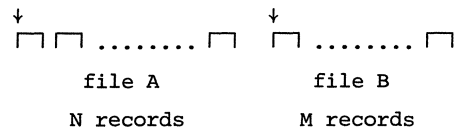
2.15. Note that the algorithm requires no third tape, and we could have used the given space even more economically had we overwritten parts which were at some stage in the algorithm not needed anymore. We seem to be close to an answer of Knuth's problem, but there is apparently no way to eliminate the need for all accounting records without thereby causing a non-linear increase of time.

It doesn't seem to be easier for random access files either, and therefore the following result of Luis Trabb Pardo [58] is of interest.

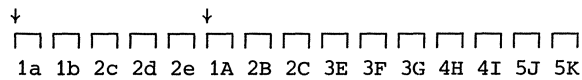
THEOREM. *There is an algorithm to merge two "rigid" random access files in a stable way in linear time and no auxiliary storage using only a bounded number of pointers.*

2.16. Trabb Pardo's algorithm contains a number of interesting contributions which we shall present in a simplified form to better demonstrate the ideas.

Assume that the files are stored as one chunk

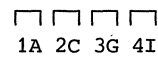


EXAMPLE.



2.17. A *block* of consecutive records is called an *internal buffer* if and only if it is fully ordered and contains only records with a distinct key.

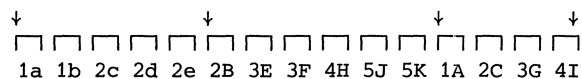
EXAMPLE.



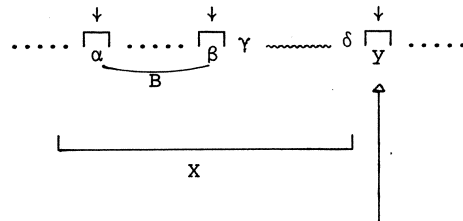
2.18. The first phase consists of *extracting* an internal buffer of size $\sim \sqrt{N+M}$ from the given files (which we always can if there are enough distinct keys).

The buffer-elements can be collected by scanning the files from left to right, picking a new larger record only if it is the *last* sample with the current key in the row (to ensure stability later).

EXAMPLE



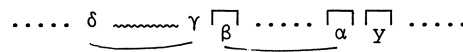
The buffer-elements are kept together in a "growing" block B sliding steadily to the right, while a pointer advances to search for a next candidate



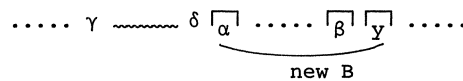
suppose this record must be assembled next.

Then:

reverse X



and reverse the marked parts



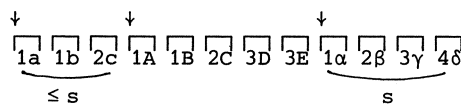
and we can proceed until B is full size. B is then exchanged in the same way to the end of the file.

The time required for this phase is bounded by

$$\sim \underbrace{N + M}_{\substack{\text{each stretch of elements} \\ \text{can be involved in rever-} \\ \text{sion once}}} + \underbrace{\left(\sqrt{N+M}\right)^2}_{\substack{\text{but the elements of the buffer} \\ \text{are involved in it perhaps as many} \\ \text{times as its number of elements.}}$$

2.19. It could very well happen that file A is "small" with respect to $s = \sqrt{N+M}$.

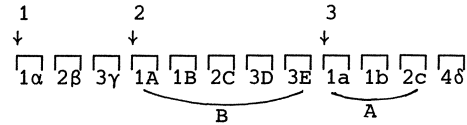
EXAMPLE.



(If it is so for B we do the next procedure "backwards".)

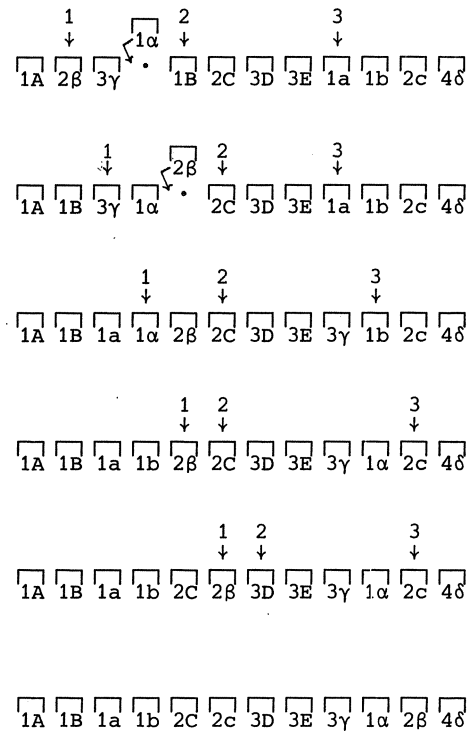
Exchange A with the first part of the buffer (in linear time using the reversion-trick):

EXAMPLE.



and merge B and A into place in a *stable* manner by moving pointers 2 and 3 steadily rightwards, each time exchanging the smallest with the (buffer-) record under 1.

EXAMPLE.



The buffer-records act as "dummies" which keep making place for records merged into position, and in the end they recollect in the buffer-zone.

Sort the buffer back into order, and complete the procedure (in linear time) by merging the buffer back into the file (how to do it will

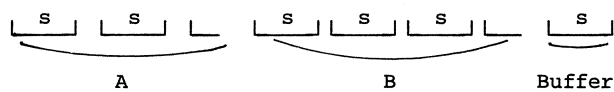
follow later).

It is now clear why an internal buffer is essential: it is a group of records which we can move around to provide a space where needed and which we can bring back into stable order by just sorting!

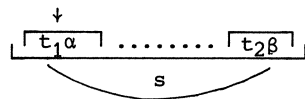
(Another application of the buffer will follow next.)

2.20. Let us now assume that both A and B are "long". Think of A and B as divided into blocks of size s .

EXAMPLE.



It is convenient to have a notation for the first and last key of a block of known size:



X

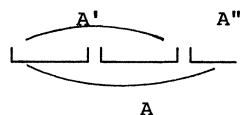
first(X) = t_1

last(X) = t_2

The values can be picked up using only two probes into the block.

2.21. If the size of A is no integer multiple of s

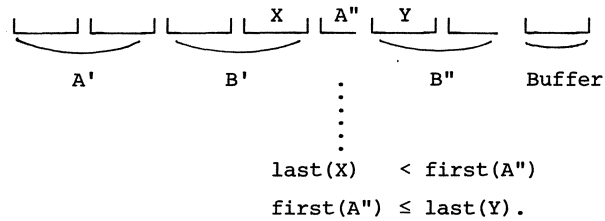
EXAMPLE.



then get rid of the left-over part A'' by merge-exchanging it (as a block) into B as far right as possible, i.e., after the last block X in B for which

last(X) < first(A'').

EXAMPLE.



X can be found by a simple scan, and using the reversion-trick the exchange follows in linear time.

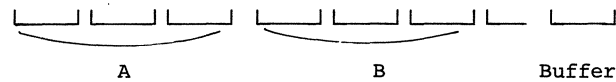
The task is now reduced to

- merge A' and B'
- merge A'' and B''

(where the last part follows as in 2.19). Note that the records will appear in final position if we do the merge in place.

2.22. Without loss of generality we can therefore assume that A has a neat block-arrangement.

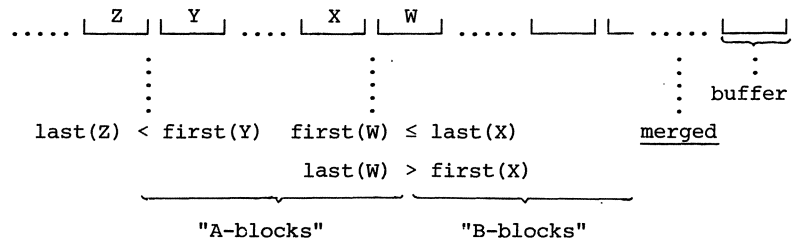
EXAMPLE.



Let the blocks of A and B be A_1, \dots, A_k and B_1, \dots, B_ℓ . Since $s \sim \sqrt{N+M}$, $k + \ell \leq s$.

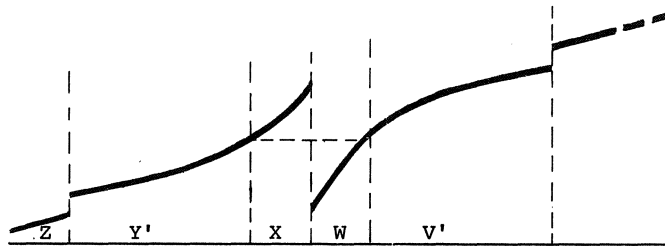
2.23. The idea is to first "merge" the A-blocks (as a block) in between the B-blocks as far to the right as possible while preserving stable order.

EXAMPLE.



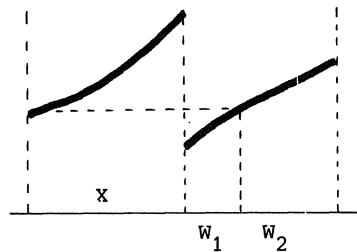
The situation can be recovered in time linear in the length of the stretch of elements involved. Note the L-block that may drag along. A graph helps to show the precise order-relation between the parts.

EXAMPLE.



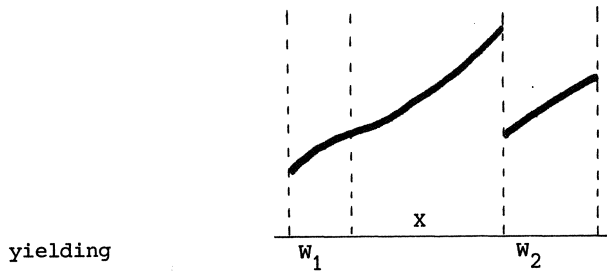
2.26. Consider XW

EXAMPLE.

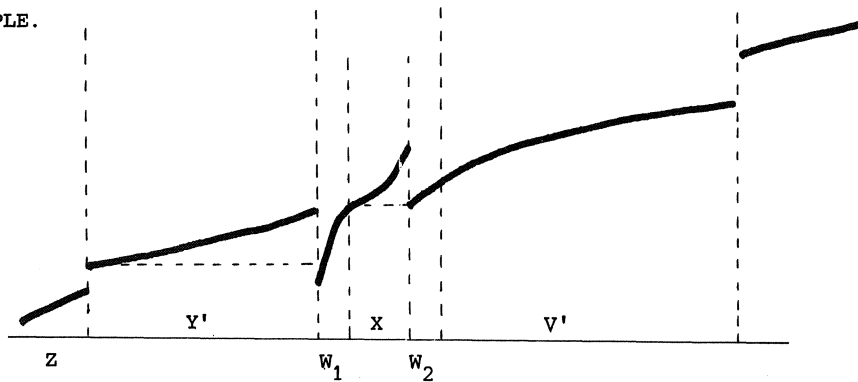


and use the reversion-trick to exchange X (as a block) as far to the right into W as possible.

EXAMPLE.

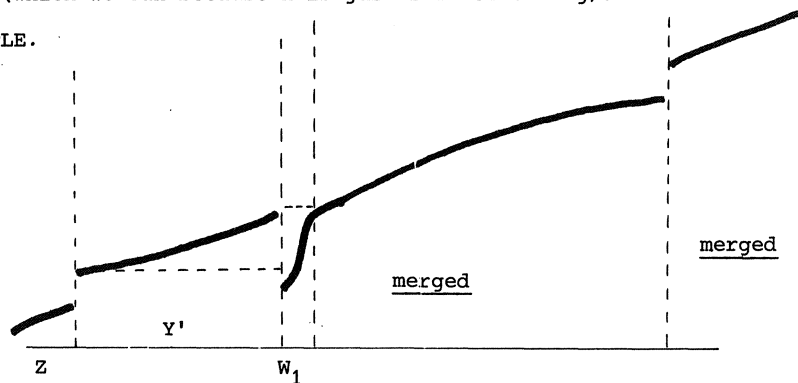


EXAMPLE.



2.27. Now merge X into W_2V' in a stable manner using the buffer-trick from 2.19 (which we can because X is just s records long).

EXAMPLE.

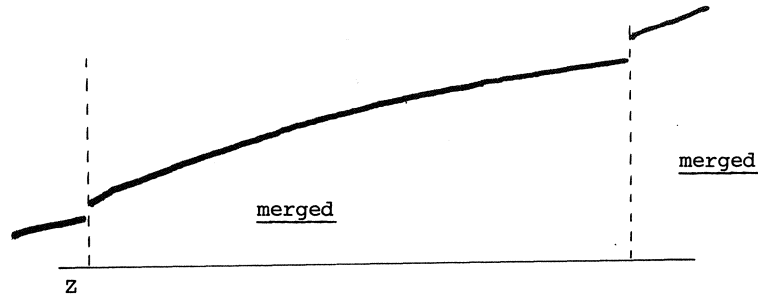


and we finished off "half" of the zone in time just linear in the number of elements involved.

Don't resort the buffer again each time we do such a phase.

With the same technique one can merge Y' and W_1 in a stable manner and "finish" the entire zone.

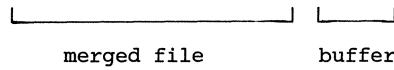
EXAMPLE.



Proceed left and locate the next zone, thus continuing and completing the merge in linear time and a bounded number of pointers.

2.28. In the last phase of the algorithm one should resort the buffer (in time $\sim s^2$).

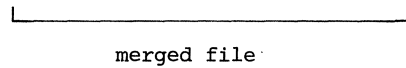
EXAMPLE.



and insert the buffer back into the file again, leaving each buffer-record at the far right of the stretch of file-records with the same key to ensure final stability.

With the reversion-trick one can steadily slide the buffer to the left and absorb it in linear time and no extra space.

EXAMPLE.



2.29. This completes the essential part of Trabb Pardo's minimal requirement stable merge procedure. (What should happen if no full size buffer can be extracted is left to the reader, but all necessary techniques to use have been presented.) Several steps were due to Kronrod and to Horvath.

2.30. Note that an "in-place" linear time merge procedure immediately yields

an "in-place" and minimal requirement $O(n \log n)$ sorting method by merging the data into adjacent blocks of a size which is an increasing power of two.

3. TABLES AND BALANCED TREES

3.1. The choice of a proper data-organisation for a task depends on how the user wants to operate on the data.

In the design-philosophy of programming languages there is a tendency to make *data-structuring* more and more automatic. Ultimately the user should be able to specify the *axioms* and *rules* for the operations on the data in some formal language, and the system should determine a consistent (and preferably efficient) internal representation. This has been studied theoretically (see Spitzen & Wegbreit [55]), and it is anticipated for PASCAL in terms of graph-specification primitives (Gehani [33]).

3.2. Perhaps the simplest requirement on a (random access) file or *table* is that one can perform

```

find (k)
    - locate the record with key k -
insert (k)
    - put the record with key k into the table -

```

for all keys quickly.

Since the universe of key-values may be much larger than the set of keys actually occurring in an application, an efficient and easy computable *key-to-address transformation* should presumably "chop down" a key into an index in a small *hash-table* where the record-address is (or will be) listed.

Such *hash-functions* $h(k)$ are almost by definition many-to-one, and it may happen that different keys are mapped onto the same *hash-address*.

The anomaly in insertion occurring when a hash-address is already occupied is called *collision*.

3.3. Collisions can be resolved by entering records in an appropriate *overflow-area*.

Another method (called *open addressing*) is to search the table for an opening "nearby" by systematically probing locations at a distance of

$$\text{incr}(k)$$

further and further down the table. The *hash-increment function* is assumed to be positive, and is presumably chosen to lead to an opening fast.

The case when

$$\text{incr}(k) = F(h(k))$$

for all keys and some easy function F is sometimes called *double-hashing*.

3.4. Suppose we maintain a table

$$\text{item}[0 : M-1]$$

in this manner, where $\text{item}[i]$ contains a key k such that $h(k) = i$ if such a key was ever added and probably 0 otherwise. (It is strictly true only until collisions have occurred.)

One can search for k with

```

visits := 0;
probe := h(k);
while visits < M do
  x := item (probe);
  if x = 0 then not found & exit for L fi;
  if x = k then found & exit for L fi;
  probe := probe - incr(k) (mod M);
  visits := visits + 1
od;
table full;

```

L:

and *insert* k with essentially the same routine after replacing "not found" by

$$\text{item}(\text{probe}) := k.$$

In order that *all* locations of the table be probed in searching we demand

that $\text{incr}(k)$ is relatively prime to M .

3.5. A hashing-method should avoid *clustering* of records (eliminating collisions from occurring too often), and is therefore usually dependent on an estimated probability distribution of the keys.

3.6. Amble & Knuth [5] proposed in 1973 to maintain *ordered hash-tables*.

It does not mean that the table itself is entirely ordered, but that large parts of the *chains* are.

One can search for k with

```

visits := 0;
probe := h(k);
step := incr(k);
while visits < M do
x := item (probe);
if x < k then not found & exit for L fi;
if x = k then found & exit for L fi;
probe := probe - step (mod M);
visits := visits + 1
od;
table full;
L:

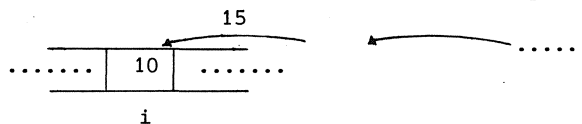
```

The kind of arrangement does not specifically make successful searches faster, but *it tends to detect unsuccessful searches much earlier*.

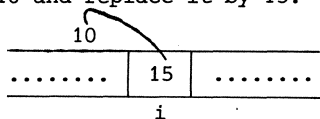
3.7. If a set of keys were inserted in an empty table from largest to smallest (in that order) with the algorithm of 3.4, then an ordered hash-table results.

The insertion of randomly presented keys into an ordered hash-table is somewhat trickier.

EXAMPLE. Suppose we want to move 15 into position i now occupied by 10.



Then we simply remove 10 and replace it by 15.



and try to re-insert 10. (Note that the replacement of an item by a *larger* item leaves all existing chains through that location consistent.)

Re-inserting 10 by back-chaining from $h(10)$ on will bring us back to i again, so we better start the re-insertion procedure from position i on right-away. Note that for consistency the stepsize must become $\text{incr}(10)$ now.

One can insert k in an ordered hash-table with

```

visits := 0;
probe := h(k);
step := incr(k);
key := k;
while visits < M do
  x := item (probe);
  if x = 0 then item (probe) := key & exit for L fi;
  if x < key then item (probe) := key;
                    step := incr(x)
                    fi;
  if x = key then found & exit for L fi;
  key := x;
  probe := probe - step (mod M);
  visits := visits + 1
od;
overflow;
L:

```

Since the x -value keeps descending M is still a valid bound on the search-length.

Verify that insertion in an ordered hash-table indeed leaves another ordered hash-table!

3.7. The arrangement of N keys ($N < M$) in an ordered hash-table is *unique*.

3.8. Entirely different requirements of a data-organisation occur in *ordered* files or in structuring the directory of indexed files where one would like to perform operations

```
find (k)
insert (k)
```

but also

```
delete (k)
    - remove the item with key k -
```

and perhaps

```
split file at k
    - separate the file in a part with items < k
      and a part with items ≥ k -
```

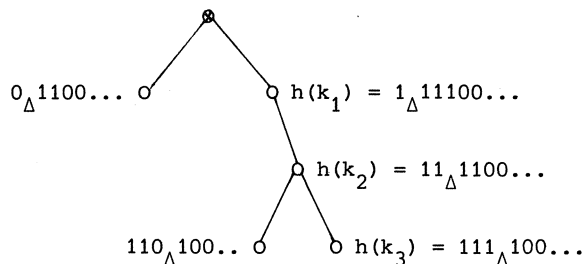
as efficiently as possible.

3.9. Deleting elements in a hash-table could break up existing chains and make items further down the chain inaccessible unless a cumbersome reorganisation procedure is applied. Most implementations therefore avoid *direct addressing schemes* and use *binary trees* instead.

3.10. An interesting intermediate file-organisation suggested by Coffman & Eve [19] results if we choose a hash-function which maps keys onto bit-strings.

One can interpret bitstrings as coding paths down a tree and maintain the table as a *hash-tree*. To insert k we follow $h(k)$ bit-after-bit and enter it into the first open node encountered, unless the key was found on the way.

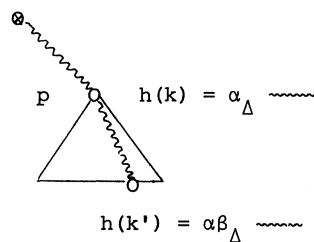
EXAMPLE.



Note the part of the hash-code used to enter k_1 , k_2 and k_3 . (0 \equiv "left", 1 \equiv "right").

Collisions are automatically avoided except when $h(k)$ is "too short" to reach a distinguishing position (in which case an overflow-area must be entered).

3.11. Deleting k from its position p in the tree is easy when it is a leaf, but otherwise



one can delete an arbitrary leaf k' from its subtree and move it up into p to keep the hash-tree *consistent*, which is very easy also (provided k' is not too far away).

3.12. The average search-time in a hash-tree with N items will be $\sim \log N$ or less (assuming a uniform distribution), but series of inserts and deletes can seriously impair it and cause long and inefficient paths in the tree.

3.13. Various kinds of *binary search tree* organisations have been proposed which remain *balanced* throughout.

In such trees the file-elements are normally arranged in left-to-right order at the leaves, and the internal nodes contain appropriate queries

if key < k then left else right

which enable one to retrieve stored items in a top-down search (see also 1.9). To maintain a balance in a tree with N items one must design insert- and delete-procedures which keep the worst case distance from the root to any leaf $\sim \log N$.

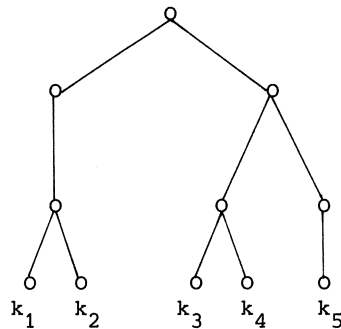
Known balanced tree-models include AVL-trees (Adelson-Velskii & Landis [1]), various kinds of B-trees (Bayer [8]), 2-3 trees (Hopcroft, see [3]), and 1-2 trees (Maurer & Wood [45]).

We shall consider a simplified kind of 1-2 tree recently proposed by Ottmann & Six [49].

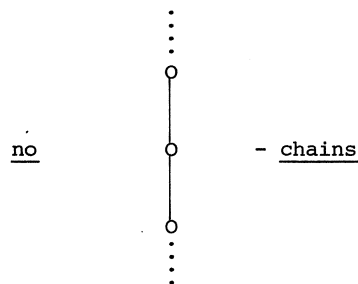
3.14. A binary (search) tree is called an *HB-tree* if and only if

- (i) all leaves have equal depth,
- (ii) each node with only one son has a brother with two sons.

EXAMPLE.



The root of an HB-tree must have two sons (unless it is a leaf). If a node has only one son, say p , then p must have two sons or else be a leaf. It follows that there can be



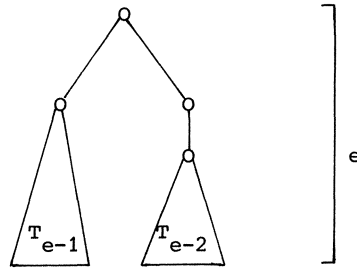
and HB-trees must always be quite "dense".

THEOREM. *The depth d of any HB-tree with N leaves satisfies*

$$\lceil \log N \rceil \leq d < 1.44 \dots \log(N+1) - 0.32 \dots$$

Proof. (Our argument differs from Ottman & Six [49]). In a binary tree of depth d we always have $N \leq 2^d \rightarrow d \geq \lceil \log N \rceil$.

Let T_e ($e \geq 0$) be an HB-tree of depth e with the smallest possible number of leaves. It follows that (up to symmetry) T_e must be of the form



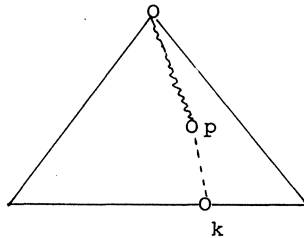
for $e \geq 2$, and T_e must have F_{e+2} leaves by induction (where F_i denotes the i -th Fibonacci-number). The depth d of a tree with N leaves satisfies therefore

$$\frac{\phi^{d+2}}{\sqrt{5}} - 1 < F_{d+2} \leq N. \quad \square$$

3.15. Observe that *HB-trees are AVL-trees in which all paths from the root down the tree are made equally long by inserting extra nodes of degree one at the appropriate places.* The effect, however, is that the operations for inserting or deleting items are much easier to explain.

3.16. One can insert k in the following manner.

Find the position in the tree where k must be inserted among the leaves.



If p had only one son then it now has two and the structure remains an HB-tree. If p had two sons already then it now has three, and we must

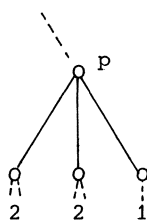
enact a procedure *split p*.

It is now advantageous to *pretend* that all leaves have two sons (for reasons of consistency only).

The *invariant* maintained throughout the algorithm is that

split p is called if and only if *p* has three sons at least two of which have degree 2.

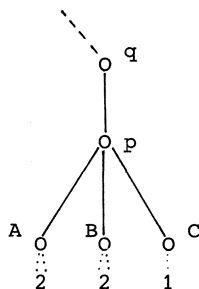
EXAMPLE.



Thus at least one of the outermost sons of *p* must have degree 2 and we shall designate one as the *active son* (on the *active side* of *p*).

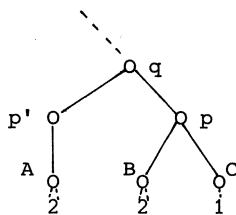
3.17. If *p* is the only son of *q* (say)

EXAMPLE



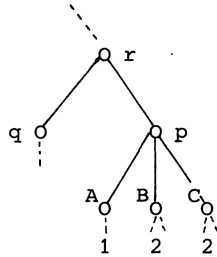
then create a new brother *p'* on the active side of *p*, connect the active son of *p* to *p'* and finish.

EXAMPLE.



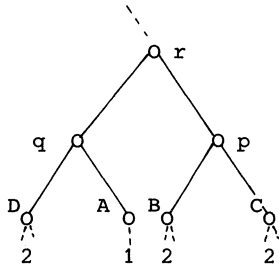
If p has a brother q and father r

EXAMPLE



then make the son of p on the q -side a son of q if q has degree one (which makes an HB-tree even if A had degree one) and finish

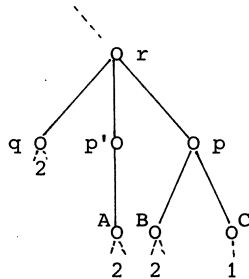
EXAMPLE



(as originally only son of q D must have degree 2, see 3.14)

and otherwise (when q has degree 2) create a new brother p' directly on the active side of p , make the active son of p the son of p' , and continue with *split* r (observing that *the invariant is preserved!*).

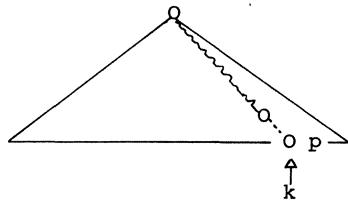
EXAMPLE



If p has no father we must have reached the root of the tree and must create a new root one level higher.

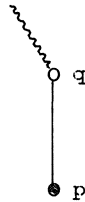
3.18. One can delete k in the following manner.

Find the position among the leaves where k is located.



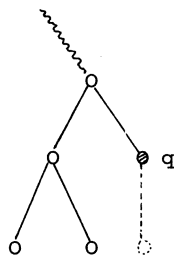
If p is the only son of q

EXAMPLE



then delete p

EXAMPLE



(note that q must have a brother of degree 2)

and continue to *delete* q .

After this initial phase the invariant maintained by the algorithm is that

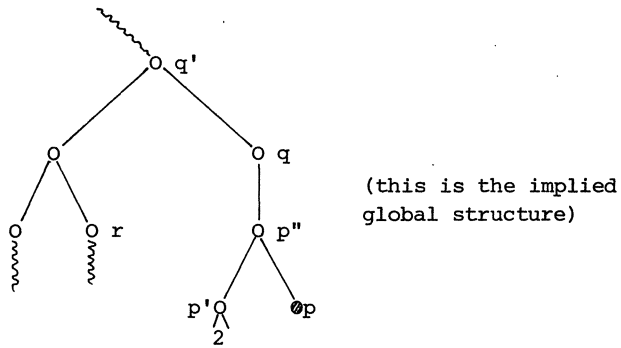
delete p is called if and only if p must be deleted and it has a brother of degree 2.

3.19. (Our deletion-procedure differs slightly from Ottman & Six.)

In a *delete* p the order between p and his brother is irrelevant, so we shall always assume the brother to be on the "convenient" side.

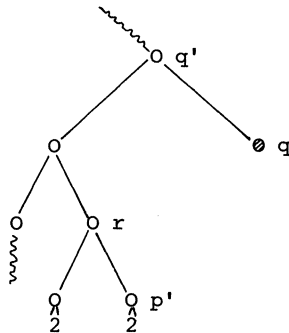
If the father of p is the only son of q

EXAMPLE



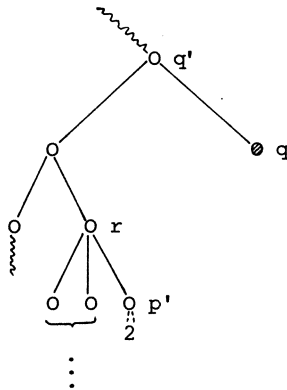
then make p' a son of r in case r has degree one, omit p and p'' , and continue with *delete* q (observing that the invariant is preserved).

EXAMPLE.



and otherwise (when r has degree 2) we make p' a son of r also while dropping p and p'' but observe that now the invariant is satisfied for *split* r .

EXAMPLE

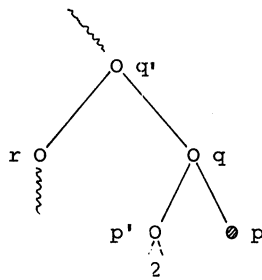


the original sons
of r cannot both
have degree one.

If *split* r (etcetera) runs out before q' is reached we continue with *delete* q (observing that the invariant will hold), otherwise we omit q and finish.

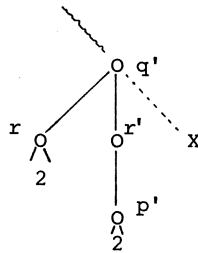
If the father of p has a brother

EXAMPLE



then make p' a son of r , omit p , and continue with *delete* q in case r had degree one, and otherwise (when r has degree 2) construct a new brother r' between r and q , make p' the son of r' and omit p and q , and finish.

EXAMPLE.



If the father of p is no son of any node we must have reached the top of the tree, and we can simply omit p and make his brother the new root.

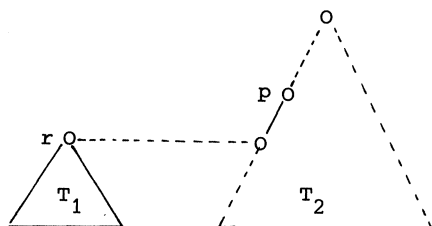
3.20. The procedures show

THEOREM. *One can execute find, insert, and delete-instructions on an arbitrary HB-tree with N leaves in $\sim \log N$ steps per instruction.*

3.21. Given two HB-trees T_1 and T_2 (in that order) one can merge the corresponding sets in the following manner.

Assume that T_1 is smaller than T_2 (and modify the algorithm accordingly otherwise). Find the element on the left-side of T_2 equally high as root (T_1).

EXAMPLE



Make r a son of p . If p had one son before it now has two and we can finish, otherwise the invariant for *split* p is enacted and an HB-tree results only after some additional moves.

Merging two trees is thus also bounded by $\sim \log N$ (or, more precisely, by $\text{depth}(T_2) - \text{depth}(T_1)$).

HB-trees can now also support instructions

split at k

in $\sim \log N$ steps as follows. Cut the tree along the path down to k in two parts, each part consisting of correctly ordered but loose hanging HB-trees. Merge the smaller trees into the larger continuously.

3.21. The idea of HB-trees can be generalised following similar steps to k -ary search trees (Maurer, Ottmann & Six [46]).

3.22. All previous tree-structures could be used also for implementing

address of the corresponding root.

3.24. Union-find programs obviously benefit from *balancing* the trees.

Let the *weight* of a node be equal to the number of leaves in its subtree.

Build the trees from items of

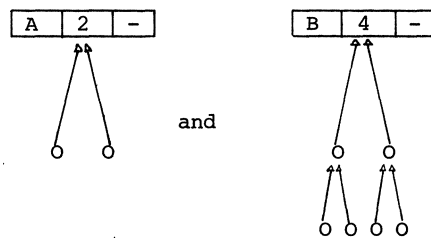
```

type node = record
    name: a character-string;
    weight: an integer;
    father: ↑ node;
end;

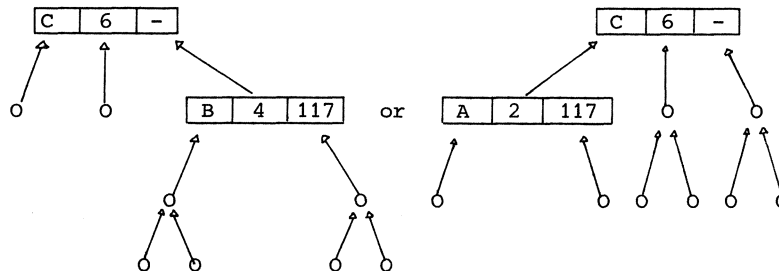
```

The idea of balancing in performing a "union" is to *never merge a tree of larger weight into a tree of smaller weight*.

EXAMPLE. Note the difference in merging.



into



(Note also that in merging some of the original set-names are automatically "lost").

which means that the resulting tree has

$$\text{weight} \geq 2^{d-1} + 2^{d-1} = 2^d$$

also.

The maximal path-length d therefore satisfies $2^d \lesssim m \rightarrow d \sim \log m$. \square

3.27. In merging, one set will always lose its original name (see 3.24) and because the result of balancing is rather unpredictable we may not know in advance which one. Note also that under circumstances the same root may be renamed more than once.

It is therefore impossible in the given implementation to perform

$\text{lca}(k_1, k_2)$

- find the *node* (or name) of the least recent ("smallest") set to which k_1 and k_2 were made to belong together, and undefined if no such set exists -

on-line.

Aho, Hopcroft, and Ullman [4] proved that one may execute $\sim m$ union-find-lca instructions online in an average of $\sim \log m$ steps per instruction.

Using a different data-structure for forests we have been able to reduce this to a less than logarithmic bound (Van Leeuwen [60]).

THEOREM. *One may execute $\sim m$ union-find-lca instructions in $\sim \log \log m$ steps per instruction on the average.*

(We note that Aho, Hopcroft, and Ullman [4] considered a slightly more general type of merge-instruction to occur in conjunction with lca's than we permitted here.)

4. PATH COMPRESSION

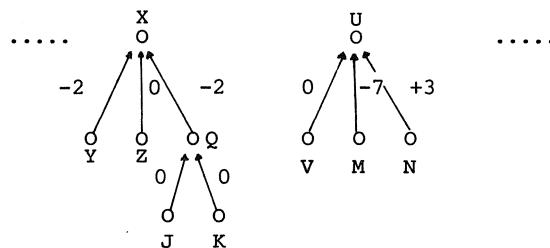
4.1. The *logical* development of set-manipulation programs often benefits from representing sets as trees, even though an actual implementation may

use different structures like arrays. It relates to a similar aspect of information management systems where the *data-model* of the user may be very different from the organisation of the physical database.

Trees were already used as data-structures in the early fifties when "sorting and searching" emerged as an area of intensive study.

In 1964 the idea was used by Galler & Fischer [32] in improving earlier work on Arden, Galler, and Graham [6] on MAD and aspects of FORTRAN-implementation. They considered how one can efficiently assign addresses to EQUIVALENCED variables in FORTRAN-programs at compile-time, and used trees which actually carried names of variables in all nodes (including array-identifiers).

EXAMPLE. EQUIVALENCE (U[3],M[10],N) is represented, for example, among



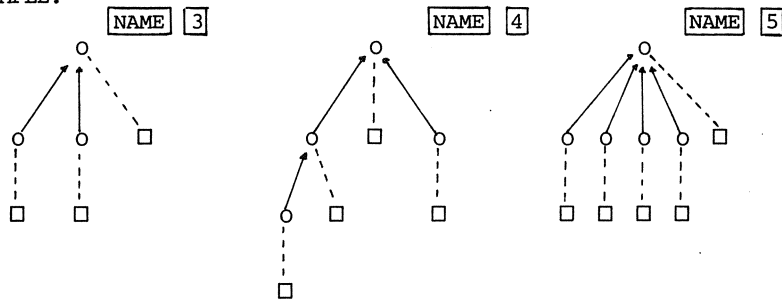
4.2. Trees have been successfully applied in programs for manipulating *disjoint sets*, which occur for instance when one is building *classes of an equivalence-relation* or in a variety of *graph-algorithms*.

4.3. There are only a few programmingtricks for making ordinary set-manipulation programs more efficient, but some of these techniques are very powerful and fundamental.

Following an approach of Aho, Hopcroft, and Ullman [3] we shall study union-find programs for more abstractly presenting an analysis of *balancing* and *path-compression*.

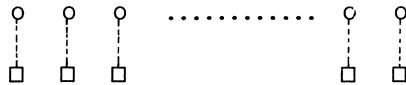
4.4. A collection of disjoint trees (representing disjoint sets) is called a *forest*. We shall assume that set-elements are attached to leaves only.

EXAMPLE.



The *weight* gives the number of set-elements stored in a tree. Directories tell where elements are stored, and where for each set-name the corresponding root is, but to execute a "find" one must explicitly follow father-links as in 3.23.

4.5. Assume (for simplicity) that we start with n singleton sets



and execute m finds and (at most) $n-1$ unions.

Unions take constant time each, but the cost for finds is to be the *number of edges* traversed in each find-path.

Let F be the set of edges which will ever appear in the finds. (Note that these edges may not all be present in the current forest at some moment.)

4.6. A direct implementation could cost time $\sim mn$ in worst case.

Balancing helps, and a rule as in 3.24 (*the weighted union rule*) for resolving unions reduces execution-time to $\sim m \log n$ in worst case.

Tritter (according to Knuth [43]) and also McIlroy & Morris (according to Hopcroft) apparently first realized how one can reduce the cost of (future) finds. In executing a $\text{find}(k)$ one should save all nodes encountered on the find-path in a stack and attach all *directly* to the root once it is located (see fig.7). It is sometimes called the *collapsing rule*, but nowadays better known as *path-compression*.

4.7. It is clear that balancing and path-compression (or even path-compression alone) should eventually lead to very low trees and thus reduce the

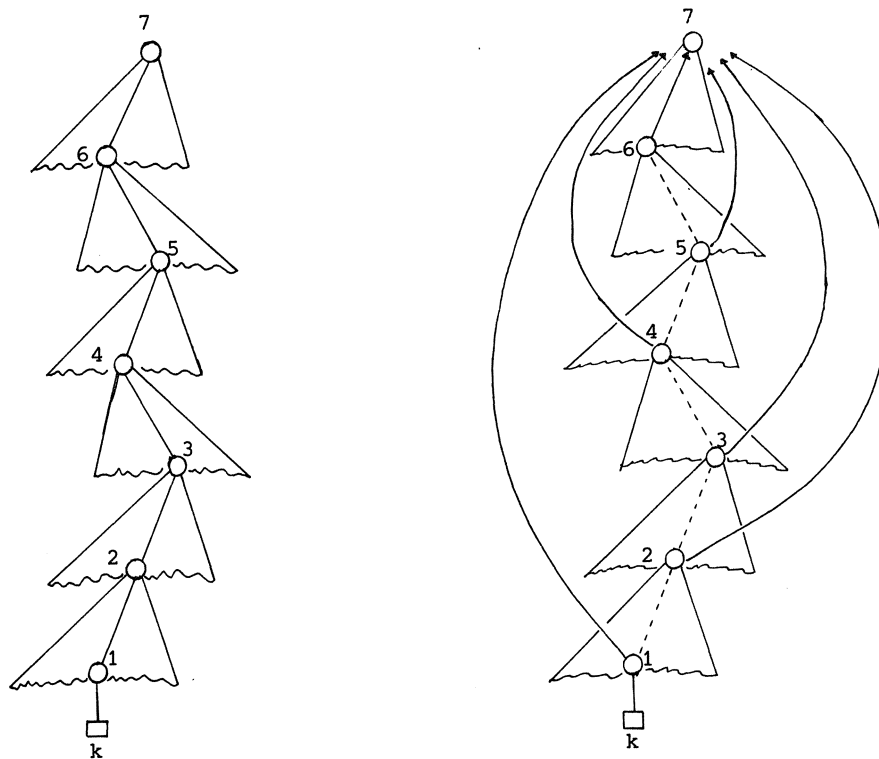


fig.7. The effect of path-compression in find(k)

cost for later finds substantially, but the precise analysis requires a new and important accounting-technique.

Paterson proved in 1972 that when only path-compression (and no balancing) is used and $m \sim n$, the execution-time is still bounded by $\sim n \log n$. Hopcroft & Ullman [39] found an elegant level-crossing argument soon after and proved that balancing and path-compression together reduce the execution-time to $\sim m \log^* n$, which is *almost* linear. The ideas were carefully analysed further by Tarjan [56] who finally succeeded showing that

The cost for m finds ($m \geq n$) and $n-1$ unions with balancing and path-compression is *practically* linear in m .

Our presentation of the key-ideas in the analysis will differ somewhat from Tarjan's.

4.8. Consider arbitrary programs of m finds and $n - 1$ unions which use *some* strategy for each instruction first.

Nodes can gradually climb to higher and higher levels.

A "find" (even when path-compression is used) does not change the current location of a root, but it can make that some interior nodes are pulled up to a shorter distance of the root.

It is therefore important to know where the roots will be during the algorithm, which is *entirely* determined by the unions!

4.9. Pretend that we do all unions first (obeying whatever strategy there was), and see where all nodes will eventually finish.

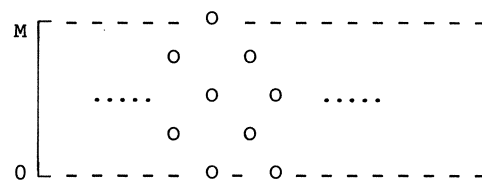
Leave all nodes where they are (hanging up in the air at their respective levels), but remove all connecting edges and start the program again.

Unions are now easy (just insert the connecting edge again), and we can concentrate entirely on what finds do in the model.

4.10. The level-number of a node v is called its rank $r(v)$.

Let the maximum rank be M .

EXAMPLE.

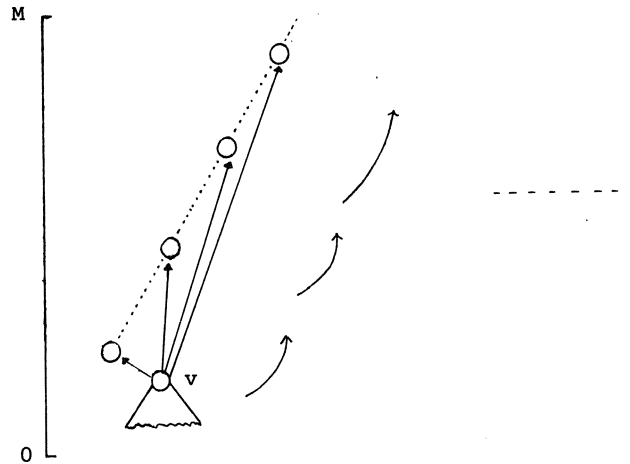


4.11. In a find with path-compression interior nodes get connected to nodes higher up, and *their (new) father always remains of higher rank*.

It means that at any moment the ranks along an upward path must form an *increasing* sequence.

4.12. We can estimate $\#F$ by counting how many times each node v ever occurring on a find-path can stretch its father-link to a higher node.

EXAMPLE.



We shall do so by counting how many times eventually occurring edges



can cross distinguished levels of a *grid*.

4.13. The crucial step in Tarjan's analysis is to work with not just one grid, but to consider $k+1$ grids

$$A(i,*) \quad (0 \leq i \leq k)$$

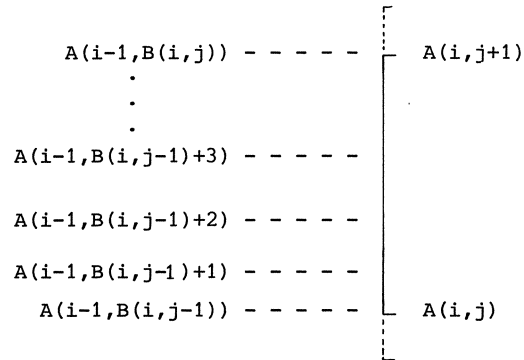
instead, where $A(i-1,*)$ is a *refinement* of $A(i,*)$. Thus we require the existence of a *monotone* function $B(*,*)$ such that

$$A(i,j) = A(i-1, B(i,j-1)) \quad (i \geq 1),$$

where we define $B(i,-1) = 0$ for consistency.

Let $a_B(i,x) = \min\{j \mid A(i,j) > x\}$. Thus $a_B(i,M)$ gives how many levels in each grid there are.

EXAMPLE.



4.14. The finest grid is $A(0,*)$ defined by

$$A(0,x) = cx \quad (c \geq 2).$$

Note that any point ever occurring on a find-path must cross this grid after being pulled up $c-1$ times.

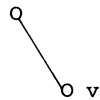
The next grids all begin with

$$A(i,0) = 0.$$

Let the number of nodes v with $A(i,j) \leq r(v) < A(i,j+1)$ be $s(i,j)$, and assume that

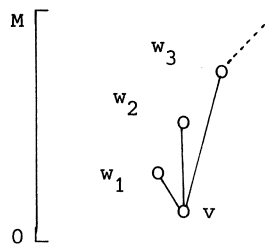
$$s(i,j) \leq \frac{f(n)}{g(i,j)} \quad (i \geq 1).$$

4.15. For each node v we should estimate how many edges



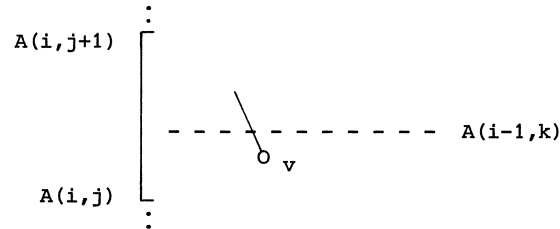
there can be on later find-paths.

EXAMPLE.



Follow v as it "bumps" into lower grids while rising first until it crosses a next grid $A(i-1,*)$ for the first time.

EXAMPLE.

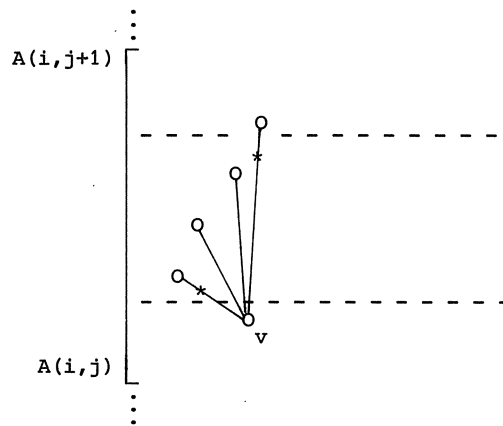


Watch v to see when it gets attached to a node across the next level of $A(i,*)$, in the meantime "charging the i -th grid" one edge whenever v gets attached to a node across yet another level of $A(i-1,*)$

- thus the i -th grid can get charged at most
 $B(i,j) - B(i,j-1)$ times -

and "charging the find" whenever v is in between but on its way

EXAMPLE



It means that we partition F into sets

$$F_0 = \{(v,w) \in F \mid \exists_j A(0,j) \leq r(v) < r(w) < A(0,j+1)\}$$

$$F_i = \{(v,w) \in F \mid \exists_j A(i,j) \leq r(v) < r(w) < A(i,j+1) \\ \text{and } \exists_\ell r(v) < A(i-1,\ell) \leq r(w)\}$$

(for $1 \leq i \leq k$)

$$F_{k+1} = \{(v,w) \in F \mid \exists_{\ell} r(v) < A(k,\ell) \leq r(w)\}$$

and divide the charges (i.e. the edge-count) between

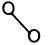
$$F_i - L_i \quad \text{and} \quad L_i \quad (0 \leq i \leq k+1),$$

where

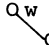
$$L_i = \{(v,w) \in F_i \mid \text{on its find-path } (v,w) \text{ is the last edge in } F_i\}$$

(thus all edges in $F_i - L_i$ are *not* last and must contribute a new crossing across the level crossed by the last).

4.16. Clearly $\#L_i \leq m$, because each find can get charged at most one last edge.

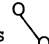
4.17. An edge  v can occur in $F_0 - L_0$ at most $c-1$ times, because after that many times v *must* be attached to a node across the first level of $A(0,*)$ (and the edges are charged to F_1 or F_2 etc.). Thus

$$\#(F_0 - L_0) \leq (c-1)n.$$

4.18. The points v on edges  v with $A(i,j) \leq r(v) < r(w) < A(i,j+1)$ in $F_i - L_i$ (see 4.15) can contribute at most $B(i,j) - B(i,j-1)$ each, and thus for each $1 \leq i \leq k$

$$\begin{aligned} \#(F_i - L_i) &\leq \sum_{j=0}^* (B(i,j) - B(i,j-1))s(i,j) \leq \\ &\leq \sum_{j=0}^* \frac{B(i,j) - B(i,j-1)}{g(i,j)} f(n) \leq h(i) \cdot f(n) \end{aligned}$$

for some function h .

4.19. For each node v there can be at most as many edges  v in $F_{k+1} - L_{k+1}$ as there are lines to cross in the k -th grid. Thus

$$\#(F_{k+1} - L_{k+1}) \leq a_B(k,M) \cdot n.$$

4.20. Adding up all estimates the total charge for m finds and $n-1$ unions

is bounded by

$$(k+2)m + cn + \sum_1^k h(i) \cdot f(n) + a_B(k, M) \cdot n ,$$

and all we have to do is choose the grids such that the expression is as small as possible.

4.21. It is important to note that the bound was derived *independent* of the strategy for unions. We can immediately obtain Paterson's result as follows.

THEOREM. *The cost for n finds and n-1 unions using only path-compression (and no balancing) is $\sim n \log n$.*

Proof. Let $B(i, j) = c \cdot (j+1)$ such that by induction we have $A(i, j) = c^i \cdot j$.

Since no balancing is used, the "best" estimate for $s(i, j)$ is just some $\frac{f(n)}{g(i, j)}$ with $f(n) = n$ and $\sum_{j=0}^* \frac{1}{g(i, j)} \sim 1$.

It follows that $h(i) = c$ for $1 \leq i \leq k$, and also $a_B(k, M) \sim M/c^k$.

By 4.20 the total cost for n finds and n-1 unions is bounded by $\sim ckn + Mn/c^k$. Let $c = 2$ and choose $k \sim \log M$. It follows that the cost is bounded by $\sim n \log M$, which could be as worse as $\sim n \log n$ for an unbalanced tree. \square

Fischer [27] proved in fact that $n \log n$ is tight to within a constant factor.

4.22. When a proper form of balancing is used we get $M \sim \log n$, and it immediately follows from 4.21 that the cost for n finds and n-1 unions is at least reduced to $\sim n \log \log n$.

Tarjan showed that one can get much better bounds from 4.20 in this case.

Remember from 3.26 that when the weighted union rule is used (or indeed any other rule for proper balancing) a node of rank r must lead a tree with at least 2^r leaves. The number of rank r nodes is therefore bounded by $n/2^r$, and it follows that

$$s(i, j) \leq \sum_{r=A(i, j)}^{A(i, j+1)-1} \frac{n}{2^r} \leq \frac{2n}{2^{A(i, j)}} \quad (= \frac{f(n)}{g(i, j)}).$$

Thus we must choose B such that we can somehow reasonably bound

$$\sum_j \frac{B(i,j) - B(i,j-1)}{2^{A(i,j)}} .$$

4.23. *Big trick*: one can choose $B(i,j) = A(i,j)$ for $j \geq 1$ and $B(i,0) = 1$. (We could have chosen even larger B's but this will do to get extremely sharp results already).

4.24. The grids we now have are defined by

$$\begin{aligned} A(0,x) &= cx \\ A(i,0) &= 0 \\ A(i,1) &= c \\ A(i,j+1) &= A(i-1, A(i,j)) \quad (i,j \geq 1) \end{aligned}$$

which yields an analog of *Ackerman's function* with an equally super-fast, not primitive recursively bounded growth-behaviour.

Observe that

$$\begin{aligned} A(1,x) &= c^x \\ A(2,x) &= c^{c^{\dots c}} \Big] |x|, \text{ etc.} \end{aligned}$$

4.25. Since we can now choose $h(i) = 2$ the total cost for m finds and $n-1$ unions is bounded by

$$\sim k(m+n) + cn + a_A(k, \log n)n$$

where one should notice that already for $k = 2$ the function $a_A(k, \log n)$ is very slowly growing.

For $k = 2$ and $c = 2$ we get Hopcroft & Ullman's result that m finds and $n-1$ unions take at most time $\sim m + n \log^* n$.

Define a "functional inverse" $\alpha(m,n)$ of the A-function in 4.24 by

$$\alpha(m,n) = \min\{j \mid A(j, \frac{m}{n}j) > \log n\}.$$

We obtain a form of Tarjan's result

THEOREM. *The cost for m finds and $n - 1$ unions using balancing and path-compression is bounded by $(m+n)\alpha(m,n)$.*

Proof. Choose $k \sim \alpha(m,n)$, and observe that $a_A(k, \log n) \leq \frac{m}{n} \alpha(m,n)$.

\rightarrow bound $\sim (m+n)\alpha(m,n)$. \square

Assuming that in practical situations usually $m \geq n$ it follows that the time-bound is

$$\sim m \alpha(m,n)$$

and that is *practically* linear. Tarjan [56] proved that the algorithm is not completely linear.

4.26. There are numerous applications of Tarjan's theorem in data-organisation, showing that all kinds of algorithms can be made practically linear with balancing and path-compression.

Hopcroft & Karp [38] proved in 1971 that the equivalence of two finite automata with at most n states can be determined in $\sim n \log n$ steps, but one may now show

THEOREM. *The equivalence of n -state finite automata can be determined in $\sim n \alpha(n,n)$ steps.*

Proof. Let $M_1 \oplus M_2 = \langle S_1 \oplus S_2, \Sigma, \delta, \{s_1, s_2\}, F \rangle$ and consider the equivalence-relation E determined by

$$p E q \iff \forall_x (\delta(p,x) \in F \iff \delta(q,x) \in F).$$

M_1 and M_2 are equivalent if and only if $s_1 E s_2$. It means that the smallest relation R satisfying

- (i) $s_1 R s_2$
- (ii) $p R q \rightarrow \forall_{\sigma \in \Sigma} \delta(p,\sigma) R \delta(q,\sigma)$

must be contained in E .

Thus M_1 and M_2 are equivalent if and only if we can build R -classes as if it were equivalence-classes without ever getting a final and a non-final state together in the same class.

Prepare classes with names $\in \{1, \dots, 2n\}$ and attributes $\in \{\text{FINAL}, \text{NON-FINAL}\}$ with the following algorithm

```

classes  $\leftarrow$  a family of  $2n$  disjoint sets of one
element each:
    seti  $\rightsquigarrow$ 

|         |       |
|---------|-------|
| state i | attr. |
|---------|-------|



queue  $\leftarrow (s_1, s_2)$ ;
while queue  $\neq \emptyset$  do
detach (p,q) from queue;
name1 := find(p);
name2 := find(q);
if name1  $\neq$  name2 then
    if name1.attrib  $\neq$  name2.attrib then
        M1 and M2 are not equivalent & signal fi
    else union(name1, name2, name1);
        attach ( $\delta(p, \sigma), \delta(q, \sigma)$ ) to queue
        for each  $\sigma \in \Sigma$ 
fi
od;
```

When M_1 and M_2 are equivalent the algorithm will produce exactly all equivalence-classes of admissible states, but if they aren't it will find the clash.

There can be at most $\sim n$ unions, and thus at most $\sim n \cdot \#\Sigma$ pairs of states will ever get onto the queue. We need to execute at most n unions and $O(n)$ finds, assuming that $\#\Sigma$ is fixed. \square

4.27. An interesting and instructive application of path-compression was found by Aho, Hopcroft, and Ullman [4] in studying programs which maintain a forest with the following instructions

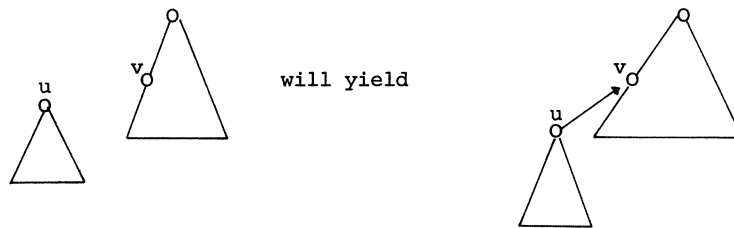
```

merge(u,v)
- to be executed only if u is a root and v is not in
  the tree with root u. The instruction makes u a son
  of v -
```


depth(v)

- determine the distance of v to the root of the tree to which it currently belongs -

EXAMPLE. A merge(u,v) applied to



where as a result $\text{depth}(u) := \text{depth}(v) + 1$.

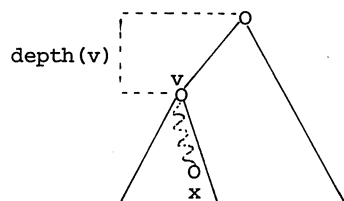
Even balancing may not seem straightforward here as it might destroy the ability to determine depth correctly if it was applied directly to the *real forest*.

In addition to the real forest we shall therefore maintain an auxiliary *shadow forest* which does support depth-instructions quickly, and we shall give a proof of

THEOREM. One can execute $n - 1$ merge- and m depth-instructions ($m \geq n$) in $\sim m \alpha(m, n)$ steps.

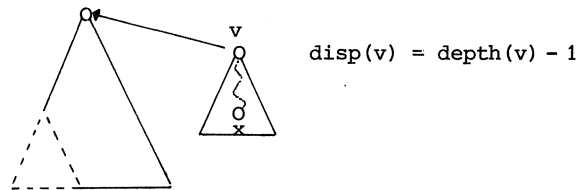
4.28. A depth-instruction looks like a find, but if you want to move an internal node v

EXAMPLE



and make it a direct son of the root (as we would in path-compression), then one must attach a *relative displacement* ("disp") to it which gets added to the depth-count when we pass v from x on our (shortened) way up to the root.

EXAMPLE.

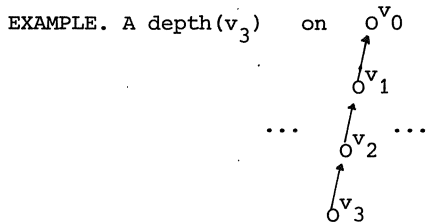


4.29. Suppose we start with a forest of n singletons, and maintain the structure with nodes of type

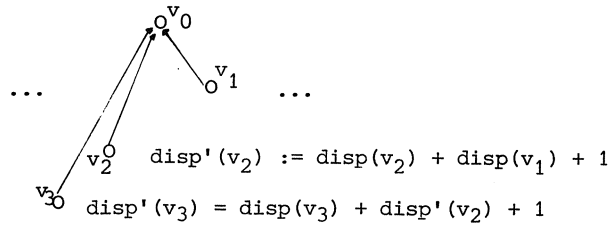
```

type node = record
    w:    number of descendant leaves;
    disp: current displacement;
    father: ↑ node;
end;
    
```

In executing a "depth(v)" we can do path-compression as in a find, provided we adjust the displacement of all nodes pulled up to the root appropriately.



should result in



with $\text{depth}(v_3) = \text{disp}'(v_3) + 1 + \text{disp}(v_0)$.

The following iterative algorithm will work

```

stack S  $\leftarrow$   $\phi$ ;
p :=  $\uparrow$ v;
if p = root then depth := p.disp & HALT fi;
while p  $\neq$  root do
  push p;
  p := p.father
od;
root := p;
p := pop S;
while stack  $\neq$   $\phi$  do
  q := pop S;
  q.father := root;
  q.disp := q.disp + p.disp + 1;
  p := q
od;
depth := p.disp + 1 + root.disp

```

(Note that the w.field of internal nodes is not up-dated and becomes unreliable, but fortunately we only need its integrity for the root.)

4.30. To execute a merge(u,v) we first check in the shadow forest that the instruction is permissible by executing

```

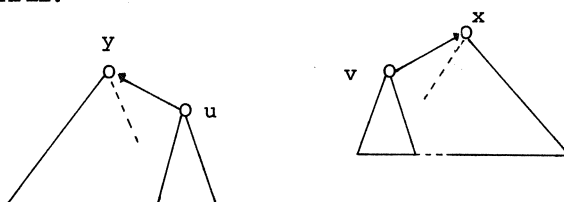
depth(u)
depth(v)

```

(with path-compression).

We should find that $\text{depth}(u) = 0$, and also get as a side-result the roots x and y of the trees where u and v currently belong. Check $x \neq y$.

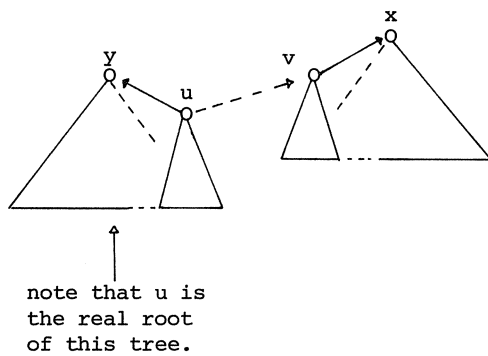
EXAMPLE.



(Note that the roots in the shadow-forest could differ from the roots in the real forest.)

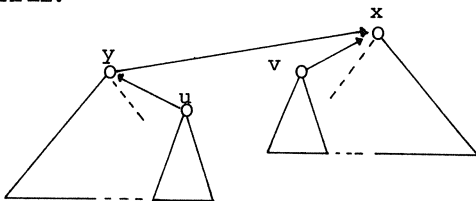
With some care the weighted union rule can be applied to the merge as reflected in the shadow forest.

EXAMPLE. The merge(u,v) should come up with a correctly displaced structure equivalent to



4.31. If $y.\text{weight} \leq x.\text{weight}$ then we may just as well make y a direct son of x, provided we adjust the displacement of y to correct for all nodes in that tree.

EXAMPLE.



where $\text{disp}'(y) = \text{disp}(y) + 1 + \text{disp}(v)$.

Note for instance that

$$\begin{aligned} \text{depth}(u) &:= \text{disp}(u) + 1 + \text{disp}'(y) + 1 + \text{disp}(x) = \\ &= \text{disp}(u) + 1 + \text{disp}(y) + 1 + \text{disp}(v) + 1 + \text{disp}(x) = \\ &= 1 + \text{disp}(v) + 1 + \text{disp}(x) \end{aligned}$$

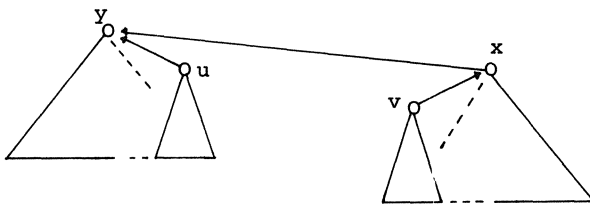
which is what is should be.

One can implement this case simply with

```
x.w := x.w + y.w;
y.disp := y.disp + 1 + v.disp;
```

4.32. If $y.\text{weight} > x.\text{weight}$ then "balancing" requires that we attach x to y instead! Adjusting the displacements is slightly trickier now.

EXAMPLE.



where now $\text{disp}'(y) = \text{disp}(y) + 1 + \text{disp}(v) + 1 + \text{disp}(x)$ and for proper correction $\text{disp}'(x) = \text{disp}(x) - \text{disp}'(y) - 1$.

We should execute

```
y.w := y.w + x.w;
y.disp := y.disp + 1 + v.disp + 1 + x.disp;
x.disp := x.disp - y.disp - 1;
```

(and the order of the instructions is important!).

4.33. The given implementation of m depths and $n - 1$ merges actually requires $\sim 2n$ extra auxiliary depth-instructions, but since $m \geq n$ the time needed is the same as for the underlying finds and unions on the shadow-forest. Apply Tarjan's theorem.

4.34. A generalisation of the techniques for merge-depth programs for other functions over trees was formulated recently by Tarjan [57]. He showed for instance that one can verify a minimal spanning tree on an n -node graph with m edges in time $\sim m \alpha(m, n)$.

5. ASSOCIATIVE SEARCH STRUCTURES

5.1. The known techniques for data-structuring can usually be combined somehow by a creative programmer into an acceptable organisation for efficiently handling records on an identifying primary key, but additional, non-trivial considerations may be needed to work it out for accessing records on *secondary* keys.

In such applications a user can issue *data-requests* (or *queries*) asking the system to produce all data-items from a file or data-base which possess a desired combination of *attributes*.

Queries are typically of the form

- list all records with property P in workspace W -

where P may be specified by some formula in relational algebra or, say, as a predicate in a calculus-based language like Codd's data-sublanguage ALPHA [18].

The process of answering data-requests is called *data-retrieval*. (One may speak of re-three-val after probing the data-base twice.)

5.2. An important task in data-base management consists of organising information such that all queries in a known language-system can be handled quickly. One may also require that in realistic applications this property must be maintained under a choice of admissible operations on the data-set.

The task usually reduces to

- resolving the nearly always present *time/*
storage trade-offs favourably -

and in close connection to it

- finding the optimal degree of *redundancy*
needed in the data-base -

5.3. There may be an advantage in keeping lists with instant answers to basic queries, and it is conceivable that in such situations one record has to occur on many lists.

Note that redundancies in the data-model can possibly be avoided in the physical data-base.

5.4. A common principle learns not to store a data-item where it is not needed, and not to duplicate a data-item when you can avoid it.

The use of *pointers* tends to be successful exactly for this reason, but one should not forget that storing many pointers may require much space also.

5.5. One may wish to avoid storing explicit pointers and try to "hack" a pointer (as a record key) in little pieces which one can search in a "small" auxiliary table first.

If at least one piece is not found then there is no point in searching the original, much larger file.

5.6. A technique of Bloom [10] suggests to maintain a long bit-string

$$\alpha_0 \dots \dots \dots \alpha_M \quad (\text{some } M)$$

and to use hashing functions h_1, \dots, h_s as follows. Mark that record k is stored by setting bits $\alpha_{h_1(k)}, \dots, \alpha_{h_s(k)}$ to 1, as if it were unique coordinates for k .

Provided $\binom{M+1}{s} > n$, there are good chances that this method can distinguish many records.

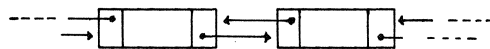
5.7. In cases where pointers are used to *link* elements together in a (linear) list, the need for extra space seems unavoidable.

EXAMPLE. Doubly linked linear lists are usually built from elements of

```

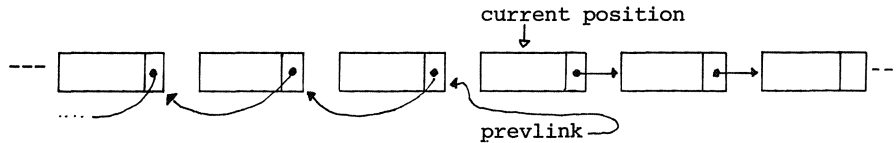
type block = record
    llink : ↑block;
    info : ...
    rlink : ↑block;
end;

```



Nevertheless (depending on the application) it may take little effort occasionally to eliminate much of the pointers even here.

EXAMPLE. Suppose a doubly linked list was chosen in order to traverse a file from left-to-right *and back*. Then one can save 50% in pointer-space by using a linked structure



One can move one block to the right with

```

if (next := pos.link; next = nil) then exit fi;
pos.link := prevlink;
prevlink := pos;
pos      := next;

```

and one can move one block back (to the left) with

```

if (next := prevlink; next = nil) then exit fi;
prevlink := next.link;
next.link := pos;
pos       := next;

```

5.8. It is common programming practice also sometimes to add a pointer into records as a means of referring to a part in memory where auxiliary information relevant to each individual record is stored.

It is probably much better to do so than to list the auxiliary information explicitly in each record in detail, especially in cases where many records *share* the same data.

Thus *pointers are used here to eliminate redundancy*, and we pay only one additional field per record (instead of many).

There are many variants of this idea; one may want to eliminate the pointers altogether and use a table at a "known" location instead in which is listed where the auxiliary data for each record can be found (but then again, at the cost of a table-search).

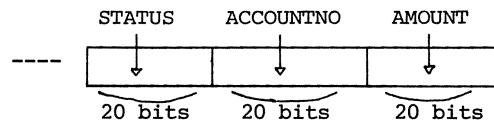
5.9. In an interesting paper Bobrow [11] pointed out very recently that the problem of locating "associated information" may be attacked by making better use of the one and only attribute of a record that comes for free: its *address*!

It is there anyway, and *it occupies no space within the record.*

The idea is that one can simply *hash the address* of a record (serving as its primary key) into a unique position of a table, and immediately find what we need.

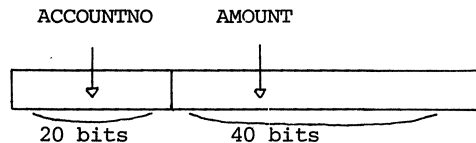
The technique of *hash-linking* can be applied in various forms.

EXAMPLE. A company may keep records



for its customers which exactly fit in one word for amounts not exceeding 20 bits of storage.

The range of AMOUNT is conceivably larger, and in programming there is always the tendency to allocate enough space to handle "the worst case". Assuming that amounts *usually* fit in 20 bits one should not allocate more words for each record, but instead *hash-link* a cell to a separate, much smaller table with "special purpose" records



whenever its amount overflows (which presumably doesn't happen very often). An "all ones"-field in the original record can signal that it is hash-linked.

As long as there are no collisions one may compress the hash-table even further because it isn't necessary to keep keys in the table (a hash-link will immediately lead to the right spot). Otherwise one might mark a record to indicate that it collided, and that one must follow a hash-link into an overflow-area where this time identifying keys are available.

5.10. The methods for locating "secondary information" in a data-file get more involved when we consider querying a data-base.

One might distinguish the following types of queries (following Rivest [53])

- exact match queries
locate the record with key k
- component match queries
locate all records in which the i-th
attribute of the key is j
- partial match queries
locate all records for which s specified
attributes match given data (and the
remaining attributes can be arbitrary)
- range match queries
locate all records for which specified
attributes are within a given range
- good match queries
locate all records which have a key
differing from k by at most a given tolerance
- best match queries
locate all records which have a key with
least distance to k

In more general, *boolean queries* one can specify records by means of an arbitrary propositional formula over the attributes.

5.11. Most queries seem to force a search through the entire data-base, and there is great need for a specific data-organisation (perhaps one for each single type of queries) to support faster retrieval. It is the object of current research in software and in hardware systems to implement *associative search* in a best possible way.

There are probably natural limits to how efficient data-retrieval can be. If one specifies more and more conditions on the attributes, then the amount of work to test each record *increases* while the actual number of records satisfying the query probably *decreases*.

5.12. Consider the following, typical *library search* problem.

Suppose records are classified on k out of many possible keywords,

and suppose we want to store records in such a way that we can later find all records with s ($\leq k$) attributed keywords quickly without having to inspect each element of the data-base individually.

How can such a goal be achieved?

Common solutions are based on an appropriate partitioning of the data-base in various *buckets*.

5.13. Apparently Gustafson [35] systematically investigated such problems first, and in 1969 he suggested the following variant of hashing for obtaining an acceptable solution.

Choose an "ordinary" hash-function

$$h: \{\text{keywords}\} \longrightarrow \{0, \dots, b-1\}$$

for some $b > k$, and build a "super" hash-function

$$H: \text{unordered } k\text{-tuples} \longrightarrow \text{b-bit numbers}$$

classifying each record with attribute-set $\{\beta_1, \dots, \beta_k\}$ in the bucket with address $\alpha_0 \dots \alpha_{b-1}$ where

$$\alpha_i = 1 \iff \exists_j h(\beta_j) = i.$$

We thus hash each key into a bit-position, and sometimes speak of *hash-coding*. Assuming a more or less uniform distribution, buckets will have $\sim N/\binom{b}{k}$ elements each.

5.14. To search for all records which contain keywords $\gamma_1, \dots, \gamma_s$ (some $s \leq k$) one should first determine the *partial* address $\alpha_0 \dots \alpha_{b-1}$ with

$$\alpha_i = \begin{cases} 1, & \text{if } \exists_j h(\gamma_j) = i \\ *, & \text{otherwise} \end{cases}$$

and then fill in the *'s with $k-s$ ones and $b-k$ zeros in all possible ways to obtain the names of all buckets which one would have to inspect only.

On the average one would have to scan through no more than

$$\sim \binom{b-s}{k-s} \cdot \frac{N}{\binom{b}{k}} = \frac{\binom{k}{s}}{\binom{b}{s}} \cdot N =$$

$$= (k-s+1) \dots (b-s) \frac{k!}{b!} N \text{ elements,}$$

which indeed tends to be a very small fraction of the entire data-base for $s \sim O(k)$.

The method is promising for small k .

5.15. We note in passing that methods for partial match retrieval based on *intersecting inverted lists* will not be discussed here (see Engles [25] or Roberts [54]).

5.16. In 1971 Rivest (see Knuth [44], Rivest [52]) suggested another, simple hash-coding algorithm for partial match retrieval in a data-base.

Suppose now that keys are ordered k -tuples, and assume it simply are length k strings over some alphabet Σ with $\#\Sigma = \sigma$.

Determine an "ordinary" hash-function

$$h: \Sigma \rightarrow \text{b-bit codes}$$

(where perhaps $2^b < \sigma$), and build a super hash-code

$$H: \Sigma^k \rightarrow \text{c-bit numbers}$$

with $c = kb$ by defining

$$H(\beta_1 \dots \beta_k) = h(\beta_1) \dots h(\beta_k).$$

There will be $\sim N/2^c$ elements per bucket on the average.

5.17. Let a partial match query be represented as a length k string

$$\gamma = \gamma_1 \dots \gamma_k$$

where for each i

$$\gamma_i \in \Sigma \text{ ("specified")} \text{ or } \gamma_i = * \text{ ("unspecified")},$$

and suppose there are s stars.

As in Gustafson's method, again, we only have to search in the $(2^b)^s$ buckets with a name matching γ in the specified positions.

It means that on the average we must inspect

$$\sim \frac{N}{2^c} \cdot 2^{bs} = 2^{b(s-k)} N (= \sigma^{s-k} N)$$

elements, which is

$$\sim 2^{s \cdot \log \sigma} = \sigma^s \sim N^{s/k}$$

for $b \sim \log \sigma$ and a full data-base with $N \sim \sigma^k$.

5.18. Provided enough extra storage is available there is a simple variant of Rivest's algorithm from which one may expect an even better behaviour.

Partition each length k key in m fields of base- σ numbers of $\frac{k}{m}$ digits

$$\boxed{A_1} \quad \boxed{A_2} \quad \dots \quad \boxed{A_{k/m}}$$

and store each record with such a key in bucket A_j of the j -th data-base for each $1 \leq j \leq m$.

We thus list each record many times in our data-model, but it will appear that *the redundancy will pay off in greater retrieval-efficiency.*

5.19. In a partial match query with a total of s stars there must be at least one field with $\leq s/m$ stars. Among these choose the field with the fewest number of stars (say field j).

One may expect that the j -th data-base will have the elements we search for concentrated in the smallest number of buckets.

Roughly, there are $\sim N/\sigma^{k/m}$ elements per bucket and we have to scan through

$$\sim N/\sigma^{k/m} \cdot \sigma^{s/m} = \sigma^{s-k/m} \cdot N$$

elements. This compares favourably with the bound $\sim \sigma^{s-k} \cdot N$ in 5.17

(although a more precise analysis would be needed to support the apparent conclusion).

5.20. One may similarly develop other methods, and it becomes of interest to study how the various hash-coding algorithms for partial-match retrieval compare in efficiency.

Rivest [52] (also [53]) answered this question for a specific class of hashing algorithms, and we shall give a simplified approach to make his result understandable.

5.21. Each hash-code storage schema consists of a super hash-function

$$H: \Sigma^k \longrightarrow \text{a finite set of buckets}$$

(the algorithm of 5.18 is *not* in this category).

The relation "having the same H-value" gives a partitioning of Σ^k , and the various equivalence-classes may be called *H-blocks*.

Now consider how we answer partial-match queries γ with s stars. The algorithm must return the names of all H-blocks which contain at least one record matching γ in the specified positions.

Rivest suggested to "measure" H by the average number of H-blocks needed to answer s -star queries.

5.22. Suppose there are Q s -star queries.

On the average we need to inspect a number of blocks equal to

$$\begin{aligned} & \frac{1}{Q} \sum_{\gamma} (\# \text{ H-blocks with a } \gamma\text{-match}) = \\ & = \frac{1}{Q} \sum_B (\# \text{ s-star } \gamma\text{'s touching B}), \end{aligned}$$

where B ranges over all H-blocks.

Each s -star query occurs at least once in this summation, and it is counted more than once if and only if it hits more than one block.

It is conceivable (and Rivest proved it for "balanced" hash-coding algorithms) *that the average behaviour will come out best once the "answer" to as many partial-match queries as is practical is contained in one block.*

It means that preferably H-blocks must allow a decomposition into

$$\Sigma_1 \times \Sigma_2 \times \dots \times \Sigma_k$$

where for each i : $\Sigma_i = \Sigma$ or $\#\Sigma_i = 1$.

It holds for the algorithms of Gustafson and Rivest, and will hold for later algorithms also.

5.23. Based on these heuristics it seems promising to consider a general type of (hash-coded) *partial-match file design* which consists of a table linking *independent key-patterns* and bucket-addresses

PATTERN	BUCKET
⋮	⋮
*1**001*1*	25
⋮	⋮

where

- (i) each row contains s stars and $k-s$ symbols from Σ ;
- (ii) any two key-patterns (rows) differ in at least one specified position (which means that buckets will be disjoint).

Technically each entry in the table is an s -star partial-match query, but obviously other queries can be answered from it also. If no entry matches a record-query then no records of that type are in the data-base.

The way stars are distributed in the key-patterns can strongly influence the average performance of the file-design. Rivest [53] and later Burkhard [14] (see also Bentley & Burkhard [9]) studied several special designs with an alleged approximately optimal behaviour.

5.24. One may assume that entries in the table (as strings over $\Sigma \cup \{*\}$) are listed in *radix-sorted* order (a general storage method suggested by Hildebrand & Isbitz [36] as early as 1959). It is likely that the best results will be implied when the entire data-base is stored in a compatible manner.

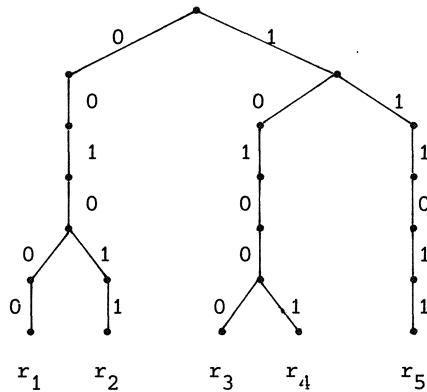
The idea is that in such a scheme it becomes unnecessary to store the table separately, and its structure (or rather, the structure of its buckets) is implied by a straightforward digit-wise search where each "*" is interpreted as a "branch in both directions"-instruction.

5.25. Assume that $\Sigma = \{0,1\}$ (for simplicity).

When length- k strings of 0's and 1's are entered into a tree where at each node "0" is interpreted as "left" and "1" as "right", each record-key becomes a code for the path to a unique leaf where the corresponding record-address can be stored.

The idea is independently due to de la Briandais [23] and Fredkin [30], and the resulting storage structure is known as "*trie memory*" (pronounced usually as "try-memory", although the editor of the section "Techniques" in the 1960-issue of the Communications pointed out in a footnote to Fredkin's paper that "trie" apparently was derived from the word re-TRIE-val).

EXAMPLE. The trie corresponding to $\{001000,001011,101000,101001,111011\}$ is



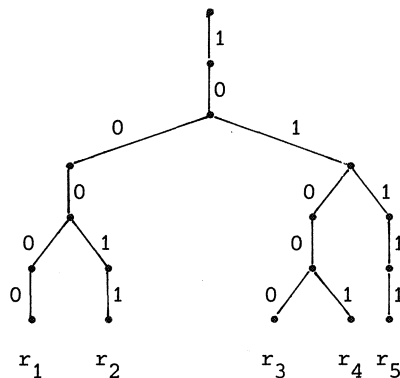
Tries are often used when there are records with variable-length keys also.

5.26. Despite the clarity of the file-structure, many programmers observed that tries can be somewhat inefficient in their use of storage.

One reason is that in entering keys we may not have made sufficient use of the observation that some strings could have had identical segments of bits which we should have followed first.

DEFINITION. A π -trie (where $\pi \in S_k$) is the storage-structure resulting after entering keys γ with bits ordered as $\gamma_{\pi(1)}\gamma_{\pi(2)}\dots\gamma_{\pi(k)}$.

EXAMPLE. The 341256-trie for $\{001000,001011,101000,101001,111011\}$ (see 5.25) is



It has only 12 internal nodes, as compared to the 123456-trie in 5.25 which

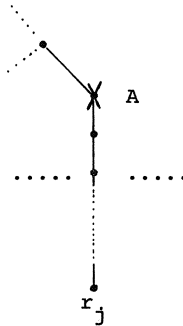
has 16.

It is unlikely that there will be a polynomial-time bounded algorithm to find an optimal π -trie for a file, since Comer & Sethi [20] recently proved the following interesting result

THEOREM. *The problem to determine a π -trie with the smallest number of internal nodes for a file is NP-complete.*

5.27. Given tries (or π -tries) can be *optimized* in various, simple ways.

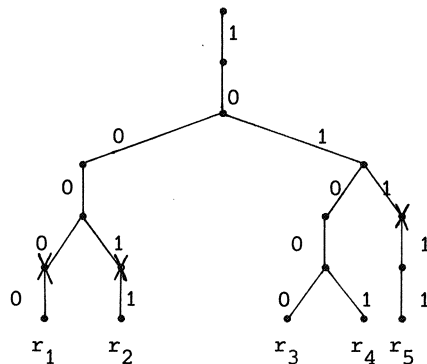
Observe, for instance, that one could eliminate chains of the form



(which Comer & Sethi call *leaf-chains*), and *prune* the tree in A at the earliest moment that r_j is uniquely identified.

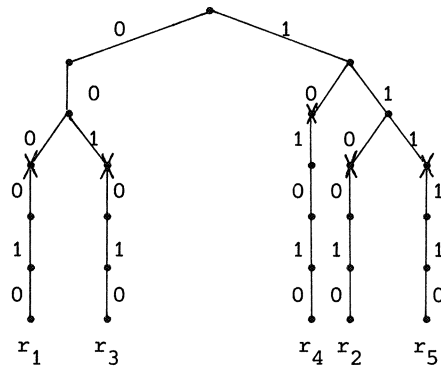
The idea is that the remaining bits of the key need not take up trie-storage and can be checked elsewhere.

EXAMPLE. The pruned 341256-trie for {001000,001011,101000,101001,111011} is



It seems to require an entirely different strategy to determine the best possible *pruned* π -trie for a file.

EXAMPLE. The 651234-trie for {001000,001011,101000,101001,111011} is

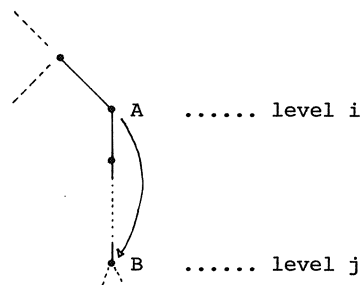


It has as many as 21 internal nodes, but the pruned version has only 5.

Comer & Sethi [20] proved also

THEOREM. *The problem to determine a pruned π -trie with the smallest number of internal nodes for a file is NP-complete.*

5.28. A next step would to collapse long chains



in the interior of the trie also, by adding an instruction at node A to skip over bits $\gamma_{\pi(i+1)}, \dots, \gamma_{\pi(j)}$ and to continue immediately with bit $\gamma_{\pi(j+1)}$.

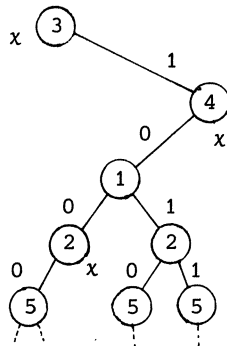
The idea was apparently first proposed by Morrison [48] in 1968, and

implemented in his system PATRICIA. Testing the skipped bits of the key afterwards avoids the need for ineffective trie-storage, and the technique leads to the most compact π -trie.

The problem to determine an "optimal" compacted π -trie is meaningless if the number of internal nodes is the only measure, and one should now also take the degree of balancing into account.

5.29. In a π -trie we could have inscribed the tested bit-positions in the corresponding nodes.

EXAMPLE.



Clearly in partial-match retrieval some bits in the input will be unspecified, and when a * is found one must descend down two branches (if there are two) simultaneously (and so on ..), thus implicitly defining what bucket is searched.

5.30. Carrying this idea one step further, Burkhard [15] recently formulated a type of generalised trie which seems an appropriate structure for partial-match retrieval from large indexed files.

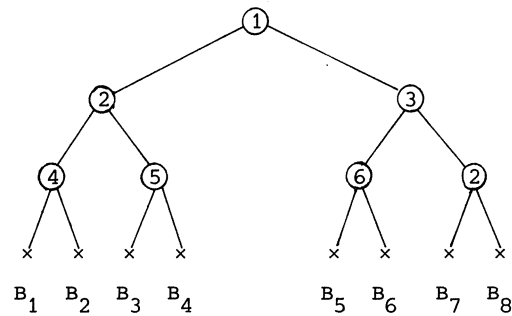
A (k,s) PM-trie is a complete binary tree with 2^{k-s} leaves in which

- (i) each leaf corresponds to a bucket;
- (ii) each internal node is labeled by one out of k possible bit-positions;
- (iii) on no path from the root to a leaf there is a bit-position which is tested more than once.

(Thus each such path is like an s -star query.)

Technically each (k,s) PM-trie directly implies a (k,s) partial-match file-design in the sense of 5.23, but not conversely.

EXAMPLE.



is a $(6,3)$ PM-trie.

To test how good a PM-trie is we shall measure how many buckets must be searched *in worst case* to answer a t -star query ($t \leq k$).

The result is likely to depend on the way tests are distributed over the tree, and the task is to find "good" (k,s) PM-tries for use on files with k -bit entries.

EXAMPLE. The query $***010$ in the given $(6,3)$ -trie leads into 5 buckets.

The query $*1****$ leads into 5 buckets also, despite the fact that it is "almost unspecified".

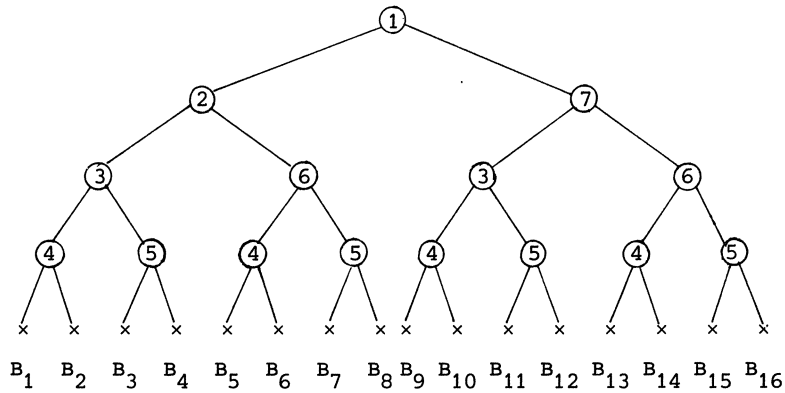
5.31. In early 1975 Burkhard [14] first formulated and analysed a type of "good" $(2n+1,n)$ PM-tries. We shall present a modified version which was shortly later suggested by Dubost & Trousse [24].

DEFINITION. A $(2n+1,n)$ PM-trie is called a BDT-trie of order n if and only if

- (i) the root (at level 1) has label 1;
- (ii) for each node in level i ($1 \leq i \leq n$) the left-son is labeled $i+1$ and the right-son is labeled $2n+2-i$.

Observe that the tested bit-positions at the nodes may be computed while descending down the tree very easily.

EXAMPLE. The BDT-trie of order 3 is

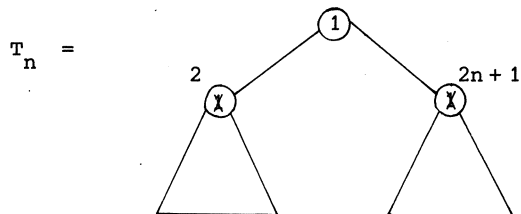


Observe that in a BDT-trie of order n all paths have length $n + 1$, and there is an obvious regularity in the assignment of tests.

The labels of two brothers always add to $2n + 3$.

Let T_n denote the BDT-trie of order n .

LEMMA. For all $n \geq 1$

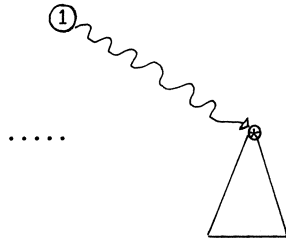


where each marked sub-tree is T_{n-1} with 1 added to all labels (and $2n$ to the root of the right copy).

5.32. Formulas for the worst case performance of BDT-tries were first given by Burkhard [14]. We shall outline an interesting proof-technique due to Dubost and Trousse.

The problem is how to visualize what buckets can be entered in a t-star query.

Suppose we first follow the trie-path until we hit the first bit-test which probes a star.



It means we narrowed the search for buckets to a small sub-trie.

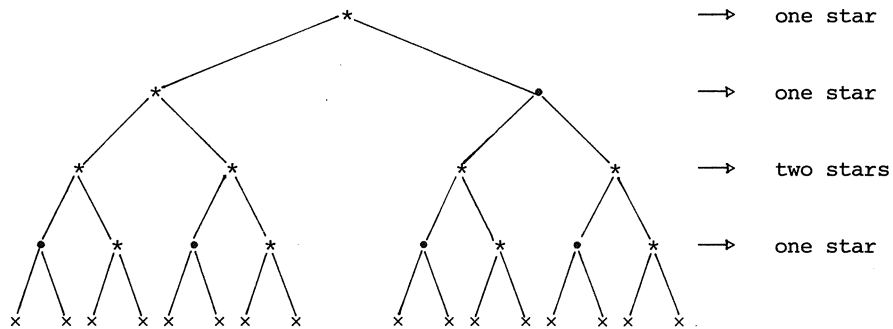
However, each sub-trie (essentially) is a BDT-trie for sub-keys of the original file.

Thus we reduced the task to finding the maximal number of buckets reached in a *t*-star query with a * in position 1 (it will turn out to be an increasing function in *n*).

5.33. A *consistent assignment* with *t* stars in a BDT-trie is a complete marking of all (internal) nodes with "stars" and "dots" such that

- (i) the root is marked with * .
- (ii) If a left-son in level *i* is marked with a *, then so are all left-sons in level *i* (which is reasonable since they probe the same bit-position).
- (iii) If it happens, then the left-sons in level *i* are said to *contribute one star*.
- (iv) If a right-son in level *i* is marked with a * then all are, and the right-sons are said to contribute one star also.
- (v) The number of contributed stars adds up to *t*.

EXAMPLE.



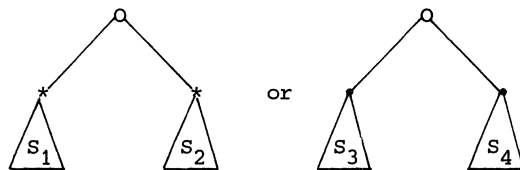
The search-structure in a BDT-trie of order 3 for a query ***** (where

dots mean *defined* bits). No matter what, 9 buckets will be visited on any 5-star query of this pattern.

The reason for introducing this concept is that *each t-star query implies a consistent assignment with t stars and conversely.*

5.34. For symmetry-reasons we may assume that if a level $i > 1$ contributes one star then it is contributed by the left-sons in that level.

Observe further that if two brothers are identically marked

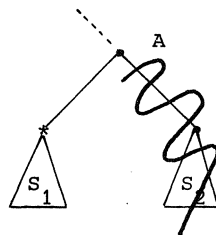


then the entire, corresponding sub-trees are identical in marking (i.e., $S_1 \equiv S_2$ and $S_3 \equiv S_4$ respectively).

5.35. Given a consistent assignment one can easily determine the largest possible number of buckets generated by a corresponding query.

Collapse and prune the tree according to the following rules.

- for each dotted node with not identically marked sons, purge the sub-tree of the dotted son



Because the bit at A is specified we would enter either S_1 or S_2 . The sub-trees are identically marked except at their root, so entering S_1 will definitely lead into a larger number of buckets.

- for each dotted node with two identically marked sons, purge one of its subtrees also.

The left- and right-subtrees are identical, and since we enter only one on the specified bit we need only count buckets in one sub-tree.

One could go one step further and collapse all chains of dotted nodes left.

The number of leaves in the resulting tree is exactly the number of buckets to be visited.

5.36. *It follows that the number of buckets visited by t-star queries of the same *-pattern is the same.*

5.37. Consider first how many buckets result when in the original, consistent assignment each level contributes one star.

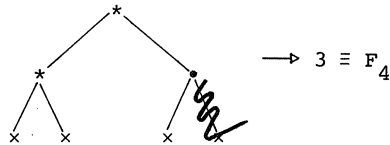
LEMMA. *For BDT-tries of order n this number is F_{n+3} .*

Proof. By induction.

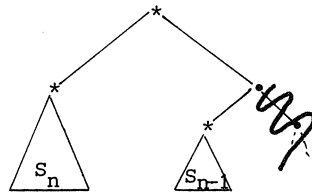
For $n = 0$ we have



and for $n = 1$ we have (remember: only one star per level)



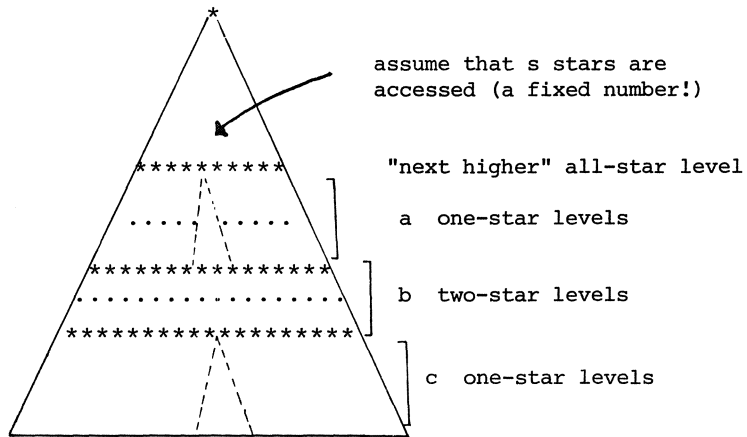
For $n + 1$ the marked trie-structure S_{n+1} of this form decomposes as



and the number of buckets adds up to $F_{n+3} + F_{n+2} = F_{n+4} \equiv F_{(n+1)+3}$. \square

5.38. Next we consider consistent assignments in which each level contributes at least one star, and the number of levels indeed contributing two stars is assumed "fixed".

Examine the "lowest" layer of two-star levels in such an assignment



The number of buckets visited is

$$s \cdot 2^b \cdot F_{a+3} \cdot F_{c+3}$$

(using the previous lemma). Assuming b fixed, and observing that

$$s \cdot 2^b \cdot F_{a+3} \cdot F_{c+3} \leq s \cdot 2^b \cdot F_3 \cdot F_{a+c+3}$$

it follows that *the number of buckets is maximal if $a = 0$* . We immediately obtain

LEMMA. *In a consistent assignment with one- and two-star levels the number of buckets is largest if all two-star levels occur together and precede all one-star levels.*

5.39. In a general consistent assignment there can be dotted levels also, and we shall now analyse where they have to occur for obtaining the largest number of buckets.

Our argument here differs from the proof of Dubost and Trousse who

tried to "collapse" non-contributing levels in the "worst possible" way.

Let

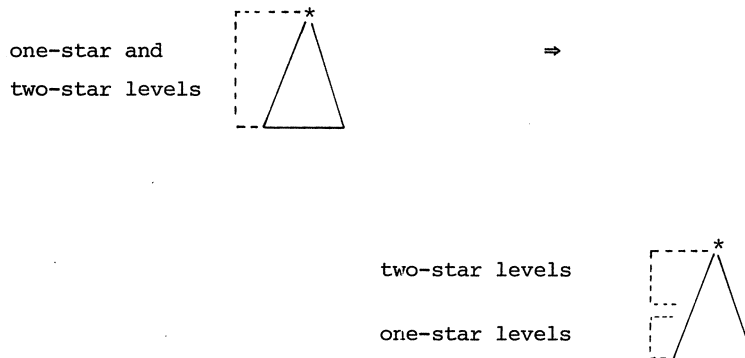
$w_n(t)$ = the maximum number of buckets visited with a
t-star query on a BDT-trie of order n.

5.40. Observe first that the number of buckets for an assignment does not change if we add non-contributing levels at the bottom of the trie.

The idea is to make a given assignment "worse" by an exchange operation on the levels *without changing the order of the trie and without changing how many contributed stars there are*. (Thus the assignment remains a t-star query).

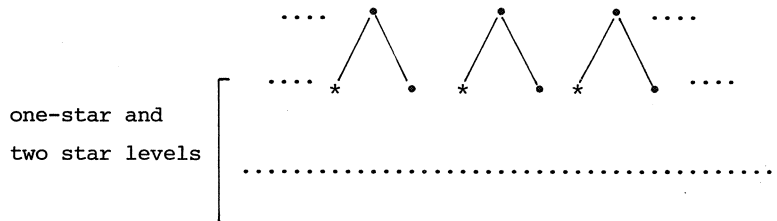
Do the following reductions whenever levels of the indicated form occur.

(i)



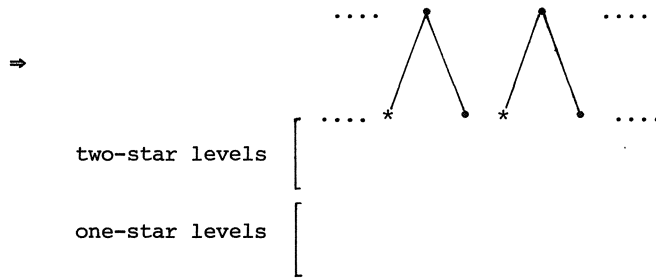
This is the simple transformation we proved in 5.38. It can be applied to subtrees also, *provided* the same transformation brings profit in all subtrees at this level (viz. when it is a two-star level).

(ii)

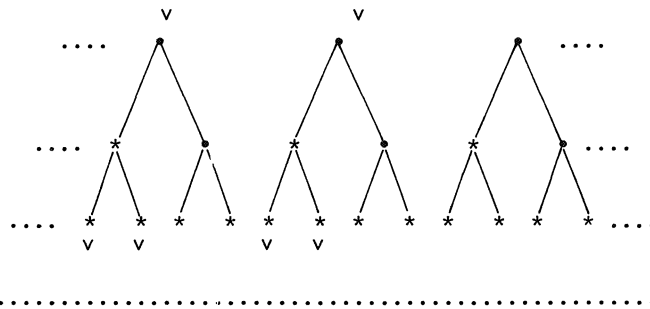


It is easy to see that the worst case occurs if from the dotted level we

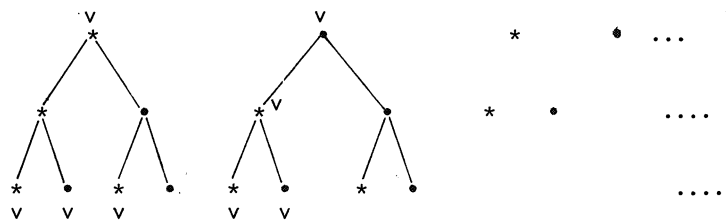
always proceed to the starred son. Thus we can apply (i) in the subtrees and should have all two-star levels immediately follow the initial one-star level



(iii)

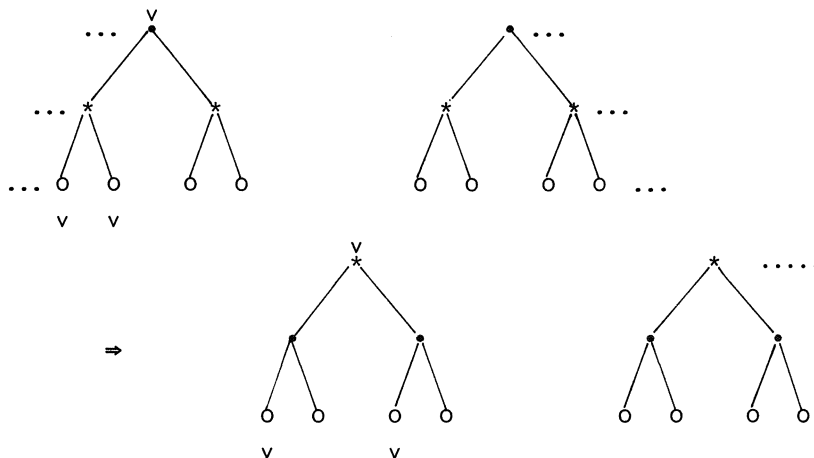


From each node in the dotted level we reach two stars two levels lower. Now borrow a star from the two-star level and make the dotted level into a one-star level:



Considering how one can enter the formerly dotted level to get a worst case now, it follows that at least as many buckets are entered now as we entered before.

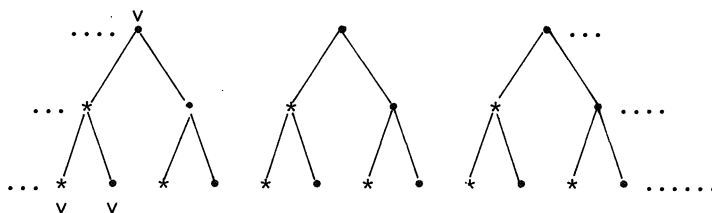
(iv)



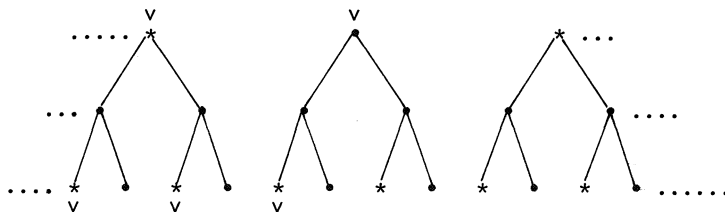
A two-star level and a dotted level can always be interchanged. (With the previous transformations it now follows that for the worst case to occur all two-star levels must occur at the top of the tree).

(v) Anywhere in the tree the number of consecutive dotted levels can be reduced to 1, provided we add the deleted levels at the bottom (to preserve the order of the tree)

(vi)

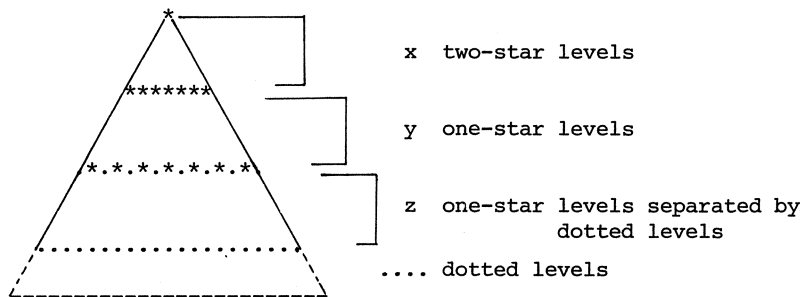


It is always "better" to have the dotted level occur in between the two one-star levels.



Note that this transformation together with the other tries to sift the dotted levels down and use them to separate as many one-star levels as it can.

5.14. If we perform these operations then we reduce a given assignment to a kind of normal form for the worst case:



where $1 + x + y + 2z \leq n + 1$ (the number of levels)

and $1 + 2x + y + z = t$ (the number of stars).

5.42. We can now combine all information and prove Burkhard's theorem

THEOREM. For t -star queries on a BDT-trie of order n

$$w_n(t) = \begin{cases} 2^t & \text{for } 0 \leq t \leq n+1 - \lceil \frac{n}{2} \rceil \\ 2^{n+1-t} F_{2t-n+1} & \text{for } n+1 - \lceil \frac{n}{2} \rceil \leq t \leq n+1 \\ 2^{t-n-1} F_{2n+4-t} & \text{for } n+1 \leq t \leq 2n+1. \end{cases}$$

Proof. For the worst case stars must be distributed as in the tree of 5.41. We reach 2^x nodes in the last all-star level (which could be the root if $x = 0$), F_{y+1} stars and F_y dots from each of these in the last of the consecutive one-star levels, and in the last part of the tree we get into 2^{z+1} buckets from a star and 2^z buckets from a dot.

It follows that we must maximize the total of

$$\begin{aligned} 2^x \cdot F_{y+1} \cdot 2^{z+1} + 2^x \cdot F_y \cdot 2^z &= 2^{x+z} (2F_{y+1} + F_y) \\ &= 2^{x+z} \cdot F_{y+3} \text{ buckets,} \end{aligned}$$

under the condition that

$$\begin{aligned}x + y + 2z &\leq n \\2x + y + z &= t - 1. \quad \square\end{aligned}$$

5.43. In a subsequent generalization Burkhard [14] formulated another, more flexible type of PM-trie in which he abstracts the essential features that make BDT-tries work.

DEFINITION. A $(2n+1, n)$ PM-trie is called a G-trie of order n if and only if

- (i) any label j ($1 \leq j \leq 2n+1$) occurs in a unique level;
- (ii) brothers carry different labels.

Burkhard proves that the expressions of 5.42 are upperbounds for order n G-tries in general.

Burkhard [16] also develops similar kinds of tries for files where keys are over an arbitrary alphabet Σ .

6. PATTERN-MATCHING

6.1. One may describe *pattern-matching* as the task of finding one or more sub-structures which meet a certain specification in given larger structures of some kind (usually in on-line manner).

Pattern-matching operations occur in lexical analysis, in searching library-files, in text-editing and in symbol-manipulation, and although it may seem easy we often need non-trivial data-organisations to perform the operations quickly without undesirable overhead.

6.2. Pattern-matching is undoubtedly the most powerful feature of the SNOBOL 4 programming language, and it often occurs in conjunction with *replacement* (see Griswold [34]).

Although we shall restrict the algorithms to strings and one-dimensional patterns, one could also consider pattern-matching problems for arrays and trees (see Karp, Miller, and Rosenberg [41]).

6.3. Let all strings considered be over the alphabet Σ with $\#\Sigma = \sigma$.

The simplest algorithm to find all length d substrings of

$$x_1 x_2 x_3 \dots x_n \quad (n \geq d)$$

makes use of a *shift-register* containing

$$\text{state} = |x_{i+1}| \sigma^{d-1} + |x_{i+2}| \sigma^{d-2} + \dots + |x_{i+d}|$$

$$\longleftrightarrow \dots x_{i+1} x_{i+2} \dots x_{i+d} \dots$$



current position of scanner

Suppose that σ^d buckets $B_{\frac{\sigma^d-1}{\sigma-1}}, \dots, B_{\frac{\sigma^d-1}{\sigma-1}}$ are allocated, one for each

possible string of length d . Then one can collect all information from x in one scan as follows

```

state := |x1|σd-1 + |x2|σd-2 + ... + |xd|;
enter 1 in Bstate;
position := d;
while position < n do
  position := position + 1;
  state := state·σ + |xposition| (mod σd);
  enter position - d+1 in Bstate
od;

```

In the end $B_{|y_1| \sigma^{d-1} + \dots + |y_d|}$ will contain all positions $i+1$ such that

$$x_{i+1} \dots x_{i+d} = y_1 \dots y_d.$$

The algorithm essentially traverses a finite state automaton with a very special transition-structure (known as a de Bruijn graph in combinatorics) which could have been computed and stored in advance. The approach is of reasonable complexity only if arithmetic modulo σ^d can be performed quickly.

6.4. More common pattern-matching algorithms of low complexity develop and use *positional information* of non-numeric nature.

A typical procedure for classifying and locating all length d substrings of x makes use of the following equivalence-relation

positions i and j of x are s -equivalent if
and only if

$$x_i \cdots x_{i+s-1} = x_j \cdots x_{j+s-1}.$$

Our representation of equivalence-classes will differ from Karp, Miller, and Rosenberg [41], and simplify the final algorithm somewhat.

6.5. Let the *name* of an equivalence-class be the smallest element it contains. $\text{CLASS}[i] = j$ means that position i belongs to the s -equivalence class named j .

The elements of an equivalence-class will be *linked* together in *ascending order* using NEXT. The end of a NEXT-chain will be marked by a negative integer which yields the beginning address of another class-chain, and 0 if all classes have been linked together.

One could enumerate all positions in equivalence classes with

```

position := 1;
repeat
alpha := CLASS[position];
print position;
while NEXT[position] > 0 do
position := NEXT[position];
print position
od;
end of class alpha;
position := - NEXT[position];
until position = 0;

```

6.6. The problem is how one can construct the s - equivalence classes in this representation quickly. The idea is an inductive approach based on the following result.

LEMMA. For all $r \leq s$ and i and j , $i \sim_{r+s} j \iff i \sim_s j \ \& \ i+r \sim_s j+r$.

Proof. By overlap. \square

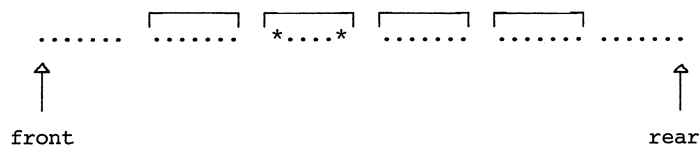
We can now show

THEOREM. *Given the representation of s -equivalence classes one can construct the $r+s$ -equivalence classes for any $r \leq s$ in linear time.*

Proof. Assume that we have n queues Q_1, \dots, Q_n available (one for each equivalence class).

Read out all s -equivalence classes (which yields the elements per class in ascending order) and attach each position i so encountered at the rear of Q_j , where $j = \text{CLASS}[i+r]$ (provided $i+r+s-1 \leq n$).

In the end each Q_j is either empty or made up of blocks



where each block simply consists of the items attached to Q_j while reading out an entire s -class. Note that each block can be identified because it is the longest stretch of elements with the same particular CLASS-value.

Observe now that i and j are in the same block if and only if $i \sim_s j$ and $\text{CLASS}[i+r] = \text{CLASS}[j+r]$, thus by the lemma if and only if $i \sim_{r+s} j$. Thus the blocks exactly are the $r+s$ -equivalence classes.

One can now obtain the representation of $r+s$ -classes by reading out the queues block after block. The first element of a block by definition is its name (since it is smallest), and it can be retained and put into CLASS while passing through the block. The classes can be linked in the order in which the queues are emptied. (Start with the block containing 1, which must occur up-front somewhere.) \square

6.7. It follows that

THEOREM. *One can construct the representation of s -equivalence classes for x in $\sim \sigma \cdot n + n \log s$ steps.*

Proof. The 1-equivalence classes follow by a simple scan of x in $\sim \sigma \cdot n$ steps. (There is one chain for each distinct symbol of x .) If the representation of the r -equivalence classes is known, then we can construct the $2r$ -equivalence classes in linear time. Thus the following algorithm will do

```

r := 1;
record the r-classes;
while 2r < s do
r := 2r;
record the r-equivalence classes
od;
record the (s-r)+r-classes from the r-classes;

```

□

6.8. To find all substrings of length d in x with the given algorithm works pleasantly in $\sim n \log d$ steps (assuming a fixed alphabet). A direct search for the classes with more than one element gives all repeated patterns of length d in about the same number of steps.

The same principle as in 6.7 can be used to show

THEOREM. *One can find all longest repeated substrings of x in $\sim \sigma \cdot n + n \log n$ steps.*

Proof. Karp, Miller, and Rosenberg [41] observed that the problem is equivalent to finding the *largest* r such that there is at least one r -class with more than one element.

Determine the 1-classes first in $\sim \sigma \cdot n$ steps, and then find r with binary search after a usual "doubling" procedure has located a lower bound d and upper bound $2d$ for it. □

One can develop similar algorithms based on s -equivalence to determine for instance whether x is of the form y^k for some $k \geq 2$ and y .

6.9. The idea of positional equivalences is clearly of great help in solving pattern-matching problems, but it seems to be efficient only when we search *all* patterns of some kind. Note that the algorithm in 6.8 does *not* simplify (for instance to complexity $\sim \sigma n$) in case we just want *any* longest repeated substring.

Another draw-back is that the algorithms become more inefficient (in storage and search-time) in case we must consider and "remember" several s -equivalences at the same time.

6.10. Weiner [65] proposed a very interesting idea for extracting exactly all the information you need to know from the various s -equivalences.

Let us assume that x always has a fixed marker at the end which is not occurring elsewhere

$$x_1x_2\dots x_n\$.$$

DEFINITION. y is called the *substring-identifier* for position i if and only if it is the smallest u with the property that u occurs nowhere else in x except at position i .

EXAMPLE. The substring-identifiers for the various positions in

bbababba\$

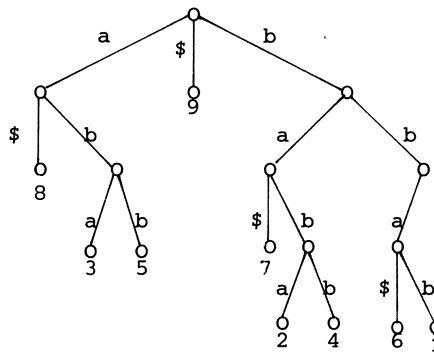
are

1	bbab
2	baba
3	aba
4	babb
5	abb
6	bba\$
7	ba\$
8	a\$
9	\$

The trick of adding \$ at the end of x guarantees that each *position* has a *substring-identifier*, and one can observe that for each i the *substring-identifier* is precisely equal to the length- r substring at i for the smallest r such that there are no other positions in the r -equivalence class containing i .

6.11. Clearly no substring-identifier can be a proper prefix of another, and one can enter all substring-identifiers consistently in a $\sigma+1$ -ary *position-tree*.

EXAMPLE. The position-tree for bbababba\$ (the previous example) is



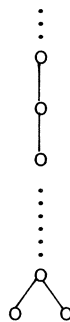
The path with y ends at the leaf with label i if and only if y is the substring-identifier of position i .

6.12. With the algorithm of Karp, Miller, and Rosenberg it follows that

THEOREM. The position-tree for x can be constructed in $\sim \sigma \cdot n \cdot L$ steps, where L is the length of the longest substring identifier.

Proof. Construct the s -equivalence classes for $s = 1, 2, \dots, L$. Each time a new relation is constructed, record the current s -value for all positions which now form a singleton equivalence class. Use the table afterwards to enter the corresponding substring identifiers in the tree in $\sim \sigma L$ steps per entry. \square

The theorem is not the best possible result though. A position-tree with many internal nodes (as many as $\sim n^2$ are possible even though it can only have n leaves) must contain long straight chains



which one can compress, resulting in a *compacted position-tree* of only $\sim O(n)$ nodes.

Weiner [65] proved that one may construct the compacted tree directly in only *linear* time. At the SWAT-meeting in Iowa City it was promptly called "the algorithm of 73", because it took that many pages to explain. (Easier descriptions are now available, see Aho, Hopcroft, and Ullman [3].)

6.13. Once the position-tree is known all kind of pattern-matching questions can be answered quickly.

Consider the problem of determining a longest repeated substring y of x . Any such y must be the *proper* prefix of some substring-identifier. Thus one must choose a y which leads to a lowest interior node in the position-tree (and any such y qualifies).

EXAMPLE. A longest repeated substring of `bbababba$` (see 6.11) is `bab`.

To find a longest common (contiguous) substring of x and y one can construct the position-tree for

$$x \text{ \textasciitilde } y \text{ \$}$$

and solve the task by searching for a lowest interior node with at least one son representing a position in x and one son representing a position in y .

6.14. The position-tree can be used also to solve the authentic pattern-matching problem to determine whether given string y is a substring of x .

Let

$$y = y[1]y[2] \dots y[m]$$

$$x = x[1]x[2] \dots x[n]$$

with $m \leq n$.

Before we see a direct algorithm for this problem it is of interest to look at the method SNOBOL 4 uses

```

position := 0;
check := 1;
while position < n - m + 1 do

```

```

while x[position+check] = y[check] do
  check := check + 1
until check > m;
if check > m then match found & exit fi;
position := position + 1;
check := 1
od;

```

The method may require in worst case $\sim n \cdot m$ steps. Morris & Pratt [47] apparently first realized that these is a *linear* time pattern-matching algorithm. At the same time ~ 1970 the result followed from a general, linear time algorithm for simulating arbitrary 2-way deterministic pushdown-automata found by Cook [21], which prompted Knuth & Pratt to write a report entitled "automata-theory can be useful" [42].

6.15. Consider how one might improve the previous algorithm.

Suppose the algorithm has been in action for a while

$$\begin{array}{ccc}
 \dots\dots \overbrace{x_{k-j+1} \dots x_k} & & \text{(match)} \\
 & & \\
 Y_1 \dots\dots\dots Y_j & &
 \end{array}$$

and now finds that $x_{k+1} \neq Y_{j+1}$

$$\begin{array}{ccc}
 \dots\dots\dots \overbrace{\dots\dots\dots x_k} & \downarrow & x_{k+1} \\
 & & \\
 Y_1 \dots\dots\dots Y_j & & Y_{j+1}
 \end{array}$$

Instead of trying again from the very beginning one should *continue at the largest i such that*

$$\begin{array}{ccc}
 \dots\dots\dots \overbrace{\dots\dots\dots x_k} & \downarrow & x_{k+1} \\
 & & \\
 Y_1 \dots\dots\dots Y_i & Y_{i+1} & \text{with } x_{k+1} = Y_{i+1}
 \end{array}$$

which is exactly the largest i such that

$$\begin{array}{ccc}
 Y_1 \dots\dots\dots \overbrace{\dots\dots\dots Y_j} \dots & & \\
 & & \\
 Y_1 \dots\dots\dots Y_i & Y_{i+1} & \text{with } x_{k+1} = Y_{i+1}
 \end{array}$$

Candidate i 's can therefore be determined based solely on y .

6.16. The following concept now appears relevant.

DEFINITION. The *overlap-identifier* $h(j)$ of position j in y is the largest $i < j$ such that

$$\begin{array}{c} \overbrace{\dots\dots\dots Y_j} \\ \underbrace{\dots\dots\dots Y_i} \end{array}$$

Observe that always $h(1) = 0$. Let $h(0) = 0$ for consistency.

EXAMPLE. The overlap-identifiers of babaaba are

b	a	b	a	a	b	a
0	0	1	2	0	1	2

LEMMA. The overlap-identifiers can be determined in $\sim m$ steps.

Proof. Suppose $h(1), \dots, h(j)$ have been determined and we want $h(j+1)$ next.

We know

$$\begin{array}{c} \overbrace{\dots\dots\dots Y_j \quad Y_{j+1}} \\ \dots\dots\dots Y_{h(j)} \end{array}$$

and can put $h(j+1) = h(j)+1$ in case $Y_{j+1} = Y_{h(j)+1}$, but otherwise we must back up further

$$\begin{array}{c} \overbrace{\overbrace{\dots\dots\dots Y_j} \quad Y_{j+1}} \\ \dots\dots\dots Y_{h(j)} \\ \dots\dots\dots Y_{h(h(j))} \\ \vdots \end{array}$$

and further until such a match is found. The algorithm can be formulated as

```

h(1) := 0;
j := 1;
while j < m do

```

```

i := h(j);
while y[j+1] ≠ y[i+1] & i > 0 do
i := h(i)
od;
h(j+1) := if y[j+1] ≠ y[i+1] then 0 else i+1;
j := j+1
od;

```

Consider the total cost of this algorithm as we keep sliding a copy of y across along y itself. As long as $y[j+1]$ and $y[i+1]$ do not match we can charge the cost for a step to position $y[j-i+1]$ (and note here that i steadily decreases). If there is a match we charge the cost to $y[j+1]$ and continue with an increased value of j . No position can get charged more than twice. Thus the algorithm takes only a linear number of steps. \square

6.17. Once the overlap-identifiers are determined we can enact the *on-line* pattern-match across x . The algorithm (see 6.15) becomes

```

match := 0;
scan := 1;
while scan ≤ n do
if x[scan] = y[match+1] then goto L fi;
while x[scan] ≠ y[match+1] do
match := h(match)
until match = 0;
L: match := if x[scan] = y[match+1] then match+1 fi;
if match = m then pattern found & exit fi
scan := scan+1
od;

```

By changing the "exit" into

```

record the match at scan;
match := h(match);

```

one can use the algorithm to locate *all* (possibly overlapping) occurrences of y . Observe that each time around the loop we either match a next posi-

tion of x or move the beginning of y up, thus "pointing" at a position of x for the second time (but we will never point at it later after we passed it). Thus the algorithm requires time bounded by $\sim 2n$.

Adding the work we obtain

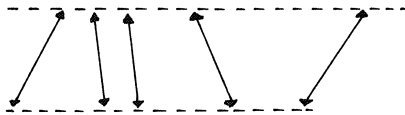
THEOREM. *One can determine whether y is a substring of x in $\sim n+m$ steps.*

6.18. Hirschberg [37] considered an interesting variant of the longest common substring problem, now known as *the maximal common subsequence problem*.

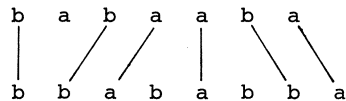
DEFINITION. z is a maximal common subsequence of x and y if and only if

- (i) $z \leq x$ & $z \leq y$;
- (ii) there is no w with $|w| > |z|$ such that $w \leq x$ & $w \leq y$.

The problem typically occurs when a maximal correspondence between two data-arrays must be found, and is encountered in comparing molecular structures and for instance in *problem-solving*. Learning systems classify problems by a series of features in order of decreasing importance, and to find out to what extent two problems are similar (or "analogical") one tries to establish the best possible sequence of common features



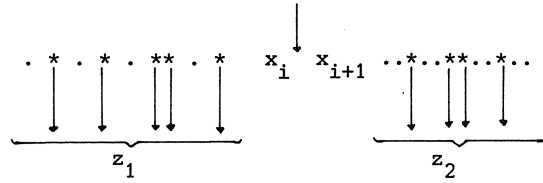
EXAMPLE. A longest common subsequence of babaaba and bbababba is bbaaba



The problem poses some interesting questions of time- and storage-utilization (the "simplicity" of the task is again deceptive). We shall prove Hirschberg's result (in a somewhat different presentation) that there is a data-organisation such that

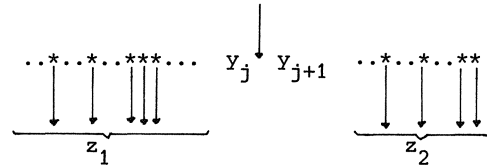
THEOREM. *The maximal common subsequence of x and y can be found in $\sim nm$ steps and $\sim n+m$ space.*

6.19. If z is to be a longest common sequence then for each position i in x



where $z_1 z_2 = z$

there must be a corresponding "split"-position j in y



such that there is no position k in y for which the maximal common subsequence of

$$x[1] \dots x[i] \quad \text{and} \quad y[1] \dots y[k]$$

concatenated with the maximal common subsequence of

$$x[i+1] \dots x[n] \quad \text{and} \quad y[k+1] \dots y[m]$$

is longer than z .

This criterion immediately suggests an algorithm to find z . Split x in two parts

$$x[1] \dots x[\lfloor \frac{n}{2} \rfloor] \quad \text{and} \quad x[\lfloor \frac{n}{2} \rfloor + 1] \dots x[n]$$

and find a position j such that the longest common subsequences of

$$x[1] \dots x[\lfloor \frac{n}{2} \rfloor] \quad \text{and} \quad y[1] \dots y[j]$$

and

$$x[\lfloor \frac{n}{2} \rfloor + 1] \dots x[n] \quad \text{and} \quad y[j+1] \dots y[m]$$

determined by a recursive application of the algorithm concatenate to a longest possible string.

6.20. The common trick in *dynamic programming* learns to compute the value of j quickly first (without actually computing what the common subsequences

are), and then do a specific recursive call to do the hard work of finding the subsequences.

The following concept is now needed.

DEFINITION. For $i \leq j$ and $k \leq \ell$ let $M_{ik}[j, \ell]$ be the length of the maximal common subsequence of

$$x[i] \dots x[j] \quad \text{and} \quad y[k] \dots y[\ell] .$$

LEMMA. One can compute $M_{ik}[j, \ell]$ in $\sim (j-i)(\ell-k)$ time and only $\sim (\ell-k)$ space.

Proof. One can determine $M_{ik}[i, k], M_{ik}[i, k+1], \dots, M_{ik}[i, \ell]$ by comparing symbol after symbol with $x[i]$ in $\sim \ell - k$ steps.

Assuming that $M_{ik}[s-1, k], M_{ik}[s-1, k+1], \dots, M_{ik}[s-1, \ell]$ have been determined, one can compute a next row with

```

Mik[s, k] := if y[k] occurs in x[i]...x[s] then 1 else 0;
for t := k+1 to ℓ do
Mik[s, t] := if x[s] = y[t] then Mik[s-1, t-1]+1
               else MAX(Mik[s-1, t], Mik[s, t-1]);
od;

```

One must compute $j - i + 1$ rows. \square

6.21. By building $M_{ik}[j, \ell]$ from "right-to-left" one can similarly compute rows

$$M_{ik}[j, \ell], M_{i, k+1}[j, \ell], \dots, M_{i, \ell}[j, \ell]$$

in the same time and space bounds.

6.22. We can now develop the recursive procedure (due to Hirschberg)

FIND $[i, j, k, \ell, \alpha]$ determining α as a longest common subsequence of

$$x[i] \dots x[j] \quad \text{and} \quad y[k] \dots y[\ell]$$

```

procedure FIND[i,j,k,l,α]
begin
  if i = j or k = l then determine α directly fi;
  half := i + (j-i)/2;
  compute row Mik[half,k], Mik[half,k+1], ..., Mik[half,l];
  compute row Mhalf+1 k[j,l], Mhalf+1 k+1[j,l], ..., Mhalf+1 l[j,l];
  split := an index k ≤ t ≤ l such that
      Mik[half,t] + Mhalf+1 t+1[j,l]
      is maximal (over all choices of t);
  FIND[i,half,k,split,α1];
  FIND[half+1,j,split+1,l,α2];
  return α := α1α2
end;

```

It is easy to make the procedure iterative and to see that only linear space is needed.

In the time-analysis we shall use that for integers $x, y \geq 1$ and $A \geq 2B + C$ (A, B , and C positive)

$$A^{xy} \geq B^{(x+y)} + C.$$

Assume as an induction-hypothesis that $\text{FIND}[i, j, k, l, \alpha]$ requires only

$$a_1(j-1)(l-k) + a_2[(j-i)+(l-k)] + a_3$$

time.

Choosing $a_2 > 1$ and $a_3 \geq 1$ appropriately it holds when $i = j$ and/or $k = l$.

For the induction-step we shall assume (by 6.20 and 6.21) that the computation of a row $M_{ik}[j, *]$ or $M_{i*}[j, l]$ takes

$$b_1(j-i)(l-k) + b_2[(j-i)+(l-k)] + b_3$$

time.

Then the time for $\text{FIND}[i, j, k, l, \alpha]$ is bounded by

$$\begin{aligned}
 & \text{const}_1 + 2b_1(j-i)(l-k) + 2b_2[(j-i)+(l-k)] + 2b_3 + \\
 & + \text{const}_2(l-k) + a_1 \frac{j-i}{2}(l-k) + a_2[(j-i)+(l-k)] + a_3 + \\
 & + \text{const}_3 =
 \end{aligned}$$

$$\begin{aligned}
&= \left(\frac{a_1}{2} + 2b_1\right)(j-i)(l-k) + (a_2 + 2b_2 + \text{const}_2)[(j-i) + (l-k)] + \\
&\quad + (a_3 + 2b_3 + \text{const}_1 + \text{const}_3) \leq \\
&\leq \left(\frac{a_1}{2} + 2b_1 + 2a_2 + 4b_2 + 2\text{const}_2 + a_3 + 2b_3 + \text{const}'\right)(j-i)(l-k).
\end{aligned}$$

If we choose a_1 such that

$$a_1 \geq \frac{a_1}{2} + 2b_1 + 2a_2 + 4b_2 + a_3 + 2b_3 + \text{const}''$$

(which we can) then the assumption is consistent. \square

6.23. One may now ask if there is perhaps a faster than quadratic algorithm to find maximal common subsequences. Aho, Hirschberg, and Ullman [2] proved under weak assumptions that in general there isn't.

Questions of *optimality* require that we somehow delimit the class of algorithms under consideration.

Consider straight-line programs which extract information from the input only by *cross-comparisons* and count how many queries

$$\text{"is } x[p] = y[q]\text{"}$$

are needed in worst case to find a longest common subsequence.

Note that if $x^{(1)}$ and $y^{(1)}$ and $x^{(2)}$ and $y^{(2)}$ yield the same answer query after query such algorithms generate the same computation in both instances and locate the longest common subsequence in identical manner. We shall exploit this fact to prove that if such algorithms query the input only a "few" times they cannot distinguish different inputs correctly.

6.24. We can immediately show

LEMMA. *Any algorithm for solving the longest common subsequence problem needs $\geq n + m - 1$ queries in worst case.*

Proof. Consider the sequence of queries (all yielding no) generated on

$$x = \underbrace{0 \ 0 \ \dots \ 0}_n \quad \text{and} \quad y = \underbrace{1 \ 1 \ \dots \ 1}_m$$

which have an "empty" solution.

Construct a *bipartite* graph G connecting the positions queried in each step.

If G has at least two connected components, then change the symbols of the positions of x and y in one component to 1 and 0 (respectively). The queries don't change and the algorithm still yields "empty" despite that there now is a non-trivial common subsequence.

It follows that G must be a connected graph with $n + m$ vertices. \square

6.25. For $\sigma = 2$ (say alphabet $\Sigma = \{0,1\}$) the given bound cannot be improved. The positions of x and y can be identified in only $n + m - 1$ comparisons as follows.

Think of $x[1]$ as "0" and ask "is $x[1] = y[1]$ " to determine if we can think of $y[1]$ as "0" also or should consider it as "1". Identify $x[2] \dots x[n]$ by asking "is $x[*] = y[1]$ ", and $y[2] \dots y[m]$ by asking "is $x[1] = y[*]$ ".

For $\sigma > 2$ it doesn't work, and we can indeed improve the lemma to a quadratic bound.

THEOREM. For $\sigma > 2$ any algorithm for solving the longest common subsequence problem requires $\geq nm$ queries in worst case.

Proof. Construct the same graph G as in 6.24, and suppose it has $< nm$ lines. Then there must be a pair p, q not connected and thus not queried in the algorithm.

Change $x[p]$ and $y[q]$ into "2" (which we can because a third symbol is available). It does not affect the answer of any query asked and the algorithm still comes up with "empty" despite that the correct result must be a non-trivial string. \square

REFERENCES (some references are not explicitly cited in the text).

- [1] ADEL'SON-VEL'SKII, G.M., and Y.M. LANDIS, An algorithm for the organization of information, Soviet Math. Dokl. 3 (1962) 1259-1262.
- [2] AHO, A.V., D.S. HIRSCHBERG, and J.D. ULLMAN, Bounds on the complexity of the longest common subsequence problem, Conf. Record 15th Annual IEEE Symp. Switching and Automata Theory, New Orleans (1974) 104-109, also: JACM 23 (1976) 1-12.

- [3] AHO, A.V., J.E. HOPCROFT, and J.D. ULLMAN, *The design and analysis of computer algorithms*, Addison-Wesley, Reading, Mass. (1974).
- [4] AHO, A.V., J.E. HOPCROFT, and J.D. ULLMAN, On finding lowest common ancestors in trees, Proc. 5th Annual ACM symp. Theory of Computing, Austin (1973) 253-265, also: SIAM J. Computing 5 (1976) 115-132.
- [5] AMBLE, O., and D.E. KNUTH, Ordered hash-tables, Stanford University, Stanford (1973) STAN-CS-73-367, also: Computer J. 17 (1974) 135-142.
- [6] ARDEN, B.W., B.A. GALLER, and R.M. GRAHAM, An algorithm for equivalence declarations, CACM 4 (1961) 310-314.
- [7] BAYER, R., Binary B-trees for virtual memory, Proc. ACM SIGFIDET Workshop, San Diego (1971) 219-235.
- [8] BAYER, R., Symmetric binary B-trees: data-structures and maintenance algorithms, Acta Inf. 1 (1972) 290-306.
- [9] BENTLEY, J.L. and W.A. BURKHARD, Heuristics for partial-match retrieval data-base design, Inf.Proc. Letters 4 (1976) 132-135.
- [10] BLOOM, B.H., Space/time trade-offs in hash-coding with allowable errors, CACM 13 (1970) 422-426.
- [11] BOBROW, D.G., A note on hash-linking, CACM 18 (1975) 413-415.
- [12] BOOTH, K.S., and G.S. LUEKER, Linear algorithms to recognize interval-graphs and test for the consecutive ones property, Proc. 7th Annual ACM Symp. Theory of Computing, Albuquerque (1975) 255-265.
- [13] BOOTH, K.S., and G.S. LUEKER, PQ-tree algorithms, Dept. of Electr. Engineering and Computer Sc., Univ. of California, Berkeley (1975) (preprint).
- [14] BURKHARD, W.A., Hashing and trie-algorithms for partial match retrieval, Dept. of Applied Physics and Information Sc., Unif. of California/San Diego, La Jolla (1975) TR-2.
- [15] BURKHARD, W.A., Partial-match queries and file-design, Proc. ACM Conf. on Very Large Data-bases (1975) 523-525.
- [16] BURKHARD, W.A., Associative retrieval trie hash-coding, Dept. of Applied Physics and Information Sc., Univ. of California/

- San Diego, La Jolla (1975) TR-6, also: Proc. 8th Annual ACM Symp. Theory of Computing, Hershey (1976) (to appear).
- [17] CODD, E.F., A Relational model of data for large shared data banks, CACM 13 (1970) 377-387.
- [18] CODD, E.F., A data-base sublanguage founded on the relational calculus, Proc. ACM SIGFIDET Workshop on data-description, access, and control (1971).
- [19] COFFMAN, E.G., and J. EVE, File-structures using hashing functions, CACM 7 (1970) 427-432,436.
- [20] COMER, D., and R. SETHI, NP-completeness of tree-structured index minimization, Dept. of Computer Sc., Pennsylvania State Univ., College Park (1975) (preprint).
- [21] COOK, S.A., Linear time simulation of deterministic two-way pushdown automata, Proc. IFIP Congress 71, TA-2, North-Holland Publ. Comp., Amsterdam (1971) 172-179.
- [22] DATE, C.J., *An introduction to data-base systems*, the Systems Progr. Series, Addison-Wesley, Reading, Mass. (1975).
- [23] DE LA BRIANDAIS, R., File searching using variable length keys, Proc. Western Joint Computer Conf. (1959) 295-298.
- [24] DUBOST, P., and J-M. TROUSSE, Software implementation of a new method of combinatorial hashing, Stanford University, Stanford (1975) STAN-CS-75-511.
- [25] ENGLES, R.W., A tutorial on data-base organisation, Annual Review in Automatic Progr. 7 (1972) 1-64.
- [26] EVE, J., On computing the transitive closure of a relation, Stanford University, Stanford (1975) STAN-CS-75-508.
- [27] FISCHER, M.J., Efficiency of equivalence algorithms, in: R.E. MILLER, and J.W. THATCHER (ed.), *Complexity of Computer Computations*, Plenum Press, New York (1972) 153-168.
- [28] FLOYD, R.W., and A.J. SMITH, A linear time two-tape merge, Stanford University, Stanford (1972) STAN-CS-72-285, also: Inform. Proc. Letters 2 (1974) 123-125.

- [29] FOSTER, D.V., Elementary file organisations: a survey, Duke University, Durham (1975) CS-1975-6.
- [30] FREDKIN, E., Trie memory, CACM 3 (1960) 490-499.
- [31] FREDMAN, M.L., Two applications of a probabilistic search technique: sorting $X + Y$ and building balanced search trees, Proc. 7th Annual ACM Symp. Theory of Computing, Albuquerque (1975) 240-244.
- [32] GALLER, B.A., and M.J. FISCHER, An improved equivalence algorithm, CACM 7 (1964) 301-303, 506.
- [33] GEHANI, N., private communication, SUNY/Buffalo, Amherst (1976).
- [34] GRISWOLD, R., *String and list processing in SNOBOL 4: techniques and applications*, PH Series in Automatic Computation, Prentice-Hall, Englewood Cliffs, N.J. (1974).
- [35] GUSTAFSON, R.A., A randomized combinatorial file-structure for storage and retrieval systems, Ph.D. thesis, Univ. of South Carolina (1969).
- [36] HILDEBRAND, P., and H. ISBITZ, Radix-exchange: an internal sorting method for digital computers, JACM 6 (1959) 156-163.
- [37] HIRSCHBERG, D.S., A linear space algorithm for computing maximal common subsequences, CACM 18 (1975) 341-343.
- [38] HOPCROFT, J.E., and R.M. KARP, An algorithm for testing the equivalence of finite automata, Cornell University, Ithaca (1971) TR-71-114.
- [39] HOPCROFT, J.E., and J.D. ULLMAN, Set merging algorithms, SIAM J. Computing 2 (1973) 294-303.
- [40] HORVATH, E.C., Efficient stable sorting with minimal extra space, Proc. 6th Annual ACM Symp. Theory of Computing, Seattle (1974) 194-215.
- [41] KARP, R.M., R.E. MILLER, and A.L. ROSENBERG, Rapid identification of repeated patterns in strings, trees, and arrays, Proc. 4th Annual ACM Symp. Theory of Computing, Denver (1972) 125-136.
- [42] KNUTH, D.E., and V.R. PRATT, Automata theory can be useful, Stanford University, Stanford (1971) STAN-CS-71- .

- [43] KNUTH, D.E., *The art of computer programming I: Fundamental algorithms*, Addison-Wesley, Reading, Mass. (1969).
- [44] KNUTH, D.E., *The art of computer programming III: Sorting and searching*, Addison-Wesley, Reading, Mass. (1973).
- [45] MAURER, H.A., and D. WOOD, Zur Manipulation von Zahlenmengen, Inst. f. Angew. Informatik u. formale Beschreibungsverfahren, Univ. Karlsruhe (1975) Bericht 34.
- [46] MAURER, H.A., TH. OTTMAN, and H.W. SIX, Implementing dictionaries using binary trees of very small height, Inst. f. Angew. Informatik u. formale Beschreibungsverfahren, Univ. Karlsruhe (1975) Bericht 37.
- [47] MORRIS, J.H., and V.R. PRATT, A linear time pattern-matching algorithm, Computing Centre, Univ. of California, Berkeley (1970) TR-40.
- [48] MORRISON, D.R., PATRICIA-practical algorithm to retrieve information coded in alphanumeric, JACM 15 (1968) 514-534.
- [49] OTTMAN, TH., and H.W. SIX, Eine neue Klasse von Ausgeglichenen Binärbäumen, Inst. f. Angew. Informatik u. formale Beschreibungsverfahren, Univ. Karlsruhe (1975) Bericht 35.
- [50] PORTER, T., and I. SIMON, Random insertion into a priority queue structure, Stanford University, Stanford (1974) STAN-CS-74-460.
- [51] RIVEST, R.L., On self-organising sequential search heuristics, Conf. Record 15th Annual IEEE Symp. Switching and Automata Theory, New Orleans (1974) 122-126, also: CACM 19 (1976) 63-67.
- [52] RIVEST, R.L., On hash-coding algorithms for partial-match retrieval, Conf. Record 15th Annual IEEE Symp. Switching and Automata Theory, New Orleans (1974) 95-103.
- [53] RIVEST, R.L., Partial-match retrieval algorithms, SIAM J. Computing 5 (1976) 19-50.
- [54] ROBERTS, D.C., File organisation techniques, in: *Adv. in Computers* 12 (1972) 115-174.
- [55] SPITZEN, J., and B. WEGBREIT, The verification and synthesis of data-structures, Acta Inf. 4 (1975) 127-144.

- [56] TARJAN, R.E., Efficiency of a good but not linear set union algorithm, JACM 22 (1975) 215-225.
- [57] TARJAN, R.E., Applications of path-compression on balanced trees, Stanford University, Stanford (1975) STAN-CS-75-512.
- [58] TRABB PARDO, L., Stable sorting and merging with optimal space and time bounds, Stanford University, Stanford (1974) STAN-CS-74-470.
- [59] VAN EMDE BOAS, P., Preserving order in a forest in less than logarithmic time, Proc. 16th Annual IEEE Symp. Foundations of Computer Sc., Berkeley (1975) 75-84.
- [60] VAN LEEUWEN, J., On finding lowest common ancestors in less than logarithmic average time, Symp. New Directions and Recent Results in Algorithms and Complexity, Carnegie-Mellon Univ., Pittsburgh (1976) 107.
- [61] VAN LEEUWEN, J., *The complexity of data-organisation*, 2nd Adv. Course on Foundations of Computer Science, Amsterdam (1976).
- [62] VAN LEEUWEN, J., On the construction of Huffman-trees, in: S. MICHAELSON & R. MILNER (ed), 3rd Int. Colloq. on Automata/languages/Programming, Edinburgh (1976), 382-410.
- [63] WARREN, H.S., A modification of Warshall's algorithm for the transitive closure of binary relations, CACM 18 (1975) 218-220.
- [64] WARSHALL, S., A theorem on Boolean matrices, JACM 9 (1962) 11-12.
- [65] WEINER, P., Linear pattern-matching algorithms, Conf. Record 14th Annual IEEE Symp. Switching and Automata Theory, Iowa City (1973) 1-11.
- [66] WIRTH, N., *Algorithms + data-structures = programs*, PH Series in Automatic Progr., Prentice-Hall, Englewood Cliffs, N.J. (1976).

OTHER TITLES IN THE SERIES MATHEMATICAL CENTRE TRACTS

A leaflet containing an order-form and abstracts of all publications mentioned below is available at the Mathematisch Centrum, Tweede Boerhaavestraat 49, Amsterdam-1005, The Netherlands. Orders should be sent to the same address.

- MCT 1 T. VAN DER WALT, *Fixed and almost fixed points*, 1963. ISBN 90 6196 002 9.
- MCT 2 A.R. BLOEMENA, *Sampling from a graph*, 1964. ISBN 90 6196 003 7.
- MCT 3 G. DE LEVE, *Generalized Markovian decision processes, part I: Model and method*, 1964. ISBN 90 6196 004 5.
- MCT 4 G. DE LEVE, *Generalized Markovian decision processes, part II: Probabilistic background*, 1964. ISBN 90 6196 005 3.
- MCT 5 G. DE LEVE, H.C. TIJMS & P.J. WEEDA, *Generalized Markovian decision processes, Applications*, 1970. ISBN 90 6196 051 7.
- MCT 6 M.A. MAURICE, *Compact ordered spaces*, 1964. ISBN 90 6196 006 1.
- MCT 7 W.R. VAN ZWET, *Convex transformations of random variables*, 1964. ISBN 90 6196 007 X.
- MCT 8 J.A. ZONNEVELD, *Automatic numerical integration*, 1964. ISBN 90 6196 008 8.
- MCT 9 P.C. BAAYEN, *Universal morphisms*, 1964. ISBN 90 6196 009 6.
- MCT 10 E.M. DE JAGER, *Applications of distributions in mathematical physics*, 1964. ISBN 90 6196 010 X.
- MCT 11 A.B. PAALMAN-DE MIRANDA, *Topological semigroups*, 1964. ISBN 90 6196 011 8.
- MCT 12 J.A.TH.M. VAN BERCKEL, H. BRANDT CORSTIUS, R.J. MOKKEN & A. VAN WIJNGAARDEN, *Formal properties of newspaper Dutch*, 1965. ISBN 90 6196 013 4.
- MCT 13 H.A. LAUWERIER, *Asymptotic expansions*, 1966, out of print; replaced by MCT 54 and 67.
- MCT 14 H.A. LAUWERIER, *Calculus of variations in mathematical physics*, 1966. ISBN 90 6196 020 7.
- MCT 15 R. DOORNBOS, *Slippage tests*, 1966. ISBN 90 6196 021 5.
- MCT 16 J.W. DE BAKKER, *Formal definition of programming languages with an application to the definition of ALGOL 60*, 1967. ISBN 90 6196 022 3.
- MCT 17 R.P. VAN DE RIET, *Formula manipulation in ALGOL 60, part 1*, 1968. ISBN 90 6196 025 8.
- MCT 18 R.P. VAN DE RIET, *Formula manipulation in ALGOL 60, part 2*, 1968. ISBN 90 6196 038 X.
- MCT 19 J. VAN DER SLOT, *Some properties related to compactness*, 1968. ISBN 90 6196 026 6.
- MCT 20 P.J. VAN DER HOUWEN, *Finite difference methods for solving partial differential equations*, 1968. ISBN 90 6196 027 4.

- MCT 21 E. WAITTEL, *The compactness operator in set theory and topology*, 1968. ISBN 90 6196 028 2.
- MCT 22 T.J. DEKKER, *ALGOL 60 procedures in numerical algebra, part 1*, 1968. ISBN 90 6196 029 0.
- MCT 23 T.J. DEKKER & W. HOFFMANN, *ALGOL 60 procedures in numerical algebra, part 2*, 1968. ISBN 90 6196 030 4.
- MCT 24 J.W. DE BAKKER, *Recursive procedures*, 1971. ISBN 90 6196 060 6.
- MCT 25 E.R. PAERL, *Representations of the Lorentz group and projective geometry*, 1969. ISBN 90 6196 039 8.
- MCT 26 EUROPEAN MEETING 1968, *Selected statistical papers, part I*, 1968. ISBN 90 6196 031 2.
- MCT 27 EUROPEAN MEETING 1968, *Selected statistical papers, part II*, 1969. ISBN 90 6196 040 1.
- MCT 28 J. OOSTERHOFF, *Combination of one-sided statistical tests*, 1969. ISBN 90 6196 041 X.
- MCT 29 J. VERHOEFF, *Error detecting decimal codes*, 1969. ISBN 90 6196 042 8.
- MCT 30 H. BRANDT CORSTIUS, *Excercises in computational linguistics*, 1970. ISBN 90 6196 052 5.
- MCT 31 W. MOLENAAR, *Approximations to the Poisson, binomial and hypergeometric distribution functions*, 1970. ISBN 90 6196 053 3.
- MCT 32 L. DE HAAN, *On regular variation and its application to the weak convergence of sample extremes*, 1970. ISBN 90 6196 054 1.
- MCT 33 F.W. STEUTEL, *Preservation of infinite divisibility under mixing and related topics*, 1970. ISBN 90 6196 061 4.
- MCT 34 I. JUHÁSZ, A. VERBEEK & N.S. KROONENBERG, *Cardinal functions in topology*, 1971. ISBN 90 6196 062 2.
- MCT 35 M.H. VAN EMDEN, *An analysis of complexity*, 1971. ISBN 90 6196 063 0.
- MCT 36 J. GRASMAN, *On the birth of boundary layers*, 1971. ISBN 90 6196 064 9.
- MCT 37 J.W. DE BAKKER, G.A. BLAAUW, A.J.W. DUIJVESTIJN, E.W. DIJKSTRA, P.J. VAN DER HOUWEN, G.A.M. KAMSTEEG-KEMPER, F.E.J. KRUSEMAN ARETZ, W.L. VAN DER POEL, J.P. SCHAAP-KRUSEMAN, M.V. WILKES & G. ZOUTENDIJK, *MC-25 Informatica Symposium*, 1971. ISBN 90 6196 065 7.
- MCT 38 W.A. VERLOREN VAN THEMAAT, *Automatic analysis of Dutch compound words*, 1971. ISBN 90 6196 073 8.
- MCT 39 H. BAVINCK, *Jacobi series and approximation*, 1972. ISBN 90 6196 074 6.
- MCT 40 H.C. TIJMS, *Analysis of (s,S) inventory models*, 1972. ISBN 90 6196 075 4.
- MCT 41 A. VERBEEK, *Superextensions of topological spaces*, 1972. ISBN 90 6196 076 2.
- MCT 42 W. VERVAAT, *Success epochs in Bernoulli trials (with applications in number theory)*, 1972. ISBN 90 6196 077 0.
- MCT 43 F.H. RUYMGAART, *Asymptotic theory of rank tests for independence*, 1973. ISBN 90 6196 081 9.
- MCT 44 H. BART, *Meromorphic operator valued functions*, 1973. ISBN 90 6196 082 7.

- MCT 45 A.A. BALKEMA, *Monotone transformations and limit laws*, 1973.
ISBN 90 6196 083 5.
- MCT 46 R.P. VAN DE RIET, *ABC ALGOL, A portable language for formula manipulation systems, part 1: The language*, 1973. ISBN 90 6196 084 3.
- MCT 47 R.P. VAN DE RIET, *ABC ALGOL, A portable language for formula manipulation systems, part 2: The compiler*, 1973. ISBN 90 6196 085 1.
- MCT 48 F.E.J. KRUSEMAN ARETZ, P.J.W. TEN HAGEN & H.L. OUDSHOORN, *An ALGOL 60 compiler in ALGOL 60, Text of the MC-compiler for the EL-X8*, 1973. ISBN 90 6196 086 X.
- MCT 49 H. KOK, *Connected orderable spaces*, 1974. ISBN 90 6196 088 6.
- MCT 50 A. VAN WIJNGAARDEN, B.J. MAILLOUX, J.E.L. PECK, C.H.A. KOSTER, M. SINTZOFF, C.H. LINDSEY, L.G.L.T. MEERTENS & R.G. FISHER (Eds). *Revised report on the algorithmic language ALGOL 68*, 1976. ISBN 90 6196 089 4.
- MCT 51 A. HORDIJK, *Dynamic programming and Markov potential theory*, 1974. ISBN 90 6196 095 9.
- MCT 52 P.C. BAAZEN (ed.), *Topological structures*, 1974. ISBN 90 6196 096 7.
- MCT 53 M.J. FABER, *Metrisability in generalized ordered spaces*, 1974. ISBN 90 6196 097 5.
- MCT 54 H.A. LAUWERIER, *Asymptotic analysis, part 1*, 1974. ISBN 90 6196 098 3.
- MCT 55 M. HALL JR. & J.H. VAN LINT (Eds), *Combinatorics, part 1: Theory of designs, finite geometry and coding theory*, 1974. ISBN 90 6196 099 1.
- MCT 56 M. HALL JR. & J.H. VAN LINT (Eds), *Combinatorics, part 2: graph theory, foundations, partitions and combinatorial geometry*, 1974. ISBN 90 6196 100 9.
- MCT 57 M. HALL JR. & J.H. VAN LINT (Eds), *Combinatorics, part 3: Combinatorial group theory*, 1974. ISBN 90 6196 101 7.
- MCT 58 W. ALBERS, *Asymptotic expansions and the deficiency concept in statistics*, 1975. ISBN 90 6196 102 5.
- MCT 59 J.L. MIJNHEER, *Sample path properties of stable processes*, 1975. ISBN 90 6196 107 6.
- MCT 60 F. GÖBEL, *Queueing models involving buffers*, 1975. ISBN 90 6196 108 4.
- * MCT 61 P. VAN EMDE BOAS, *Abstract resource-bound classes, part 1*. ISBN 90 6196 109 2.
- * MCT 62 P. VAN EMDE BOAS, *Abstract resource-bound classes, part 2*. ISBN 90 6196 110 6.
- MCT 63 J.W. DE BAKKER (ed.), *Foundations of computer science*, 1975. ISBN 90 6196 111 4.
- MCT 64 W.J. DE SCHIPPER, *Symmetric closed categories*, 1975. ISBN 90 6196 112 2.
- MCT 65 J. DE VRIES, *Topological transformation groups 1 A categorical approach*, 1975. ISBN 90 6196 113 0.
- MCT 66 H.G.J. PIJLS, *Locally convex algebras in spectral theory and eigenfunction expansions*, 1976. ISBN 90 6196 114 9.

- * MCT 67 H.A. LAUWERIER, *Asymptotic analysis, part 2.*
ISBN 90 6196 119 x.
- MCT 68 P.P.N. DE GROEN, *Singularly perturbed differential operators of second order*, 1976. ISBN 90 6196 120 3.
- MCT 69 J.K. LENSTRA, *Sequencing by enumerative methods*, 1977.
ISBN 90 6196 125 4.
- MCT 70 W.P. DE ROEVER JR., *Recursive program schemes: semantics and proof theory*, 1976. ISBN 90 6196 127 0.
- MCT 71 J.A.E.E. VAN NUNEN, *Contracting Markov decision processes*, 1976.
ISBN 90 6196 129 7.
- MCT 72 J.K.M. JANSEN, *Simple periodic and nonperiodic Lamé functions and their applications in the theory of conical waveguides*, 1977.
ISBN 90 6196 130 0.
- MCT 73 D.M.R. LEIVANT, *Absoluteness of intuitionistic logic*, 1979.
ISBN 90 6196 122 x.
- MCT 74 H.J.J. TE RIELE, *A theoretical and computational study of generalized aliquot sequences*, 1976. ISBN 90 6196 131 9.
- MCT 75 A.E. BROUWER, *Treelike spaces and related connected topological spaces*, 1977. ISBN 90 6196 132 7.
- MCT 76 M. REM, *Associons and the closure statement*, 1976. ISBN 90 6196 135 1.
- MCT 77 W.C.M. KALLENBERG, *Asymptotic optimality of likelihood ratio tests in exponential families*, 1977 ISBN 90 6196 134 3.
- MCT 78 E. DE JONGE, A.C.M. VAN ROOIJ, *Introduction to Riesz spaces*, 1977.
ISBN 90 6196 133 5.
- MCT 79 M.C.A. VAN ZUIJLEN, *Empirical distributions and rankstatistics*, 1977.
ISBN 90 6196 145 9.
- MCT 80 P.W. HEMKER, *A numerical study of stiff two-point boundary problems*,
1977. ISBN 90 6196 146 7.
- MCT 81 K.R. APT & J.W. DE BAKKER (eds), *Foundations of computer science II*,
part I, 1976. ISBN 90 6196 140 8.
- MCT 82 K.R. APT & J.W. DE BAKKER (eds), *Foundations of computer science II*,
part II, 1976. ISBN 90 6196 141 6.
- * MCT 83 L.S. VAN BENTEM JUTTING, *Checking Landau's "Grundlagen" in the
AUTOMATH system*, ISBN 90 6196 147 5.
- MCT 84 H.L.L. BUSARD, *The translation of the elements of Euclid from the
Arabic into Latin by Hermann of Carinthia (?) books vii-xii*, 1977.
ISBN 90 6196 148 3.
- MCT 85 J. VAN MILL, *Supercompactness and Wallman spaces*, 1977.
ISBN 90 6196 151 3.
- MCT 86 S.G. VAN DER MEULEN & M. VELDHORST, *Torrix I*, 1978.
ISBN 90 6196 152 1.
- * MCT 87 S.G. VAN DER MEULEN & M. VELDHORST, *Torrix II*,
ISBN 90 6196 153 x.
- MCT 88 A. SCHRIJVER, *Matroids and linking systems*, 1977.
ISBN 90 6196 154 8.

- MCT 89 J.W. DE ROEVER, *Complex Fourier transformation and analytic functionals with unbounded carriers*, 1978.
ISBN 90 6196 155 6.
- * MCT 90 L.P.J. GROENEWEGEN, *Characterization of optimal strategies in dynamic games*, . ISBN 90 6196 156 4.
- * MCT 91 J.M. GEYSEL, *Transcendence in fields of positive characteristic*, . ISBN 90 6196 157 2.
- * MCT 92 P.J. WEEDA, *Finite generalized Markov programming*,
ISBN 90 6196 158 0.
- MCT 93 H.C. TIJMS (ed.) & J. WESSELS (ed.), *Markov decision theory*, 1977.
ISBN 90 6196 160 2.
- MCT 94 A. BIJLSMA, *Simultaneous approximations in transcendental number theory*, 1978 . ISBN 90 6196 162 9.
- MCT 95 K.M. VAN HEE, *Bayesian control of Markov chains*, 1978.
ISBN 90 6196 163 7.
- * MCT 96 P.M.B. VITÁNYI, *Lindermayer systems: structure, languages, and growth functions*, . ISBN 90 6196 164 5.
- * MCT 97 A. FEDERGRUEN, *Markovian control problems; functional equations and algorithms*, . ISBN 90 6196 165 3.
- MCT 98 R. GEEL, *Singular perturbations of hyperbolic type*, 1978.
ISBN 90 6196 166 1
- MCT 99 J.K. LENSTRA, A.H.G. RINNOOY KAN & P. VAN EMDE BOAS, *Interfaces between computer science and operations research*, 1978.
ISBN 90 6196 170 X.
- MCT 100 P.C. BAAYEN, D. VAN DULST & J. OOSTERHOFF (Eds), *Proceedings bicentennial congress of the Wiskundig Genootschap, part 1*, 1979.
ISBN 90 6196 168 8.
- MCT 101 P.C. BAAYEN, D. VAN DULST & J. OOSTERHOFF (Eds), *Proceedings bicentennial congress of the Wiskundig Genootschap, part 2*, 1979.
ISBN 90 9196 169 6.
- MCT 102 D. VAN DULST, *Reflexive and superreflexive Banach spaces*, 1978.
ISBN 90 6196 171 8.
- MCT 103 K. VAN HARN, *Classifying infinitely divisible distributions by functional equations*, 1978 . ISBN 90 6196 172 6.
- * MCT 104 J.M. VAN WOUWE, *Go-spaces and generalizations of metrizability*,
. ISBN 90 6196 173 4.
- * MCT 105 R. HELMERS, *Edgeworth expansions for linear combinations of order statistics*, . ISBN 90 6196 174 2.

AN ASTERISK BEFORE THE NUMBER MEANS "TO APPEAR"

