

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.

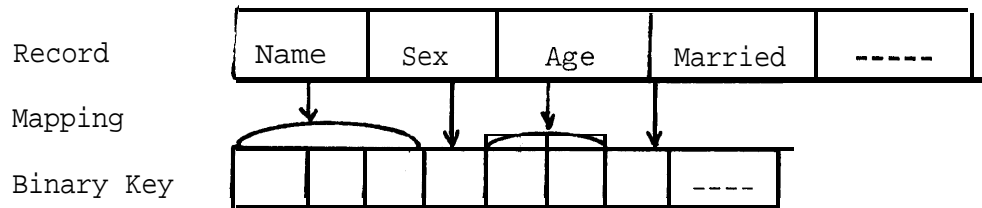
1. Retrieval on Secondary Key [2]

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.

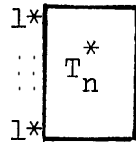
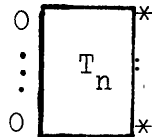
The mapping of binary keys into the bucket addresses is described by a table. Each bucket corresponds to a given value of some of the bits, the others remaining unspecified. Each entry in the tables gives a description of the keys which might be in the corresponding bucket. For example, the bucket corresponding 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: $\begin{matrix} & & 0 \\ & 0 & 0 \\ & 1 & c^1_1 \end{matrix}$

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.

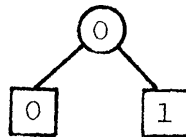
The representation of each table T_n by a tree S_n is then very simple. A node of the tree S_n contains the position of a bit in the query to be checked and a leaf contains the address of a bucket.

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

$n = 0$

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

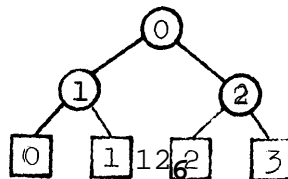
corresponding bucket address is 0 and if the bit is 1, the corresponding bucket address is 1. We can represent this procedure by a very simple tree S_0



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

$n = 1$

T_1 is $\begin{matrix} & 0 & 1 & 2 \\ 0 & 0 & 0 & * \\ 1 & 0 & 1 & * \\ 2 & 1 & * & 0 \\ 3 & 1 & * & 1 \end{matrix}$ the corresponding tree S_1 is



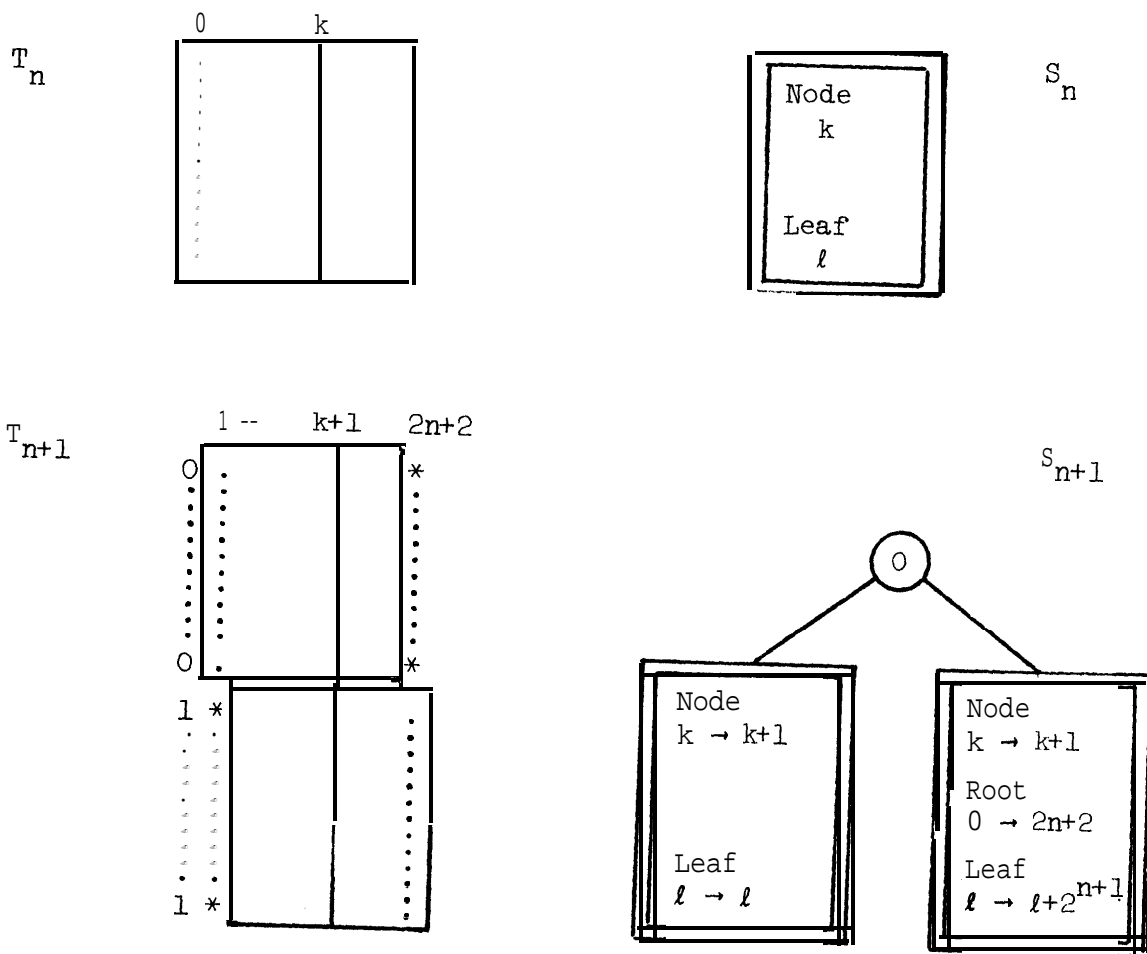
We have seen the passage from T_n to T_{n-1} . Similarly, we can define a transformation on the corresponding trees between S_n and S_{n+1} . Let us suppose that we have our tree S_n corresponding to the table 'I' :

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 $\textcircled{0}$) 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 $\textcircled{0}$ changed into $\textcircled{2n+2}$ 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_4 and the corresponding tree S_4 on the next page)

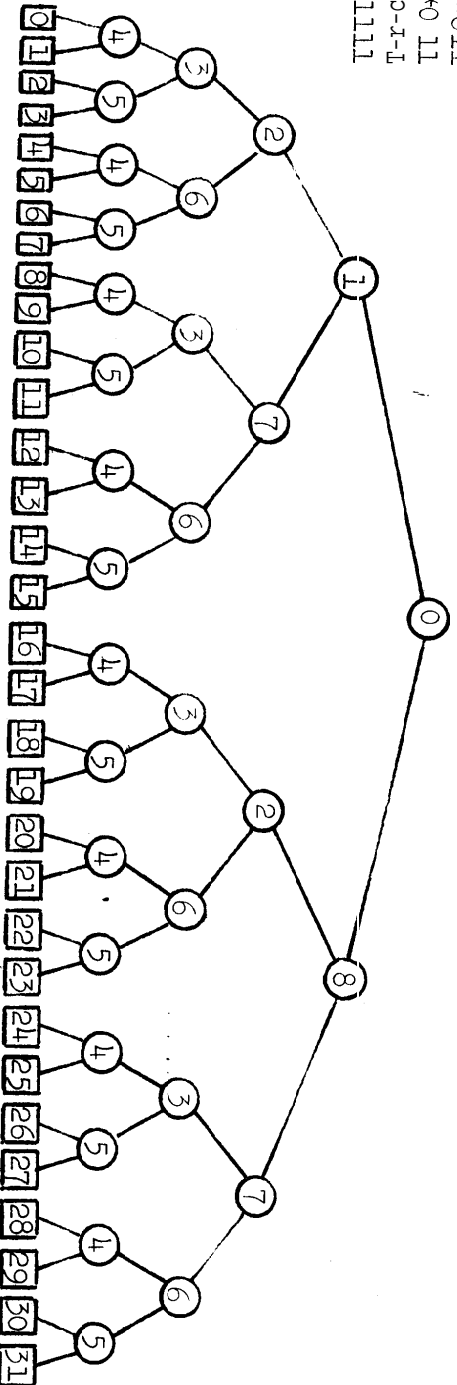
- A -- S_n is perfectly balanced.
- B -- At level i ($0 \leq i < n$), the content of the nodes is either i (denoted by \min_i) or $2n+1-i$ (denoted by \max_i). \min_i is always a left son and \max_i a right son.
- C -- The content of a leaf can be computed directly by the path used from the root adding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

012345678
 00000*****
 00001*****
 0001*0*****
 0001*1*****
 001*0*0***
 001*1*0***
 001**01**
 001**11**
 01*00**0*
 01*01**0*
 01*1*0*0*
 01*1*1*0*
 01**0*01*
 01**1*01*
 01***011*
 01***111*
 1*000**0*
 1*001**0*
 1*01*0**0
 1*01*1**0
 1*1*0*0*0
 1*1*1*0*0
 1*1**01*0
 1*1**11*0
 1**00**01
 1**0 1**01
 1**1*0*01
 1**1*1*01
 1***0*011
 1***1*0 11
 1***0 c-r-T
 1*****1111

Γ_4

7

0
1
2
3
4



S_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:

$$w_n(k) = 2^k \quad \text{for } 0 \leq k \leq n+1 - \left\lceil \frac{n}{2} \right\rceil$$

$$w_n(k) = 2^{n+1-k} F_{2k-n+1} \quad \text{for } n+1 - \left\lceil \frac{n}{2} \right\rceil \leq k \leq n+1$$

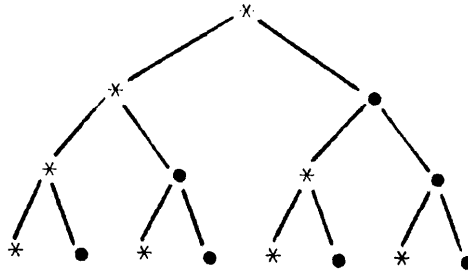
$$w_n(k) = 2^{k-(n+1)} F_{2n+4-k} \quad \text{for } n+1 < k \leq 2n+1$$

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:

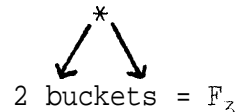


Since in our study of the worst case we considered first the left son to be unspecified, a dot corresponds, in fact, to a bit 0.

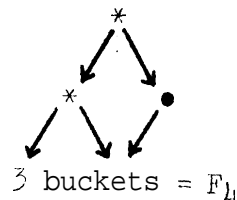
We can now study three different facts:

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^*

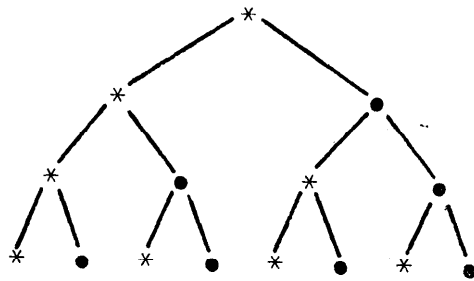


For $n = 1$, we have S_1^*

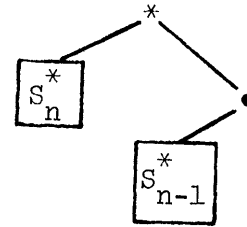


Let us assume that the formula for the number of buckets found, F_{n+3} holds for n and $n-1$, and let us study the case for $n+1$.

S_{n+1}^*



which can be expressed by:



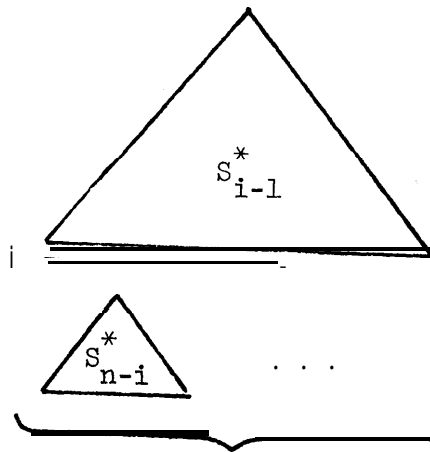
So, # of buckets for S_{n+1}^* = # of buckets for S_n^* + # of buckets for S_{n-1}^*

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

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

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:



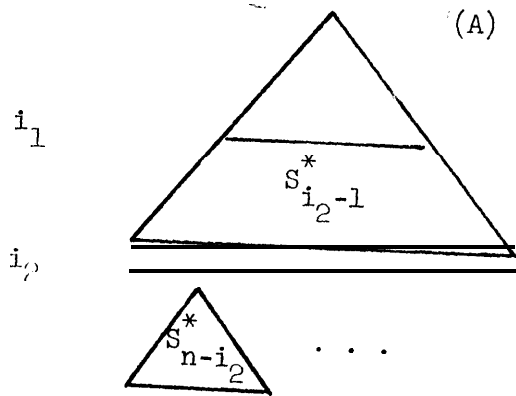
F_{i+2} trees S_{n-i}^* are explored and we have $F_{i+2} F_{n-i+3}$ buckets.

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

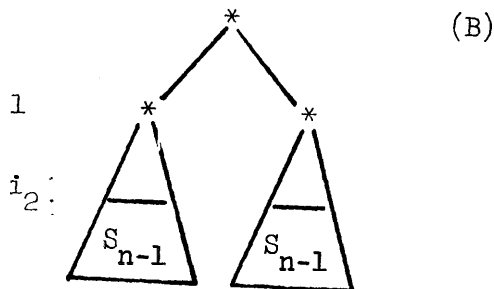
Then it can be easily shown that this product is maximum for $i = 1$ (or $i = n+1$) since $F_i F_{n+1-i}$ varies like $(-1)^{i+1}(F_{n-2i} + 2F_{n-2i+1})$. So, in a tree S_n^* a level i with two *s gives the worst case for $i = 1$ or $n+1$.

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:

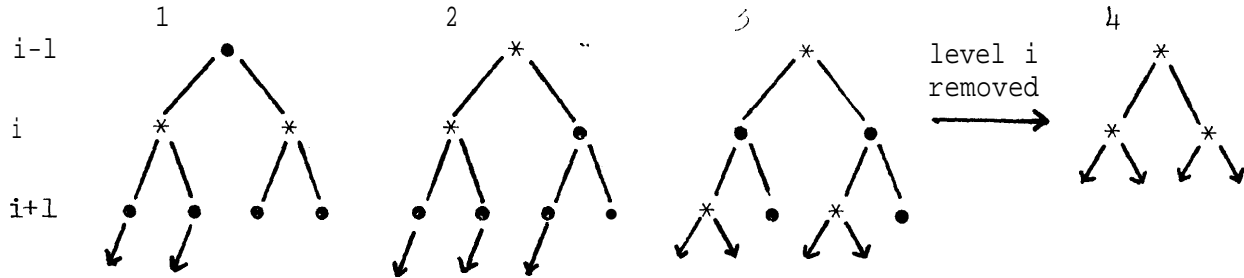


For any i_2 , the tree $S_{i_2-1}^*$ gives the maximum of entries at level i_2 if the level i_1 is equal to 1, then we obtained the following representation: see (B).



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

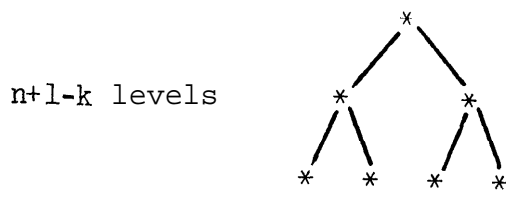
us to place all the *s one level apart so, going down the tree, each new * affects every path generated by the others. So

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

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

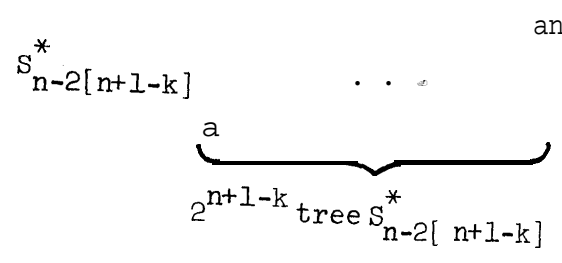
For $n+1 - \lceil \frac{n}{2} \rceil \leq k \leq n+1$, the number of *s is such that $n+1-k$

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:

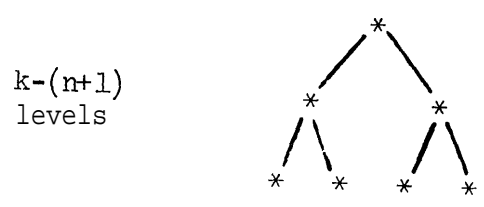
$$w_n(k) = 2^{n+1-k} F_{n-2[n+1-k]+3}$$



and finally,

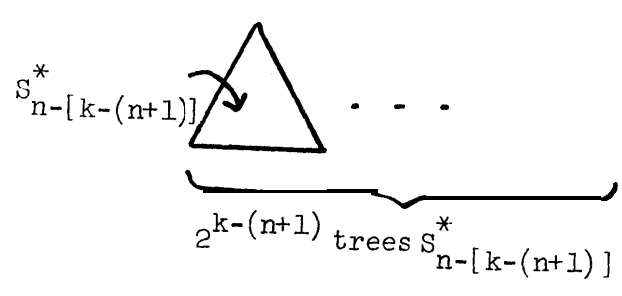
$$w_n(k) = 2^{n+1-k} F_{2k-n+1}$$

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 $i \leftarrow i-1$.
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 $C_D = 8$ and add the extra cost when getting to a leaf, to C_D . Then, we will take in the following study:

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

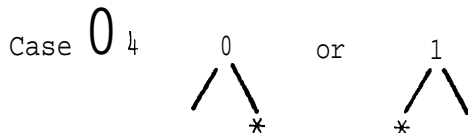
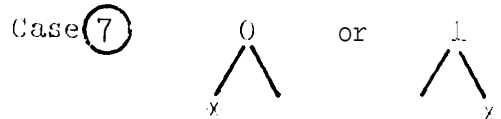
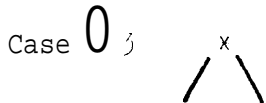
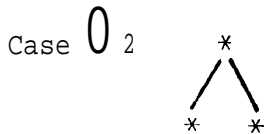
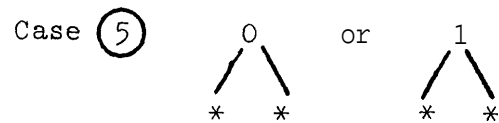
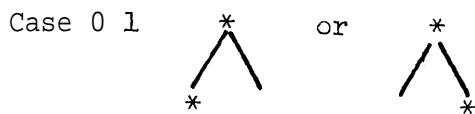
Let us now try to express $U_n(k)$ as a function of k and n . For a given n , k can take all the values between 0 and $2n+1$. We have seen the direct representation of Algorithm S by the tree S_n and the relation of the trees for n and $n+1$. Those preceding observations lead us to express the average 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 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) :

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

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

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

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

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

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

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

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

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

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

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

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

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:

$$W_{n+1}(k) = \alpha_{n+1k} + 2W_n(k-2) + 3W_n(k-1) + W_n(k)$$

with:

$$\alpha_{nk} = \binom{2n}{k-1} (2C_D + C_S) + \binom{2n}{k} C_D$$

and the initial condition:

$$W_0(0) = C_D, \quad W_0(1) = 2C_D + C_S$$

Assuming

$$W_n(i) = 0 \quad \forall i < 0 \quad \text{or} \quad \forall i > 2n+1$$

(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:

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

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:

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

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:

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

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

$$U_n(k) = \alpha \binom{2n+1}{k}^{-1} W_n(k) + \frac{2n+1-k}{k+1} \beta \binom{2n+1}{n, k+1}^{-1}$$

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:

$$U_n(0) = \beta \binom{n+1}{1}$$

$$U_n(2n+1) = \alpha \binom{2n+1}{2n+1}$$

Asymptotic Behavior

Using the big- O notation we have:

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

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

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

Using the first and second derivatives of the generating function

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

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 $U_n(k)$

$$U_n(k) = A(k)n + B(k) + o(n^{-1})$$

with $A(k) = \frac{2\beta}{k+1} \left(\left(\frac{3}{2}\right)^{k+1} - 1 \right)$

and $B(k) = \alpha \left(\left(\frac{3}{2}\right)^k - 1 \right) + \frac{1-k}{1+k} \beta \left(\left(\frac{3}{2}\right)^{k+1} - 1 \right) + \frac{\beta}{8} (6-k) \left(\frac{3}{2}\right)^{k-1} .$

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.

HEX	ASSEMBLY	LOCATION	OP	ADDRESS
00000000	00000000	0000	ENT3	4
00000001	00000001	0001	ENT3	5
00000002	00000002	0002	ENT3	6
00000003	00000003	0003	ENT3	7
00000004	00000004	0004	ENT3	8
00000005	00000005	0005	ENT3	9
00000006	00000006	0006	ENT3	10
00000007	00000007	0007	ENT3	11
00000008	00000008	0008	ENT3	12
00000009	00000009	0009	ENT3	13
0000000A	0000000A	000A	ENT3	14
0000000B	0000000B	000B	ENT3	15
0000000C	0000000C	000C	ENT3	16
0000000D	0000000D	000D	ENT3	17
0000000E	0000000E	000E	ENT3	18
0000000F	0000000F	000F	ENT3	19
00000010	00000010	0010	ENT3	20
00000011	00000011	0011	ENT3	21
00000012	00000012	0012	ENT3	22
00000013	00000013	0013	ENT3	23
00000014	00000014	0014	ENT3	24
00000015	00000015	0015	ENT3	25
00000016	00000016	0016	ENT3	26
00000017	00000017	0017	ENT3	27
00000018	00000018	0018	ENT3	28
00000019	00000019	0019	ENT3	29
0000001A	0000001A	001A	ENT3	30
0000001B	0000001B	001B	ENT3	31
0000001C	0000001C	001C	ENT3	32
0000001D	0000001D	001D	ENT3	33
0000001E	0000001E	001E	ENT3	34
0000001F	0000001F	001F	ENT3	35
00000020	00000020	0020	ENT3	36
00000021	00000021	0021	ENT3	37
00000022	00000022	0022	ENT3	38
00000023	00000023	0023	ENT3	39
00000024	00000024	0024	ENT3	40
00000025	00000025	0025	ENT3	41
00000026	00000026	0026	ENT3	42
00000027	00000027	0027	ENT3	43
00000028	00000028	0028	ENT3	44
00000029	00000029	0029	ENT3	45
0000002A	0000002A	002A	ENT3	46
0000002B	0000002B	002B	ENT3	47
0000002C	0000002C	002C	ENT3	48
0000002D	0000002D	002D	ENT3	49
0000002E	0000002E	002E	ENT3	50
0000002F	0000002F	002F	ENT3	51
00000030	00000030	0030	ENT3	52
00000031	00000031	0031	ENT3	53
00000032	00000032	0032	ENT3	54
00000033	00000033	0033	ENT3	55
00000034	00000034	0034	ENT3	56
00000035	00000035	0035	ENT3	57
00000036	00000036	0036	ENT3	58
00000037	00000037	0037	ENT3	59
00000038	00000038	0038	ENT3	60
00000039	00000039	0039	ENT3	61
0000003A	0000003A	003A	ENT3	62
0000003B	0000003B	003B	ENT3	63
0000003C	0000003C	003C	ENT3	64
0000003D	0000003D	003D	ENT3	65
0000003E	0000003E	003E	ENT3	66
0000003F	0000003F	003F	ENT3	67
00000040	00000040	0040	ENT3	68
00000041	00000041	0041	ENT3	69
00000042	00000042	0042	ENT3	70
00000043	00000043	0043	ENT3	71
00000044	00000044	0044	ENT3	72
00000045	00000045	0045	ENT3	73
00000046	00000046	0046	ENT3	74
00000047	00000047	0047	ENT3	75
00000048	00000048	0048	ENT3	76
00000049	00000049	0049	ENT3	77
0000004A	0000004A	004A	ENT3	78
0000004B	0000004B	004B	ENT3	79
0000004C	0000004C	004C	ENT3	80
0000004D	0000004D	004D	ENT3	81
0000004E	0000004E	004E	ENT3	82
0000004F	0000004F	004F	ENT3	83
00000050	00000050	0050	ENT3	84
00000051	00000051	0051	ENT3	85
00000052	00000052	0052	ENT3	86
00000053	00000053	0053	ENT3	87
00000054	00000054	0054	ENT3	88
00000055	00000055	0055	ENT3	89
00000056	00000056	0056	ENT3	90
00000057	00000057	0057	ENT3	91
00000058	00000058	0058	ENT3	92
00000059	00000059	0059	ENT3	93
0000005A	0000005A	005A	ENT3	94
0000005B	0000005B	005B	ENT3	95
0000005C	0000005C	005C	ENT3	96
0000005D	0000005D	005D	ENT3	97
0000005E	0000005E	005E	ENT3	98
0000005F	0000005F	005F	ENT3	99
00000060	00000060	0060	ENT3	100


```

0206
0206 0045 00 02 49 0206
0001 00 00 52 0207
/-(CC1 00 30 29 0204
CC00 01 18 37 0209
0024 00 01 53 0210
0024 00 03 49 0211
0206 00 02 45 0212
CC21 00 18 37 0213
0045 00 02 49 0214
-(CC1 01 00 07 0215
0156 00 00 39 0216
0158 0217
0000
0000

```

```

0000).....
QUERY*11110*10 0000
013 014 021 022

```

```

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

```

```

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

```

```

QUERY****10111 001 002 006 014 030

```

```

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

```

BUCKET ADDRESSES :

BUCKET ADDRESSES :

BUCKET ADDRESSES :

BUCKET ADDRESSES :

BUCKET ADDRESSES :

Appendix 2

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

EE = 8.000000 CS = 11.00000

K	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
0	0	24	30	40	48	56	64	72	80	88											
1	0	40	50	69	89	79	89	99	101	119											
2	0	63	74	86	99	111	123	136	143	151											
3	0	81	103	124	139	155	171	187	203	219											
4	0	135	155	175	195	215	236	257	277	296											
5	0	189	217	244	271	298	325	352	380	408											
6	0	0	299	337	374	411	448	485	521	558											
7	0	0	405	461	514	565	616	666	716	765											
8	0	0	0	625	701	774	844	914	983	1051											
9	0	0	0	837	949	1054	1155	1252	1349	1444											
10	0	0	0	0	1276	1428	1573	1712	1843	1972											
11	0	0	0	0	1701	1924	2134	2334	2529	2713											
12	0	0	0	0	0	2577	2882	3172	3451	3723											
13	0	0	0	0	0	3429	3873	4294	4697	5093											
14	0	0	0	0	0	0	5178	5788	6373	6937											
15	0	0	0	0	0	0	6885	7769	8615	9431											
16	0	0	0	0	0	0	0	10378	11501	12731											
17	0	0	0	0	0	0	0	13797	15558	17261											
18	0	0	0	0	0	0	0	0	20775	23276											
19	0	0	0	0	0	0	0	0	27521	31133											
20	0	0	0	0	0	0	0	0	0	41567											
21	0	0	0	0	0	0	0	0	0	55269											

	8.000000					11.00000					
K	11	12	13	14	15	16	17	18	19	20	21
0	56	104	112	120	123	136	144	152	160	168	176
1	129	139	146	159	169	179	189	199	209	219	229
2	174	186	199	211	224	237	249	262	275	287	300
3	235	251	267	283	300	316	332	348	364	381	397
4	319	340	361	382	403	424	445	466	487	508	529
5	435	463	490	513	545	573	608	628	656	683	711
6	595	631	668	705	742	778	815	852	888	925	962
7	815	855	914	963	1013	1062	1112	1161	1210	1260	1309
8	1119	1187	1254	1321	1389	1456	1523	1590	1657	1724	1790
9	1538	1631	1724	1817	1909	2001	2093	2185	2277	2369	2462
10	2114	2244	2374	2503	2631	2759	2886	3013	3140	3269	3392
11	2905	3089	3271	3452	3631	3809	3987	4164	4340	4516	4692
12	3989	4259	4503	4763	5019	5267	5517	5765	6013	6259	6505
13	5465	5793	6210	6573	6932	7288	7641	7992	8341	8688	9034
14	7486	8022	8545	9067	9579	10086	10588	11086	11581	12074	12564
15	10223	10996	11753	12497	13231	13955	14673	15384	16089	16793	17487
16	13925	15140	16132	17203	18257	19298	20326	21345	22354	23356	24351
17	18914	20524	22099	23644	25163	26660	28139	29602	31051	32488	33914
18	25609	27934	30207	32436	34628	36787	38917	41023	43108	45174	47224
19	34561	37500	41188	44405	47565	50683	53757	56794	59799	62775	65726
20	46478	51292	56012	60651	65214	69719	74146	78529	82865	87159	91415
21	62273	69178	75959	82636	89212	95697	102100	108427	114686	120885	127028
22	83139	92464	102703	112292	121755	131109	140333	149444	158500	167452	176325
23	110565	124559	139434	152169	165747	179199	192500	205666	218707	231631	244447
24	156270	184003	205615	225615	245378	264937	284301	282463	301264	319913	339417
25	221157	249105	277009	304797	332430	359885	387149	414221	441104	467803	

K	11	12	13	14	15	16	17	18	19	20	21
26	0	0	332507	372057	411508	451069	490380	529506	568424	607124	689605
27	0	0	442341	458171	554265	610444	665596	722604	778442	834067	899455
28	0	0	0	664937	744195	822912	903857	983860	1063801	1143599	1223193
29	0	0	0	884709	996260	1108983	1222471	1336417	1450591	1564824	1678993
30	0	0	0	0	132713	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	3817748	5935210	6605113	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	15935872	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

