

SOFTWARE IMPLEMENTATION OF A NEW METHOD
OF COMBINATORIAL HASHING

by

P. Dubost
J. -M. Trousse

STAN-CS-75-511
SEPTEMBER 1975

COMPUTER SCIENCE DEPARTMENT
School of Humanities and Sciences
STANFORD UNIVERSITY



Software Implementation of a New Method of Combinatorial Hashing

Pierre Dubost and Jean-Michel Trousse

Stanford University

Abstract

This is a study of the software implementation of a new method of searching with retrieval on secondary keys.

A new **family** of partial match file designs is presented, the 'worst case' is determined, a detailed algorithm and program are given and the average execution time is studied.

This research was supported in part by National Science Foundation grant DCR72-03752 A02 and by the Office of Naval Research contract NR 044-402. Reproduction in whole or in part is permitted for any purpose of the United States Government.

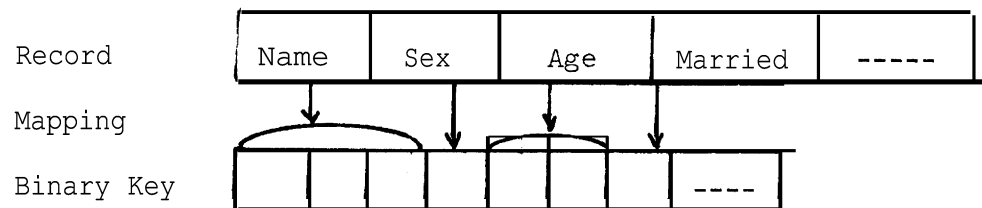
1. Retrieval on Secondary Key [2]

Common searching techniques use 'primary keys' which uniquely define a record. But, it is sometimes necessary to make a search on other fields of a record, called 'secondary keys'. We might want to retrieve some records from a file, given the values of some of these secondary keys. These values may define zero, one, or several records. A set of values is called a 'query'.

For example, if we considered an hypothetical CIA file containing information about all the American people, we might be interested in knowing which men are married, own two cars, and have been to France last year. We do not specify their age, residence, etc.

We now may assume that each secondary key is mapped by a common hashing technique into a short bit string. We then need a fast method of retrieval on a reasonably short binary key with some unspecified bits (noted *). This technique must map any specified value of this binary key into a shorter address. This address will correspond to a group of records called a 'bucket'. This technique must differ from a common hashing technique in that it allows some of the bits of the key to remain unspecified. In this case, it is clear that a query may lead to different buckets and the better the method is, the fewer buckets there will be for a given query.

Example: American People File



The previous query might be represented, for example, by: ****0**111---** .

2. A New Family of Partial Match File Designs and its Binary Tree Representation.

W. A. Burkhard has recently presented a new family of partial match file designs [1]. The interesting aspect of these designs is that they can be represented by a binary tree. The binary tree leads directly to a simple software implementation. To obtain an even simpler implementation, we have slightly modified the Burkhard partial match file (PMF) and introduced a new family of PMF which has, as we show in the next section, the same worst case.

The mapping of binary keys into the bucket addresses is described by a table. Each bucket corresponds to a given value of some of the bits, the others remaining unspecified. Each entry in the tables gives a description of the keys which might be in the corresponding bucket. For example, the bucket corresponding to the entry $*10*1$ might contain the following keys:

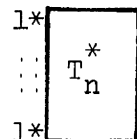
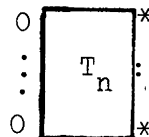
01001 , 01011 , 11001 , 11011 .

Any specified key must be mapped into one, and only one, bucket. This technique works for any query with an odd number of bits. The tables T_n are constructed by induction.

T_0 corresponds to a one bit key and two buckets:

| | |
|---|---------|
| | 0 |
| 0 | 0 |
| 1 | c^1_1 |

T_{n+1} is deduced from T_n by the following method:



T_n^* being T_n circularly shifted by one position (instead of the symmetric image of T_n with columns in the reversing order, as in the Burkhard technique). It is clear that T_n will have $2n+1$ columns and 2^{n+1} lines. It then maps $(2n+1)$ -bit keys into 2^{n+1} buckets.

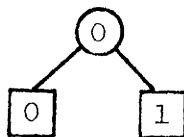
The representation of each table T_n by a tree S_n is then very simple. A node of the tree S_n contains the position of a bit in the query to be checked and a leaf contains the address of a bucket.

Let us consider the trees for $n = 0$ and 1 .

$n = 0$

The table T_0 is $\begin{array}{cc} & 0 \\ 0 & 0 \\ 1 & 1 \end{array}$ which means, if the bit is 0, the

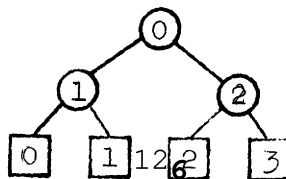
corresponding bucket address is 0 and if the bit is 1, the corresponding bucket address is 1. We can represent this procedure by a very simple tree S_0



with the convention that -- if the bit checked at a node is 0, we go to the left subtree (here the leaf 0) and to the right subtree for 1 (here the leaf 1). If the bit is a *, we go down in both subtrees.

$n = 1$

T_1 is $\begin{array}{c|ccc} & 0 & 1 & 2 \\ \hline 0 & 0 & 0 & * \\ 1 & 0 & 1 & * \\ 2 & 1 & * & 0 \\ 3 & 1 & * & 1 \end{array}$ the corresponding tree S_1 is



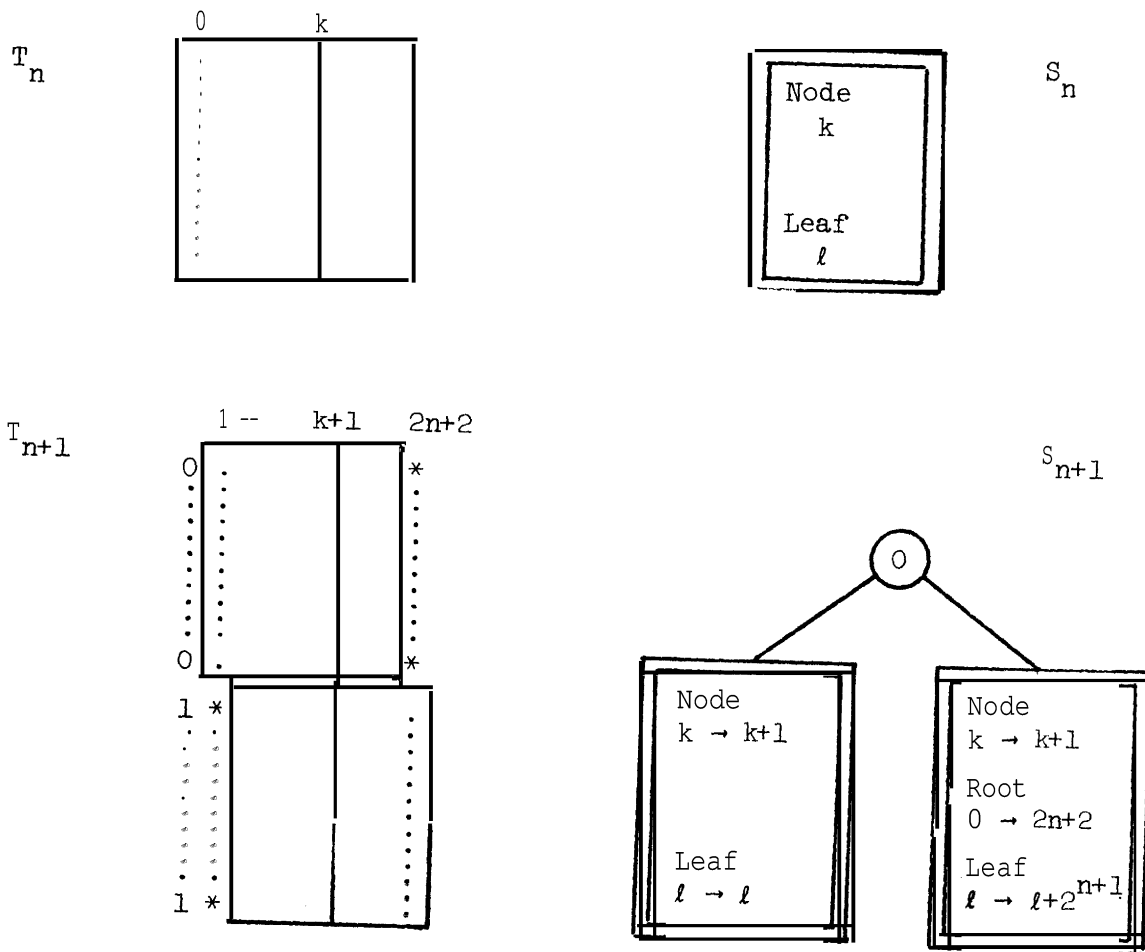
We have seen the passage from T_n to T_{n-1} . Similarly, we can define a transformation on the corresponding trees between S_n and S_{n+1} . Let us suppose that we have our tree S_n corresponding to the table 'I' _n:

The first bit separates T_{n+1} into two halves. If the first bit is 0, the corresponding bucket address is going to be contained in the first-half; or, if it is 1, in the second-half. Correspondingly, the root of S_{n+1} (node 0) separates S_{n+1} into a left and right subtree.

In the first-half of T_{n+1} , T_n is directly inserted, which means if a query leads to the bucket l in T_n , the query constructed with the preceding query (shifted by 1 position) with a 0 leading bit, is going to lead to the same bucket l in T_{n+1} . So, the left subtree of S_{n+1} is S_n but with the content of each node increased by one.

In the second-half of T_{n+1} , T_n is inserted after a circular shift of one position. In the same way, the right subtree of S_{n+1} is S_n but with the content of each node increased by one, except the root 0 changed into 2n+2 and the content of each leaf increased by 2^{n+1} .

Let us represent on the same figure, the induction on T_n and S_n :



It is easy to check the transformation between S_0 and S_1 .

Let us now determine the various properties of the tree S_n (see T_h and the corresponding tree S_h on the next page)

- A -- S_n is perfectly balanced.
- B -- At level i ($0 \leq i < n$), the content of the nodes is either i (denoted by \min_i) or $2n+1-i$ (denoted by \max_i).
 \min_i is always a left son and \max_i a right son.
- C -- The content of a leaf can be computed directly by the path used from the root adding at each level i , 0 if we go to the left and 2^{n-i} , if we go to the right.

```

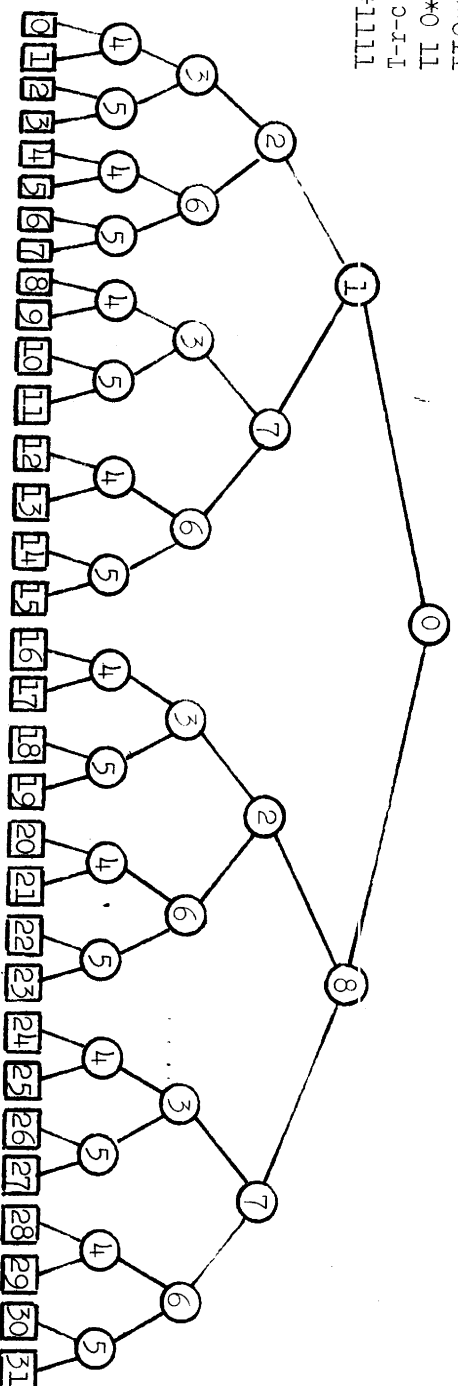
0 012345678
1 00000*****
2 00001*****
3 0001*0*****
4 0001*1****
5 001*0*0***
6 001*1*0**
7 001**01**
8 01*00**0*
9 01*01**0*
10 01*1*0*0*
11 01*1*1*0*
12 01**0*01*
13 01**1*01*
14 01***011*
15 01***111*
16 1*000*****
17 1*001****
18 1*01*0***
19 1*01*1**
20 1*1*0*0*0
21 1*1*1*0*0
22 1*1**01*0
23 1*1**11*0
24 1**00**01
25 1**0 1**01
26 1**1*0*01
27 1**1*1*01
28 1***0*011
29 1***1*0 11
30 1***0 c-r-T
31 1***1111

```

Γ_4

7

S_4



0
1
2
3
4

Properties A , B , C hold for S_0, S_1 . Suppose that A , B , C hold for S_n :

Obviously, A holds for S_{n+1} .

The level i of S_n corresponds to the level $i+1$ of S_{n+1} . For $i > 0$, the node i of S_n gives $i+1$ and the node $2n+1-i$ gives $2n+2-i = 2(n+1)+1-(i+1)$. For $i = 0$, the node 0 of S_n gives 1 in the left subtree and $2n+2 = 2(n+1)+1-1$ in the right subtree. So, the property B holds for S_{n+1} .

The increase at level $i+1$ in S_{n+1} is either 0 or $2^{n+1-(i+1)}$ exactly the same as in level i in S_n . If we follow the same path in the left subtree of S_{n+1} as in S_n , we find l as in S_n . If we follow the same path in the right subtree of S_{n+1} , we find $l+2^{n+1}$ since we had to go right in S_{n+1} at level 0 and add 2^{n+1} . . . , C holds for S_{n+1} .

The representation of this family of partial match file designs, by so simple a binary tree, leads to an easier proof of the worst case and a very simple search algorithm.

3. Worst Case.

The worst case is expressed in terms of the number of buckets $w_n(k)$ found when k bits are unspecified. W. A. Burkhard has shown that the worst case for his PMF family can be expressed in terms of the Fibonacci number in the following way:

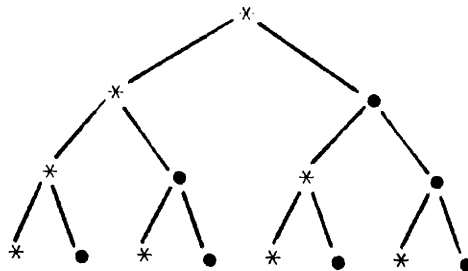
$$\begin{aligned} w_n(k) &= 2^k & \text{for } 0 \leq k \leq n+1 - \left\lceil \frac{n}{2} \right\rceil \\ w_n(k) &= 2^{n+1-k} F_{2k-n+1} & \text{for } n+1 - \left\lceil \frac{n}{2} \right\rceil \leq k \leq n+1 \\ w_n(k) &= 2^{k-(n+1)} F_{2n+4-k} & \text{for } n+1 < k \leq 2n+1 \end{aligned}$$

Using the tree, we can prove those same results for our new PMF family. Let us set the notation: If the position $i \leq n$ (or $i > n$) of the query contains a * , all the nodes which are left sons (or right sons)

result in a * in the level i (or $2n+1-i$) of S_n . In such a case, we are going to say that the level i has one star. Consequently, if all the nodes at a level i are *, we are going to say that the level i has two stars.

Finally, we represent S_n in the following way: a node is denoted by a * if it corresponds to an unspecified bit and by a • if it corresponds to a specified bit.

As an example, we can have the following tree:

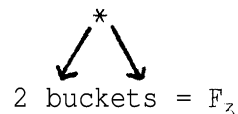


Since in our study of the worst case we considered first the left son to be unspecified, a dot corresponds, in fact, to a bit 0.

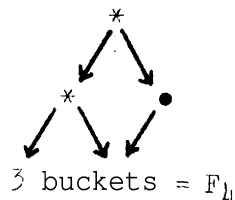
We can now study three different facts:

Fact 1. Let us see how the Fibonacci number arises in the worst case in studying the special case where all the levels contained one star. Such a tree with $n+1$ levels is denoted S_n^* .

For $n = 0$, we have S_0^*

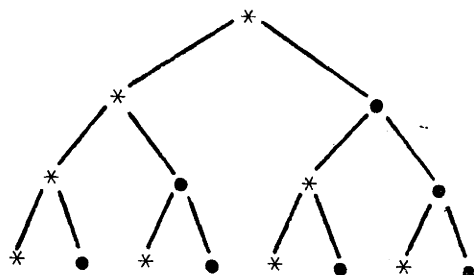


For $n = 1$, we have S_1^*

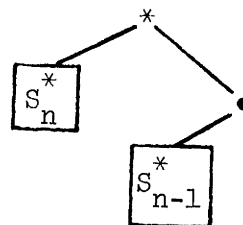


Let us assume that the formula for the number of buckets found, F_{n+2} holds for n and $n-1$, and let us study the case for $n+1$.

S_{n+1}^*



which can be expressed by:



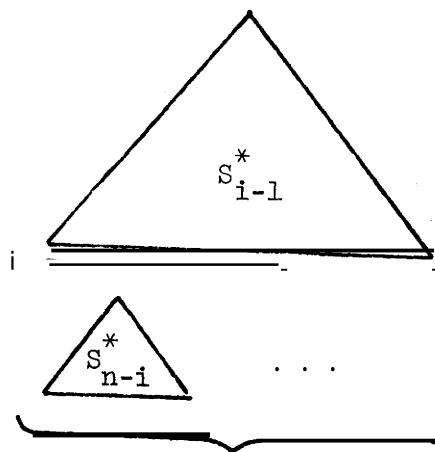
So, # of buckets for S_{n+1}^* = # of buckets for S_n^* + # of buckets for S_{n-1}^*

$$= F_{n+3} + F_{(n-1)+3} = F_{n+3} + F_{n+2}$$

$$= F_{n+4} = F_{(n+1)+3}$$

and the number of buckets for $S_n^* = F_{n+3}$.

Fact 2. Suppose a level i in S_n^* contains two stars, let us see at which level i the number of buckets is maximum. Such a tree can be represented by the following figure:



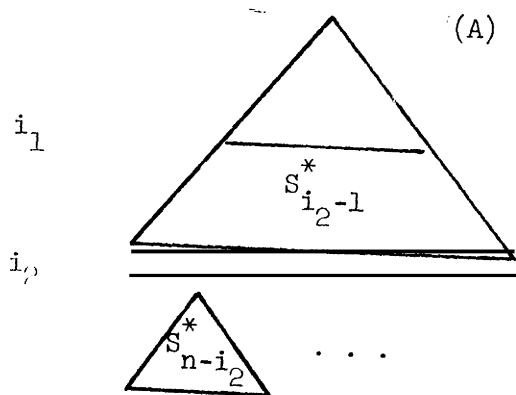
F_{i+2} trees S_{n-i}^* are explored and we have $F_{i+2} F_{n-i+3}$ buckets.

Now we have: $F_{i+2} F_{n-i+3} = F_{n+3} \cdot F_i F_{n+1-i}$.

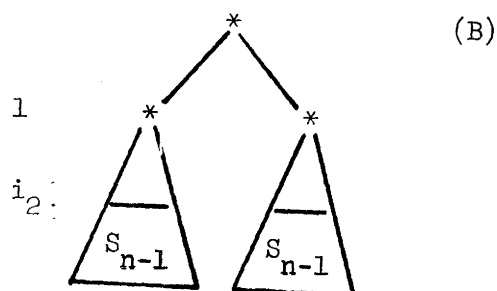
Then it can be easily shown that this product is maximum for $i = 1$ (or $i = n+1$) since $F_i F_{n+1-i}$ varies like $(-1)^{i+1} (F_{n-2i} + 2F_{n-2i+1})$. So, in a tree S_n^* a level i with two *s gives the worst case for $i = 1$ or $n+1$.

We can now claim that if we have j levels with two stars, the worst case is obtained when those j levels are assembled either at the top of the tree or at the bottom.

Let us consider the case where $j = 2$. We can represent the tree S_n^* in the following way:

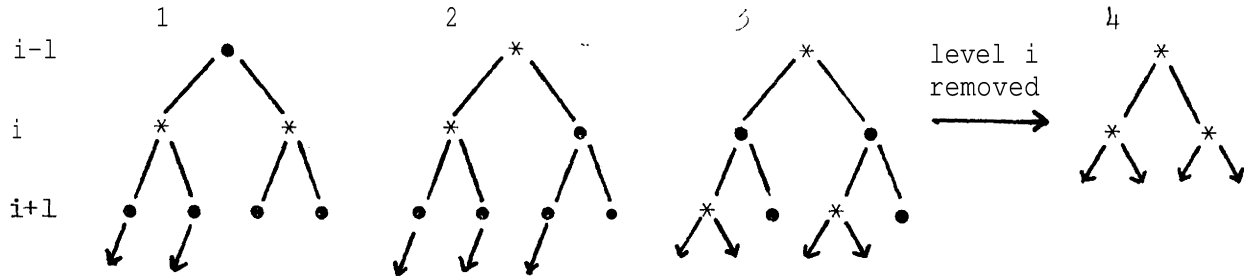


For any i_2 , the tree $S_{i_2-1}^*$ gives the maximum of entries at level i_2 if the level i_1 is equal to 1, then we obtained the following representation: see (B).



Now S_{n-1}^* gives the maximum of buckets if i_2 in S_{n-1}^* equals 1, so in S_n^* , $i_2 = 2$. This procedure can be clearly extended to any number j .

Fact 3. The # of paths generated by the respective positions of two stars, is the following:



1. On the same level, does not affect each other.
2. On two consecutive levels, the second * affects one of the paths generated by the other.
3. At least one level apart, each of the paths generated by the first * is affected by the second one.
4. We can also remark that any level without *s can be removed from the tree without changing the final number of buckets.

The analysis of the worst case for any k follows easily.

For $0 \leq k \leq n+1 - \lceil \frac{n}{2} \rceil$, the number of stars in this range allows

us to place all the *s one level apart so, going down the tree, each new * affects every path generated by the others. So

$$\underline{w_n(k) = 2^k} \quad (\text{equals the upper bound}).$$

We can remark, also, that if we cancel all the levels without *s, we get a tree with only *s and k levels.

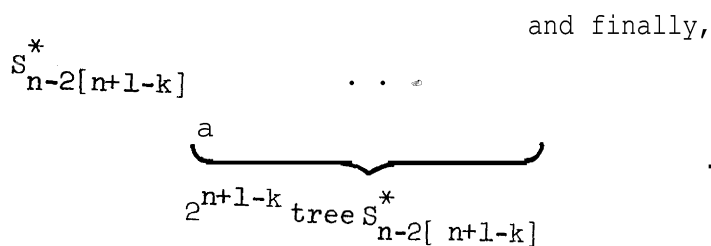
For $n+1 - \lceil \frac{n}{2} \rceil \leq k \leq n+1$, the number of *s is such that $n+1-k$

levels have no *s since no levels have two *s for the worst case (see Fact 3). Thus, we can cancel those $n+1-k$ levels and we are left with a tree with $n-(n+1-k)$ levels in which $n+1-k$ levels have two stars. By Fact 2, these $n+1-k$ levels with two *s have to be at the top of the tree for the worst case. We can then represent the tree in the following way:



so we get:

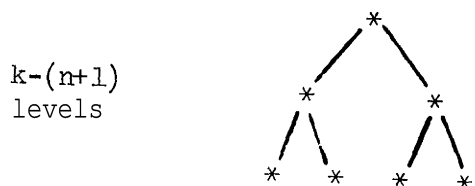
$$w_n(k) = 2^{n+1-k} F_{n-2[n+1-k]+3}$$



$$w_n(k) = 2^{n+1-k} F_{2k-n+1}$$

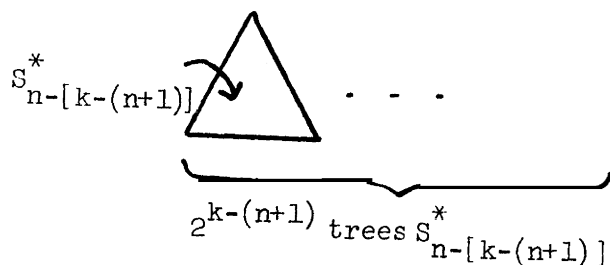
For $n+1 \leq k \leq 2n+1$, the number of *s is such that $k-(n+1)$ levels

have two stars and the others have one *. By Fact 2, the levels with two *s have to be at the top for the worst case. So, the tree can be represented in the following way:



We get:

$$w_n(k) = 2^{k-(n+1)} F_{n-[k-(n+1)]+3}$$



$$w_n(k) = 2^{k-(n+1)} F_{2n+4-k}$$

4. Searching Algorithm.

As we have previously seen, it is not necessary to keep the tree in memory because the value of the nodes may be computed when going down this tree.

We will use the following variables:

L = the level in the tree, beginning at zero.

B = the bucket address which is computed during the search.

$QUERY(i)$ = the i -th bit of the query (starting at zero).

The algorithm is similar to a binary search. When getting to a node, the bit of the query specified in this node is tested. If it is zero, the left subtree is explored and if it is one, the right subtree is explored. In the case of an unspecified bit, the current level and bucket address at this node are saved on a stack, then the left subtree is explored. When getting to a leaf, the bucket address is stored in a table. If the stack is empty, the search is completed; otherwise, the level and bucket address of a node are popped from the stack and the right subtree of this node is explored.

Algorithm S [Searching all the bucket addresses corresponding to a given query.]

1. [INITIALIZE.] Set $i \leftarrow 0$, $B \leftarrow 0$, $L \leftarrow 0$.
2. [TEST BIT OF QUERY.] set $L \leftarrow L+1$, set $B \leftarrow 2B$. If $query(i) = '1'$, go to 7.
3. [UNSPECIFIED BIT.] If $query(i) = '*'$, push (L, B) on the stack.
4. [MOVE LEFT.] Set $i \leftarrow -L$.
5. [TEST FOR LEAF.] If $L < N$, go to 2; otherwise, store B in the bucket table.
6. [TEST FOR DONE.] If stack empty, the algorithm terminates; otherwise, pop (L, B) from the stack.
7. [MOVE RIGHT.] Set $i \leftarrow 2N+1-L$, set $B \leftarrow B+1$, go to 5.

Algorithm S has been implemented in MIX. We give the MIX program and, in particular, the inner loop corresponding to Algorithm S in Appendix 1.

5. Detailed Study of the Average Search Time.

We will now study the average execution time $U_n(k)$ for computing all the bucket addresses corresponding to a query with k unspecified bits for a given n .

We will consider, in the following argument, two time costs:

C_D is the time cost for testing a node and getting to the next node or leaf (left or right son).

C_S is the additive cost when encountering a "*" for saving the parameters in the stack and then restoring them.

From the MIX program, we can see that C_D has a value of 8 when the son is not a leaf, and 9 when the son is a leaf. We can also see that C_S is equal to 10.

By considering that any time we get to the leaf, except the last time, we have to pop new parameters from the stack, we may take $C_D = 9$ and add the extra cost when getting to a leaf, to C_S . Then, we will take in the following study:

$$\left. \begin{array}{l} C_S = 11 \\ C_D = 8 \end{array} \right\} \text{ Units of MIX time}$$

Let us now try to express $U_n(k)$ as a function of k and n . For a given n , k can take all the values between 0 and $2n+1$. We have seen the direct representation of Algorithm S by the tree S_n and the relation of the trees for n and $n+1$. Those preceding observations lead us to express the average time for $n+1$ in terms of the average time for n . We can remark that if the key 0 is not specified, the average time spent in the left and right sub-tree cannot be considered equal since the root of each of those subtrees are different.

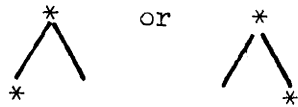
In respect of this last remark, we have isolated two types of average time:

$A_n(k)$ = Average time for n if k keys are unspecified, the first key being specified.

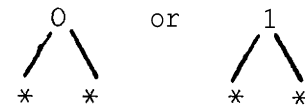
$B_n(k)$ = Average time for n if k keys are unspecified, the first key being unspecified.

Seven different cases, regarding the diverse combinations of the two first keys and the last one, are going to be considered:

Case 0 1



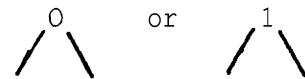
Case ⑤



Case 0 2



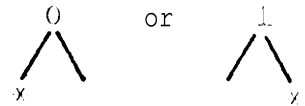
Case ⑥



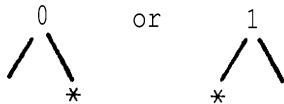
Case 0 3



Case ⑦



Case 0 4



Case ①

$$B_{n+1}^{①}(k) = A_n(k-2) + B_n(k-1) + 2C_D + C_S$$

since we have to explore both subtrees, the average time in one is $A_n(k-2)$ and $B_n(k-1)$ in the other one. We have also to follow the two edges corresponding to the time $2C_D$ and to pop in and out the stack corresponding to C_S . In the same way we can compute the

$A_{n+1}^{(i)}(k)$ (or $B_{n+1}^{(i)}(k)$) in each case (i) :

$$(2) \quad B_{n+1}^{(2)}(k) = 2B_n^{(k-2)} + 2C_D + C_S$$

$$(3) \quad B_{n+1}^{(3)}(k) = 2A_n^{(k-1)} + 2C_D + C_S$$

$$(4) \quad A_{n+1}^{(4)}(k) = A_n^{(k-1)} + C_D$$

$$(5) \quad A_{n+1}^{(5)}(k) = B_n^{(k-1)} + C_D$$

$$(6) \quad A_{n+1}^{(6)}(k) = A_n^{(k)} + C_D$$

$$(7) \quad A_{n+1}^{(7)}(k) = B_n^{(k)} + C_D$$

We have now to compute the probability for each case to occur.

[Remark: for $n+1$ we have $2(n+1)+1 = 2n+3$ bits]

Case 1

$$\left\{ \begin{array}{l} \text{Number of possibilities} \\ \text{to have } k \text{ stars with} \\ \text{one star in } 0 \text{ and one} \\ \text{star in } 1 \end{array} \right\} =$$

$$\left\{ \begin{array}{l} \text{Number of} \\ \text{symmetrical} \\ \text{cases} \end{array} \right\} \times \left\{ \begin{array}{l} \text{Number of ways to select} \\ (k-2) \text{ items out of} \\ (2n+3)-3 = 2n \text{ (since the} \\ \text{items } 0, 1 \text{ and } 2n+2 \\ \text{are already chosen} \end{array} \right\} \times \left\{ \begin{array}{l} \text{Number of ways} \\ \text{to fill the} \\ 2n+3-k \text{ left} \\ \text{bits with } 1 \\ \text{or } 0 \end{array} \right\}$$

\uparrow \uparrow \uparrow
 2 $\frac{2n}{k-2}$ 2^{2n+3-k}

so

$$(1) \quad N_{n+1}^{(1)} = 2 \binom{2n}{k-2} 2^{2n+3-k}$$

In the same way we obtain:

$$\begin{aligned} \textcircled{2} N_{n+1} &= \binom{2n}{k-3} 2^{2n+3-k} & \textcircled{3} N_{n+1} &= \binom{2n}{k-1} 2^{2n+3-k} & \textcircled{4} N_{n+1} &= \binom{2n}{k-1} 2^{2n+3-k} \\ \textcircled{5} N_{n+1} &= \binom{2n}{k-2} 2^{2n+3-k} & \textcircled{6} N_{n+1} &= \binom{2n}{k} 2^{2n+3-k} & \textcircled{7} N_{n+1} &= \binom{2n}{k-1} 2^{2n+3-k} . \end{aligned}$$

The total number of ways N_{n+1}^B when the first bit is unspecified is given by:

$$\textcircled{B} N_{n+1} = \binom{2n+2}{k-1} 2^{2n+3-k}$$

The total number of ways N_{n+1}^A when the first bit is specified is given by

$$\textcircled{A} N_{n+1} = \binom{2n+2}{k} 2^{2n+3-k} .$$

We can easily verify that

$$\textcircled{B} N_{n+1} = \sum_{i=1}^3 \textcircled{i} N_{n+1} \quad \textcircled{A} N_{n+1} = \sum_{i=4}^7 \textcircled{i} N_{n+1}$$

and finally that

$$N_{n+1} = \textcircled{B} N_{n+1} + \textcircled{A} N_{n+1} = \binom{2n+3}{k} 2^{2n+3-k} = \left\{ \begin{array}{l} \text{total number of} \\ \text{queries having } k \\ \text{bits unspecified} \\ \text{for } n+1 \end{array} \right\} .$$

We can derive $A_{n+1}(k)$ and $B_{n+1}(k)$

$$A_{n+1}(k) = \frac{1}{\textcircled{A} N_{n+1}} \sum_{i=4}^7 \textcircled{i} N_{n+1} A_{n+1}^{\textcircled{i}}(k) , \quad B_{n+1}(k) = \frac{1}{\textcircled{B} N_{n+1}} \sum_{i=1}^3 \textcircled{i} N_{n+1} B_{n+1}^{\textcircled{i}}(k)$$

and finally $U_{n+1}(k)$

$$U_{n+1}(k) = \frac{1}{N_{n+1}} \left[\overset{\textcircled{A}}{r+1} A_{n+1}(k) + N_{n+1} \overset{\textcircled{B}}{B_{n+1}(k)} \right]$$

If we write $W_n(k) = \binom{2n+1}{k} U_n(k)$, we get after computation the following recurrence formula:

$$\begin{aligned} W_{n+1}(k) &= \alpha_{n+1k} + 2W_n(k-2) + 3W_n(k-1) + W_n(k) \\ \text{with:} \\ \alpha_{nk} &= \binom{2n}{k-1} (2C_D + C_S) + \binom{2n}{k} C_D \\ \text{and the initial condition:} \\ W_0(0) &= C_D, \quad W_0(1) = 2C_D + C_S \\ \text{Assuming} \\ W_n(i) &= 0 \quad \forall i < 0 \quad \text{or} \quad \forall i > 2n+1 \end{aligned} \tag{1}$$

Let us now define the generating function $G_n(z) = \sum_{k \geq 0} W_n(k) z^k$.

From (1) we directly deduce that:

$$G_n(z) = \sum_{k \geq 0} \alpha_{nk} z^k + (2z^2 + 3z + 1) G_{n-1}(z)$$

Using the formula $(1+z)^r = \sum_{k \geq 0} \binom{r}{k} z^k$ we can simplify the sum in α_{nk} . We get:

$$\sum_{k \geq 0} \alpha_{nk} z^k = (\alpha z + \beta) (1+z)^{2n} \quad \text{with } \alpha = 2C_D + C_S \quad \text{and } \beta = C_D$$

When we express $G_{n-1}(z)$ in terms of $G_{n-2}(z)$ and $G_{n-2}(z)$ in terms of $G_{n-3}(z)$, etc., we finally obtain:

$$G_n(z) = (\alpha z + \beta) \sum_{i=0}^{n-1} (2z^2 + 3z + 1)^i (1+z)^{n-i} + (2z^2 + 3z + 1)^n G_0(z)$$

but

$$G_0(z) = \alpha z + \beta$$

so

$$G_n(z) = (\alpha z + \beta) \sum_{i=0}^n (2z^2 + 3z + 1)^i (1+z)^{2(n-i)}$$

$$2z^2 + 3z + 1 = (1+z)(1+2z)$$

so we have

$$G_n(z) = (\alpha z + \beta)(1+z)^{2n} \sum_{i=0}^n \left(\frac{1+2z}{1+z} \right)^i$$

Now using the formula $\sum_{0 \leq k \leq n} z^k = \frac{1-z^{n+1}}{1-z}$ we obtain after computation:

$$G_n(z) = \left(\alpha + \frac{\beta}{z} \right) [(1+z)^n (1+2z)^{n+1} - (1+z)^{2n+1}]$$

We can remark that $(1+2z)^{n+1} = [(1+z)+z]^{n+1} = \sum_{j \geq 0} \binom{n+1}{j} (1+z)^{n+1-j} z^j$.

Then the product $(1+z)^n (1+2z)^{n+1}$ is:

$$(1+z)^n (1+2z)^{n+1} = \sum_{k \geq 0} \left[\sum_{j \geq 0} \binom{n+1}{j} \binom{2n+1-j}{k-j} \right] z^k$$

and finally

$$G_n(z) = \left(\alpha + \frac{\beta}{z} \right) \sum_{k \geq 0} \left[\sum_{j \geq 0} \binom{n+1}{j} \binom{2n+1-j}{k-j} - \binom{2n+1}{k} \right] z^k$$

We get after Computation a direct formula for

$$U_n(k) = \sum_{k=0}^{2n+1} w_n(k) :$$

$$U_n(k) = \alpha(c_{nk} - 1) + \frac{2n+1-k}{k+1} \beta(c_{n, k+1} - 1)$$

with

$$c_{nk} = \frac{1}{\binom{2n+1}{k}} \sum_{j \geq 0} \binom{n+1}{j} \binom{2n+1-j}{2n+1-k} .$$
(2)

We can remark that this formula gives for $k = 0$ and $2n+1$ results that we could have expected:

$$U_n(0) = \beta(n+1)$$

$$U_n(2n+1) = \alpha(2^{n+1} - 1) .$$

Asymptotic Behavior

Using the big- O notation we have:

$$\binom{n+1}{j} = \frac{n^j}{j!} \left(1 + \frac{1+0-1 \dots - (j-2)}{n} + o(n^{-2}) \right)$$

$$= \frac{n^j}{j!} \left(1 + \frac{1 - \frac{1}{2}(j-1)(j-2)}{n} + o(n^{-2}) \right) .$$

Similarly

$$\binom{2n+1}{k} = \frac{(2n)^k}{k!} \left(1 + \frac{1+0-1 \dots - (k-2)}{2n} + o(n^{-2}) \right)$$

and

$$\binom{2n+1-j}{k-j} = \frac{(2n)^{k-j}}{(k-j)!} \left(1 + \frac{\frac{1}{2}(k-1)(k-2) - \frac{1}{2}(j-1)(j-2)}{2n} + o(n^{-2}) \right) .$$

So we get for c_{nk} after computation

$$c_{nk} = \left(\frac{3}{2}\right)^k + \frac{1}{4n} \sum_{j \geq 0} \binom{k}{j} \left(\frac{1}{2}\right)^j \dots\dots\dots + o(n^{-2}) .$$

Using the first and second derivatives of the generating function

$$\sum_{j \geq 0} \binom{k}{j} \left(\frac{1}{2}\right)^j = \left(\frac{3}{2}\right)^k \quad \text{we can give a direct formula for the sum}$$

so

$$c_{nk} = \left(\frac{3}{2}\right)^k + \frac{k}{16n} (7-k) \frac{3}{2}^{k-2} + o(n^{-2})$$

and we finally get for $U_n(k)$

$$\begin{aligned} U_n(k) &= A(k)n + B(k) + o(n^{-1}) \\ \text{with } A(k) &= \frac{2\beta}{k+1} \left(\left(\frac{3}{2}\right)^{k+1} - 1 \right) \\ \text{and } B(k) &= \alpha \left(\left(\frac{3}{2}\right)^k - 1 \right) + \frac{1-k}{1+k} \beta \left(\left(\frac{3}{2}\right)^{k+1} - 1 \right) + \frac{\beta}{8} (6-k) \left(\frac{3}{2}\right)^{k-1} . \end{aligned}$$

We can verify that for $k = 0$, we get $A(0) = B(0) = \beta$

A program has been written to compute the values of $U_n(k)$. The recurrence (1) has been used instead of the final formula (2) for an easier and more efficient implementation. The results for $n \leq 20$ and graphs showing the variation of $U_n(k)$ for a given k or a given n are shown in Appendix 2.

The method is very efficient when the number of *s stays within a reasonable limit. Assuming a 1 μ s cycle time, the algorithm for $n = 15$ (31 bits), will take only 65 ms with 20 *s but will take 1.77 s for $k = 30$.

In any case, the searching time will be negligible compared to the time spent for retrieving the records themselves from secondary storage. As we have already seen in the asymptotic behavior of $U(k)$, the graphs show the almost linearity in n for $U(k)$, k being fixed.

Conclusions

The study of this new technique of retrieval on secondary keys shows two different aspects:

- A very 'good' worst case.
- A very simple and efficient software implementation.

Acknowledgment

The authors would like to thank D. E. Knuth for his advice and encouragement in this work.

References

- [1] W. A. Burkhard, "Partial Match Retrieval File Designs," Computer Science Division, University of California, San Diego. January, 1975.
- [2] D. E. Knuth, The Art of Computer Programming, Vol. 3, Sorting and Searching (Addison-Wesley, 1973), Section 6.5.

Appendix 1

MIX Program and in particular
the inner loop for Algorithm S.

ENTY-1414

| PC | Op | OpC | OpD | OpE | OpF | OpG | OpH | OpI | OpJ | OpK | OpL | OpM | OpN | OpO | OpP | OpQ | OpR | OpS | OpT | OpU | OpV | OpW | OpX | OpY | OpZ | OpAA | OpAB | OpAC | OpAD | OpAE | OpAF | OpAG | OpAH | OpAI | OpAJ | OpAK | OpAL | OpAM | OpAN | OpAO | OpAP | OpAQ | OpAR | OpAS | OpAT | OpAU | OpAV | OpAW | OpAX | OpAY | OpAZ | OpBA | OpBB | OpBC | OpBD | OpBE | OpBF | OpBG | OpBH | OpBI | OpBJ | OpBK | OpBL | OpBM | OpBN | OpBO | OpBP | OpBQ | OpBR | OpBS | OpBT | OpBU | OpBV | OpBW | OpBX | OpBY | OpBZ | OpCA | OpCB | OpCC | OpCD | OpCE | OpCF | OpCG | OpCH | OpCI | OpCJ | OpCK | OpCL | OpCM | OpCN | OpCO | OpCP | OpCQ | OpCR | OpCS | OpCT | OpCU | OpCV | OpCW | OpCX | OpCY | OpCZ | OpDA | OpDB | OpDC | OpDD | OpDE | OpDF | OpDG | OpDH | OpDI | OpDJ | OpDK | OpDL | OpDM | OpDN | OpDO | OpDP | OpDQ | OpDR | OpDS | OpDT | OpDU | OpDV | OpDW | OpDX | OpDY | OpDZ | OpEA | OpEB | OpEC | OpED | OpEE | OpEF | OpEG | OpEH | OpEI | OpEJ | OpEK | OpEL | OpEM | OpEN | OpEO | OpEP | OpEQ | OpER | OpES | OpET | OpEU | OpEV | OpEW | OpEX | OpEY | OpEZ | OpFA | OpFB | OpFC | OpFD | OpFE | OpFF | OpFG | OpFH | OpFI | OpFJ | OpFK | OpFL | OpFM | OpFN | OpFO | OpFP | OpFQ | OpFR | OpFS | OpFT | OpFU | OpFV | OpFW | OpFX | OpFY | OpFZ | OpGA | OpGB | OpGC | OpGD | OpGE | OpGF | OpGG | OpGH | OpGI | OpGJ | OpGK | OpGL | OpGM | OpGN | OpGO | OpGP | OpGQ | OpGR | OpGS | OpGT | OpGU | OpGV | OpGW | OpGX | OpGY | OpGZ | OpHA | OpHB | OpHC | OpHD | OpHE | OpHF | OpHG | OpHH | OpHI | OpHJ | OpHK | OpHL | OpHM | OpHN | OpHO | OpHP | OpHQ | OpHR | OpHS | OpHT | OpHU | OpHV | OpHW | OpHX | OpHY | OpHZ | OpIA | OpIB | OpIC | OpID | OpIE | OpIF | OpIG | OpIH | OpIJ | OpIK | OpIL | OpIM | OpIN | OpIO | OpIP | OpIQ | OpIR | OpIS | OpIT | OpIU | OpIV | OpIW | OpIX | OpIY | OpIZ | OpJA | OpJB | OpJC | OpJD | OpJE | OpJF | OpJG | OpJH | OpJI | OpJJ | OpJK | OpJL | OpJM | OpJN | OpJO | OpJP | OpJQ | OpJR | OpJS | OpJT | OpJU | OpJV | OpJW | OpJX | OpJY | OpJZ | OpKA | OpKB | OpKC | OpKD | OpKE | OpKF | OpKG | OpKH | OpKI | OpKJ | OpKK | OpKL | OpKM | OpKN | OpKO | OpKP | OpKQ | OpKR | OpKS | OpKT | OpKU | OpKV | OpKW | OpKX | OpKY | OpKZ | OpLA | OpLB | OpLC | OpLD | OpLE | OpLF | OpLG | OpLH | OpLI | OpLJ | OpLK | OpLL | OpLM | OpLN | OpLO | OpLP | OpLQ | OpLR | OpLS | OpLT | OpLU | OpLV | OpLW | OpLX | OpLY | OpLZ | OpMA | OpMB | OpMC | OpMD | OpME | OpMF | OpMG | OpMH | OpMI | OpMJ | OpMK | OpML | OpMM | OpMN | OpMO | OpMP | OpMQ | OpMR | OpMS | OpMT | OpMU | OpMV | OpMW | OpMX | OpMY | OpMZ | OpNA | OpNB | OpNC | OpND | OpNE | OpNF | OpNG | OpNH | OpNI | OpNJ | OpNK | OpNL | OpNM | OpNN | OpNO | OpNP | OpNQ | OpNR | OpNS | OpNT | OpNU | OpNV | OpNW | OpNX | OpNY | OpNZ | OpOA | OpOB | OpOC | OpOD | OpOE | OpOF | OpOG | OpOH | OpOI | OpOJ | OpOK | OpOL | OpOM | OpON | OpOO | OpOP | OpOQ | OpOR | OpOS | OpOT | OpOU | OpOV | OpOW | OpOX | OpOY | OpOZ | OpPA | OpPB | OpPC | OpPD | OpPE | OpPF | OpPG | OpPH | OpPI | OpPJ | OpPK | OpPL | OpPM | OpPN | OpPO | OpPP | OpPQ | OpPR | OpPS | OpPT | OpPU | OpPV | OpPW | OpPX | OpPY | OpPZ | OpQA | OpQB | OpQC | OpQD | OpQE | OpQF | OpQG | OpQH | OpQI | OpQJ | OpQK | OpQL | OpQM | OpQN | OpQO | OpQP | OpQQ | OpQR | OpQS | OpQT | OpQU | OpQV | OpQW | OpQX | OpQY | OpQZ | OpRA | OpRB | OpRC | OpRD | OpRE | OpRF | OpRG | OpRH | OpRI | OpRJ | OpRK | OpRL | OpRM | OpRN | OpRO | OpRP | OpRQ | OpRR | OpRS | OpRT | OpRU | OpRV | OpRW | OpRX | OpRY | OpRZ | OpSA | OpSB | OpSC | OpSD | OpSE | OpSF | OpSG | OpSH | OpSI | OpSJ | OpSK |
|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|

```

0206
0206 0045 00 02 49 0206 0000
0001 00 00 52 0207 0000
-CC(1 00 30 29 0204 0000
0000 01 18 37 0209 0000
0024 00 01 53 0210 0000
0024 00 01 53 0211 0000
0024 00 02 45 0212 0000
0024 00 02 45 0213 0000
CC(1 00 18 37 0213 0000
0045 00 02 49 0214 0000
-CC(1 01 00 07 0215 0208
0155 00 00 39 0215 0208
0159 0217 0217

0000).....0000.....0000
QUERY*11110*10 013 014 021 022
BUCKET ADDRESSES :

QUERY*0*0*0*0* 000 001 004 005 006 016 017 020 021 022 024 025
BUCKET ADDRESSES :

QUERY***** 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022
BUCKET ADDRESSES :

QUERY****10111 001 002 006 014 030
BUCKET ADDRESSES :

QUERY*1*1*1* 012 013 015 017 020 021 023 029 031
BUCKET ADDRESSES :

0000).....0000.....0000
QUERY*11110*10 013 014 021 022
BUCKET ADDRESSES :

QUERY*0*0*0*0* 000 001 004 005 006 016 017 020 021 022 024 025
BUCKET ADDRESSES :

QUERY***** 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022
BUCKET ADDRESSES :

QUERY****10111 001 002 006 014 030
BUCKET ADDRESSES :

QUERY*1*1*1* 012 013 015 017 020 021 023 029 031
BUCKET ADDRESSES :

```

Appendix 2

Results and Graphs for $U, (k)$
the average execution time.

| | | CD = | | 8.000000 | | 5 = | | 11.00000 | | | |
|--|--|------|--|----------|--|--------|--|----------|--|--------|--|
| | | K | | 11 | | 12 | | 13 | | 14 | |
| | | 0 | | 56 | | 104 | | 112 | | 120 | |
| | | 1 | | 129 | | 139 | | 146 | | 159 | |
| | | 2 | | 174 | | 186 | | 195 | | 211 | |
| | | 3 | | 235 | | 251 | | 267 | | 283 | |
| | | 4 | | 319 | | 340 | | 361 | | 382 | |
| | | 5 | | 435 | | 463 | | 490 | | 513 | |
| | | 6 | | 535 | | 631 | | 666 | | 705 | |
| | | 7 | | 815 | | 855 | | 914 | | 963 | |
| | | 8 | | 1115 | | 1187 | | 1254 | | 1321 | |
| | | 9 | | 1538 | | 1631 | | 1724 | | 1817 | |
| | | 10 | | 2114 | | 2244 | | 2374 | | 2503 | |
| | | 11 | | 2905 | | 3089 | | 3271 | | 3452 | |
| | | 12 | | 3989 | | 4253 | | 4503 | | 4763 | |
| | | 13 | | 5465 | | 5763 | | 6210 | | 6573 | |
| | | 14 | | 7486 | | 8022 | | 8545 | | 9067 | |
| | | 15 | | 10223 | | 10996 | | 11753 | | 12497 | |
| | | 16 | | 13925 | | 15040 | | 16132 | | 17203 | |
| | | 17 | | 18914 | | 20524 | | 22059 | | 23644 | |
| | | 18 | | 25609 | | 27934 | | 30207 | | 32436 | |
| | | 19 | | 34561 | | 37500 | | 41168 | | 44405 | |
| | | 20 | | 46476 | | 51292 | | 56012 | | 60651 | |
| | | 21 | | 62273 | | 69170 | | 75955 | | 82636 | |
| | | 22 | | 83139 | | 92564 | | 102703 | | 112292 | |
| | | 23 | | 110565 | | 124556 | | 139434 | | 152169 | |
| | | 24 | | 0 | | 156270 | | 184003 | | 205615 | |
| | | 25 | | 0 | | 221157 | | 249105 | | 277009 | |

CC = 8.000000 CS " 11 . 000000

| K | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 |
|----|----|----|--------|--------|---------|---------|---------|----------|----------|----------|-----------|
| 26 | 0 | 0 | 332507 | 372057 | 411508 | 451069 | 490380 | 529506 | 568424 | 607124 | 645608 |
| 27 | 0 | 0 | 442341 | 498171 | 554265 | 610444 | 666596 | 722604 | 778442 | 834067 | 889455 |
| 28 | 0 | 0 | 0 | 664937 | 744195 | 822912 | 903857 | 983860 | 1063801 | 1143559 | 1223193 |
| 29 | 0 | 0 | 0 | 884703 | 998260 | 1108983 | 1222471 | 1336417 | 1450591 | 1564824 | 1678993 |
| 30 | 0 | 0 | 0 | 0 | 1329713 | 1488535 | 1649117 | 1810930 | 1973570 | 2136724 | 2300149 |
| 31 | 0 | 0 | 0 | 0 | 1769445 | 1992363 | 2218807 | 2447897 | 2678936 | 2911378 | 3144794 |
| 32 | 0 | 0 | 0 | 0 | 0 | 2659144 | 2977343 | 3300639 | 3627890 | 3958178 | 4290763 |
| 33 | 0 | 0 | 0 | 0 | 0 | 3538917 | 3984445 | 4439186 | 4901315 | 5369335 | 5842039 |
| 34 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5935210 | 6608773 | 7267077 | 7937137 |
| 35 | 0 | 0 | 0 | 0 | 0 | 0 | 7077861 | 7968380 | 8681314 | 9812955 | 10760203 |
| 36 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 10634509 | 11911410 | 13219955 | 14555308 |
| 37 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 14155749 | 15935672 | 17768128 | 19545216 |
| 38 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 21267216 | 23824736 | 26455600 |
| 39 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 28311520 | 31870176 | 35546624 |
| 40 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 42531188 | 47653200 |
| 41 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 56623055 | 63737632 |
| 42 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 85055464 |
| 43 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 113246128 |

000.32 SECONDS IN EXECUTION

