# PERFORMING SET OPERATIONS
# BY USING HASHING TECHNIQUES

by

Seiichi Nishihara

Hiroshi Hagiwara

September 12, 1977

INSTITUTE
OF
INFORMATION SCIENCES AND ELECTRONICS

UNIVERSITY OF TSUKUBA

PERFORMING SET OPERATIONS BY USING HASHING TECHNIQUES

Seiichi Nishihara

Institute of Information
Sciences and Electronics
University of Tsukuba
Niihari-gun, Ibaraki 300-31
Japan

Hiroshi Hagiwara

Department of Information
Science
Kyoto University
Sakyo-ku, Kyoto 606
Japan

September 12, 1977

Abstract

Performing set operations is one of the basic techniques in the fields of information retrieval, data structure and data base management.

In this paper, it is shown that hashing techniques can effectively be applied to performing set operations, where each set is a set of keys. Each entry of a hash table contains a key field, a pointer field and a match level indicator field. The last field is used to indicate how well the key satisfies the set formula under consideration.

Some algorithms to process set formulas containing no complementary set are given and the efficiency is proved by some experiments.

# 1. Introduction

One of the purposes of recent data management is the centralized control of many files, so that the redundancy and inconsistency in the stored data may be avoided. Further, queries concerning more than one files can be accepted by unifying files. Most of such operations basically contain set operations especially in information retrieval systems. For instance, when two sets of records satisfy different conditions, the intersection of the two sets is the set of records satisfying the both conditions.

In this paper, a method to perform set operations by using hashing techniques is proposed. First a method for set formulas in disjunctive normal form is described, and then the method is extended to general set formulas. Simple experiments are also executed to estimate the efficiency of the method.

# 2. Performing Set Operations
## 2.1 Definition of Terms

Before describing the method, we shall introduce the terms necessary for the algorithms. The sets appearing in expression of set operations (shortly <u>set formula</u>) are expressed as S or $S_i$ (i=1,2,$\cdots$). Each set is a finite set of keys. We assume the operation to get each key in a set one after another without repetition is available. Let card(S) be the cardinal number of set S. Intersection or union of two sets $S_i$ and $S_j$ are written as $S_i \cdot S_j$ or $S_i + S_j$, respectively. Further, elemen-

tary intersection or elementary union is defined as

$$\bigcap_{i=1}^{m} S_i = S_1 \cdot S_2 \cdots \cdot S_m \quad \text{or} \quad \bigcup_{i=1}^{m} S_i = S_1 + S_2 + \cdots + S_m \,,$$

respectively.    Then a set formula is called to be in disjunctive normal form if it is a union of elementary intersections, i.e.

$$\bigcup_{i=1}^{m} \bigcap_{j=1}^{n(i)} S_{ij} = S_{11} \cdot S_{12} \cdots \cdot S_{1n(1)} + \cdots$$
$$+ S_{m1} \cdot S_{m2} \cdots \cdot S_{m(n)}. \tag{1}$$

In the formula (1), the first set of each elementary intersection (i.e. $S_{11}$, $S_{21}$, $\cdots$, $S_{m1}$) is called a <u>candidate set</u>. Conversely, the last one (i.e. $S_{1n(1)}$, $S_{2n(2)}$, $\cdots$, $S_{mn(m)}$) is called a <u>determinating set</u>.    The keys in the resulting set of a given set formula are called the <u>matched keys</u>.

Up to now, several kinds of hash method have been proposed [1], whose detailed explanation is entirely omitted here. However, we just claim that the hash method adopted in our algorithms works correctly even if a given query key is not in the table.    Thus a hash method such as the separate chaining method[1], the conflict flag method[2] or the predictor method [3] is preferable.

Each entry of the hash table contains at least three fields, as is shown in Fig.1.    A key is hold in the key field.

| match level | key | pointer |
|---|---|---|

Fig.1  Structure of an entry.

The match level field(ML-field) is used to indicate how well the key in the key field agrees with the given set formula. The pointer field, holding a pointer of the chaining method, is of course replaced by the conflict flag or the predictor field in the case other method is adopted.

Now our problem is to get the matched keys of a given set formula by using a hash table.

[Example]

We give here a simple example. Assuming each entry to be initially empty, the algorithm to perform the set formula $S_1 \cdot S_2 \cdot S_3$ is described as follows:

Step 1. Store each element $K_1$ of set $S_1$ into the hash table, setting the ML-field to 1.

Step 2. For each element $K_2$ of set $S_2$, execute the following operation: Search the element $K_2$ in the table. If $K_2$ is found and its ML-field is equal to 1, then change the ML-field to 2.

Step 3. For each element $K_3$ of set $S_3$, execute the following operation: Search the element $K_3$ in the table. If $K_3$ is found and its ML-field is equal to 2, then change the ML-field to 3.

As the conclusion of the algorithm, the key in the entry whose ML-field is equal to 3 is a matched key of the set formula $S_1 \cdot S_2 \cdot S_3$. In this example the necessary and sufficient length of the ML-field is 2 bits.

## 2.2 Method for Disjunctive Normal Form

In this section we give a method to process set formulas in disjunctive normal form. The method consists of two phases as follows.

Phase 1 (Preprocessing- assigning a match value to each set)

Assign serial numbers to all the sets in the set formula from left to right except the determining sets. The number assigned to each set is called the match value of the set. Then assign to each determining set a same value called the final match value, which is the least integer greater than any match value.

For example,

$$S_1 \cdot S_2 \cdot S_3 \cdot S_4 + S_5 \cdot S_6 + S_7 \cdot S_8 \cdot S_9 \; ,$$
$$1 \quad 2 \quad 3 \quad 7 \quad 4 \quad 7 \quad 5 \quad 6 \quad 7$$

where the final match value is 7.

Phase 2 (Execution)

Before giving the algorithm of phase 2, we define some wordings used throughout the paper.

First, "storing set S" means "to store each element of set S into the hash table while initializing the ML-field with the match value assigned to the set. But notice that the entry whose ML-field is not equal to the final match value is treated as empty." In this operation, if the key to be stored already exists in the table and its ML-field is equal to the final match value, then there is no need to store the key again.

Next, "filtering x-valued keys according to set S" means

the following operation: "For each element key of set S, if the key exists in the hash table and the ML-field is greater than or equal to x and less than the match value (say y) of S, then update the ML-field by y. Otherwise, leave as it is."

Here we give the phase 2 algorithm to process the set formula (1):

Step 1. Set $i=1$;

Step 2. Store the i-th candidate set $S_{i1}$;

Step 3. Set $j=2$;

Step 4. Set x equal to the match value of set $S_{ij-1}$;

Step 5. Filter x-valued keys according to set $S_{ij}$;

Step 6. Set $j=j+1$; Is $j>n(i)$? If so go to step 7, if not go back to step 4;

Step 7. Set $i=i+1$; Is $i>m$? If not go back to step 2, if so we are done.

As the result of the algorithm, the key in the entry whose ML-field is equal to the final match value is a mathced key of (1).

In short, the algorithm first stores the keys belonging to a candidate set as candidates of matched keys (step 2), and then reduces them gradually by checking with the sets following after the candidate set (step 5).

The irreducible minimum size of the hash table does not exceed card($\bigcup\limits_{i=1}^{m} S_{i1}$).

Under the situation that the table size is fixed, the several ways to reduce the processing time are considered as,

i) When a set is stored in step 2, choose a set whose cardi-

(5)

nality is as small as possible.    In other words, place

the smallest set at the first position of each elementary

intersection in formula (1).

ii)  Arrange the sets in each elementary intersection in set

formula (1) in such a way that the number of remaining

keys which passed the filtering process of step 5 is

reduced as fast as possible.

iii)  Arrange the elementary intersections of formula (1) in

an ascending order of the size card($\bigcap_{j=1}^{n(i)} S_{ij}$), ($1 \le i \le n$).

Ingeneral, requirement ii) and iii) are hard to insight

in advance.    On the other hand, requirement i) is relatively

easy to satisfy by modifying the algorithm.    Further, the

effect of requirement i) is greater than that of the rest, as

is proved by experiments in the following section.


3.  Some Experiments

3.1  Simulations of a Simple Intersection Operation

In the experiment, a basic set operation to get the inter-

section of three sets $S_A$, $S_B$ and $S_C$ is simulated and evaluated

by employing the separate chaining method with overflow area

[1].    The table size is 2000.    Varying not only the cardi-

nality of each set (, which influences the load factor[1])

but also set formula (, which influences the filtering

sequence of sets),  six cases(case1.1 − case2.3) shown in

Table 1 are executed.    Computer generated pseudorandom numbers

are used as keys.    Before using them, we made $\chi^2$-test for

Poisson distribution at the 5% significance level.

Table 1    The cases executed by simulations.

| cardinal number of each set | set formula | case no. |
|---|---|---|
| $card(S_A)=1000$, $card(S_B)=500$, $card(S_C)=200$, | $S_A \cdot S_B \cdot S_C$ | case 1.1 |
| $card(S_A \cdot S_B)=200$, $card(S_B \cdot S_C)=100$, | $S_C \cdot S_B \cdot S_A$ | case 1.2 |
| $card(S_C \cdot S_A)=40$,    $card(S_A \cdot S_B \cdot S_C)=30$ | $S_C \cdot S_A \cdot S_B$ | case 1.3 |
| $card(S_A)=1800$, $card(S_B)=900$, $card(S_C)=360$, | $S_A \cdot S_B \cdot S_C$ | case 2.1 |
| $card(S_A \cdot S_B)=360$, $card(S_B \cdot S_C)=180$, | $S_C \cdot S_B \cdot S_A$ | case 2.2 |
| $card(S_C \cdot S_A)=72$,    $card(S_A \cdot S_B \cdot S_C)=54$ | $S_C \cdot S_A \cdot S_B$ | case 2.3 |

The efficiency of the algorithm may be expressed in terms of the average number of table access operations (i.e. probes) that occur in hashing processes included in step 2 and step 5. Simulations were programmed and run ten times for each case. The results of the simulations are listed in Table 2, where 'average' columns indicate the values averaged by dividing by the total number of keys, i.e. $card(S_A)+card(S_B)+card(S_C)$.

## 3.2 Analysis of the Experiments

It is easy to estimate analytically the average number of probes needed to get the intersection of three sets. Here we estimate the number of probes needed to process the set formula $S_1 \cdot S_2 \cdot S_3$. Assume that the adopted hash method is the separate chaining method, and each entry in the table is hit as frequently as any other. Then, using Poisson approximation, we can expect that the probability $P(i,x)$ of a cluster of length i is $e^{-x} \cdot x^i/i!$, where x is the load factor[3].

Let M be the table size, and let

$$card(S_1)=k_1, \quad card(S_2)=k_2, \quad card(S_3)=k_3,$$
$$card(S_1 \cdot S_2)=k_2', \quad card(S_1 \cdot S_3)=k_3', \qquad (2)$$
$$card(S_1 \cdot S_2 \cdot S_3)=k,$$

see Fig.2, where k is the number of matched keys. Let $\alpha=k_1/M$.

First estimate the number of probes to store set $S_1$. For an empty entry, probing occurs two times, i.e. to check and to store. Similarly for a chain of length $\ell$, probing occurs $\ell+2$ times where $\ell$ indicates the number of probings to trace the chain. As described above, the probability of a

(8)

Table 2  Summary of results of simulations and analysis.

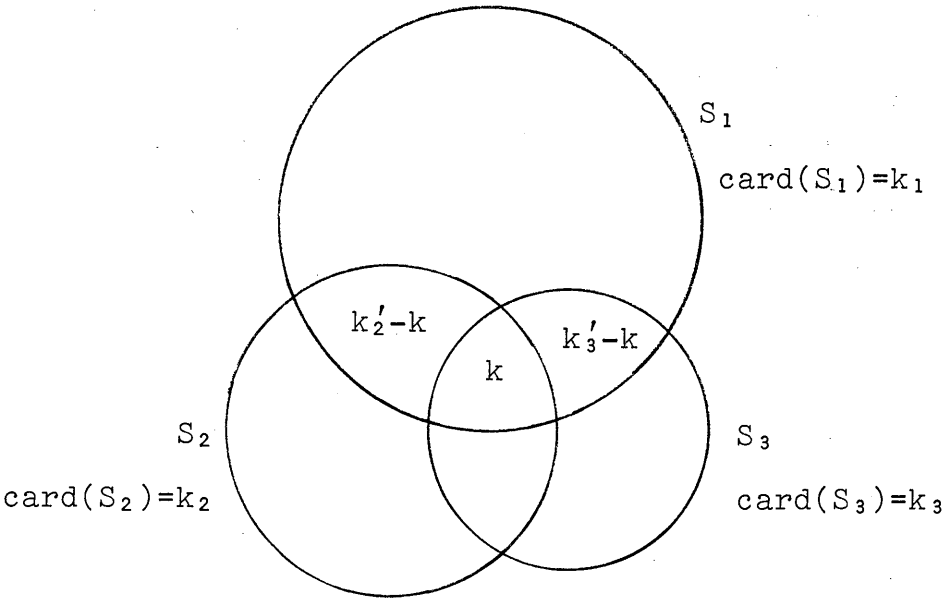| | observed value | | theoretical value | |
|---|---|---|---|---|
| | total | average | total | average |
| case 1.1 | 3331 | 1.96 | 3289 | 1.94 |
| case 1.2 | 2086 | 1.23 | 2054 | 1.21 |
| case 1.3 | 2026 | 1.19 | 1994 | 1.17 |
| case 2.1 | 6612 | 2.16 | 6532 | 2.14 |
| case 2.2 | 3807 | 1.24 | 3746 | 1.22 |
| case 2.3 | 3699 | 1.21 | 3638 | 1.19 |



Fig.2  Venn diagram of $S_1 \cdot S_2 \cdot S_3$.

chain of length $\ell$ is $P(\ell,x)$, where x is the load factor.
Thus the average number of probes needed to store a key when
the load factor is x is given as

$$2 \cdot P(0,x) + \sum_{\ell=1}^{\infty} (\ell+2) \cdot P(\ell,x) = 2+x.$$

Let $T_1$ be the average number of probes needed to store each
key of set $S_1$. Then, by integrating and averaging:

$$T_1 = \frac{1}{\alpha} \int_0^{\alpha} (2+x)dx = 2+\frac{\alpha}{2} , \qquad (3)$$

where $\alpha$ is the load factor after storing process of set $S_1$.

Next consider the keys belonging to set $S_2$. For each
key belonging to the intersection of $S_2$ and $S_1$, the average
number of probes to search is $1+\alpha/2$, and further one more
probing occurs to update the ML-field. Thus average number
of probes is as

$$T_2 = 1+\alpha/2+1 = 2+\alpha/2 . \qquad (4)$$

On the other hand, the average number of probes $T_3$ for the
keys belonging to $S_2-S_1$ is equal to the average number for
the reject operation (i.e. unsuccessful search)[2]:

$$T_3 = P(0,\alpha)+ \sum_{\ell=1}^{\infty} \ell \cdot P(\ell,\alpha) = e^{-\alpha}+\alpha . \qquad (5)$$

Finally, consider the keys belonging to set $S_3$. For
the keys in $S_3$ and in $S_1 \cdot S_2$ (i.e. matched keys), the average
number of probes is equal to $T_2$. For the keys in $S_1-S_2$,
however, the update operation is not necessary. Thus,
average number of probes is as

$$T_4 = 1+\alpha/2 , \qquad (6)$$

which is given in [3]. For the keys in $S_3-S_1$, the average

number of probes is equal to $T_3$.

From the definition (2) and the results (3), (4), (5) and (6), the total number of probing operations T is given as follows:

$$T = T_1 \cdot k_1 + T_2 \cdot (k_2' + k) + T_3 \cdot (k_2 - k_2' + k_3 - k_3') + T_4 \cdot (k_3' - k)$$

$$= k_1 \cdot (2 + \alpha/2) + (k_2 + k_3) \cdot (\alpha + e^{-\alpha})$$

$$+ (k_2' + k_3') \cdot (2 - \alpha/2 - e^{-\alpha}) + k - k_3' . \tag{7}$$

Then the average number of probes E for each key is given as

$$E = T/(k_1 + k_2 + k_3) . \tag{8}$$

The results of theoretical evaluation (7) and (8) are presented in Table 2.

Comparing case 1.1 or case 2.1 with case 1.2 or case 2.2 respectively, the effect of requirement i) is proved. The difference between case 1.2 and 1.3 or between case 2.2 and 2.3 indicates the effect of requirement ii).


## 4. Extending to General Set Formula

### 4.1 Necessity of Extension

Every set formula can be rewritten in an equivalent disjunctive normal form. Thus the algorithm given in section 2 is theoretically applicable to any set formula. Consider, however, an example set formula $S_1 \cdot (S_2 + S_3)$, which may be transformed to $S_1 \cdot S_2 + S_1 \cdot S_3$. Then the processing speed will be considerably slowed down, since set $S_1$ should be stored twice. Therefore, it is desirable that there is an algorithm to execute any set formula in the form as it is,

(11)

which we call <u>direct execution</u>.

In the following section, we give a direct execution algorithm for general set formula containing no complementary set. The fundamental idea is similar to that of section 2.

Here we extend and redefine the term <u>determining set</u>. When a given set formula contains parenthesized subformulas, assume each of them to be a single set. Then the original set formula can be regarded as a disjunctive normal form. Therefore, the determining sets are determined by using the definition given in section 2.1. If the determining set is a parenthesized subformula, then apply the above rule again recursively.

Similarly the term <u>candidate set</u> can also be extended and redefined, but the manner is omitted here.

For example, consider the set formula:

$$(S_1 + S_2 \cdot S_3) \cdot (S_4 + S_5 \cdot (S_6 + S_7)) + S_8 \ , \qquad (9)$$

where the determining sets are $S_4$, $S_6$, $S_7$ and $S_8$, and the candidate sets are $S_1$, $S_2$ and $S_8$. Especially paying attention to subformula $(S_1 + S_2 \cdot S_3)$, the determining sets are $S_1$ and $S_3$, and the candidate sets are $S_1$ and $S_2$.

4.2 Preprocessing of Set Formula (Phase 1)

The rule for assigning a match value to each set is similar to that given in section 2. Roughly speaking, assign serial number from left to right under the restriction that the determining sets in each parenthesized subformula

(12)

should be assigned the same value.

For example, the match values assigned to set formula (9) are as:

$$(S_1+S_2 \cdot S_3) \cdot (S_4+S_5 \cdot (S_6+S_7))+S_8 \qquad (9')$$
$$\quad\ 2 \quad 1 \quad 2 \qquad 4 \quad 3 \quad\ 4 \quad 4 \qquad 4$$

where the final match value is 4.

In section 2, the match value of $S_{ij-1}$ ($1 \le i \le m$, $2 \le j \le n(i)$) is used to filter candidate keys according to $S_{ij}$ in step 5. In the case of general set formula, however, this does not hold.   Therefore newly a value, called <u>check value</u>, is introduced, which is assigned to each set so that the filtering process may work correctly.   The basic rule of assigning check values is as follows:  with respect to each intersection operator (i.e. '·'), the final match value of the left-hand subformula of the operator becomes the check value of the candidate sets of the right-hand subformula. The set that cannot be assigned a check value by the basic rule must be a candidate set of the original set formula and is assigned zero.

For example, the match values and the check values of the set formula (9) are as:

$$\left.\begin{array}{l} (S_1+S_2 \cdot S_3) \cdot (S_4+S_5 \cdot (S_6+S_7))+S_8 \\ \text{match value}\quad 2 \quad 1 \quad 2 \qquad 4 \quad 3 \quad 4 \quad 4 \qquad 4 \\ \text{check value}\quad\ 0 \quad 0 \quad 1 \qquad 2 \quad 2 \quad 3 \quad 3 \qquad 0 \end{array}\right\} (9'')$$

In conclusion, what phase 1 should do is to assign a check value and a match value to each set of the given set formula.   A concrete algorithm of phase 1 is presented in Appendix.

## 4.3 Execution by Using a Hash Table (Phase 2)

After the completion of phase 1, the main execution process performed on a hash table is started. Let $S_i$ be the i-th set from left in the set formula and let check($S_i$) be the check value assigned to set $S_i$. Let m be the number of sets appears in the set formula. Then the algorithm of phase 2 takes a simple form as follows:

[Algorithm of Phase 2]

Step 1. Set i=1;

Step 2. Set x=check($S_i$);

Step 3. If x≠0, then go to step 4. Otherwise, store $S_i$ and go to step 5;

Step 4. Filter x-valued keys according to set $S_i$;

Step 5. Set i=i+1; If i≤m, then go back to step 2. Otherwise, we are done.

As the result of the algorithm, the key in the entry whose ML-field is equal to the final match value is a matched key of the given set formula.

Now let k be the number of intersection operator appearing in a set formula. Then, notice that the final match value is equal to k+1. Thus $\lceil \log_2(k+1) \rceil$ bits are needed for the ML-field to process the set formula.


## 5. Conclusion

We have proposed methods to perform set operations by using a hash table. Two algorithms for disjunctive normal

form and general set formulas are presented.

In this note, the influence of complementary set to the algorithm has not been considered at all, which is the future problem.

References

1) Knuth,D.E. The Art of Computer Programming, Vol.3: Sorting and Searching, Addison-Wesley(1973).

2) Furukawa,K. Hash addressing with conflict flag, Information Processing in Japan, Vol.13(1973),pp.13-18.

3) Nishihara,S. & Hagiwara,H. An open hash method using predictors, ibid.,Vol.15(1975),pp.6-10.

APPENDIX.    An Algorithm of Phase 1

A stack is used as the work area.    Fig.A shows the structure of each entry of the stack, where the fields of <u>set id.</u>, <u>match</u> and <u>check</u> are used to hold a set identifier, a match or final match value and a check value, respectively. The handling of perentheses is performed by using delimiter fields.

Let $p$ indicate the position in the set formula where the process is in progress, and let $a$ indicate the address of the stack.    The position of the first $V$ is 0.    The initial values of $p$, $a$ and $v$ are 0, 1 and 1, respectively.

The algorithm of phase 1 is shown in Table A.    In the algorithm, if the symbols placed at the $p$-th and $(p+1)$-th positions agree with the symbols in the columns of 'present' and 'next' of Table A, then the corresponding operations in 'operation' column are applied.

As an example, the results of the processing of set formula (9) is shown in Fig.B, which coincide with (9'').

| address | set id. | match | check | delimiter |
|---|---|---|---|---|

Fig.A   Structure of an entry of the stack.

Table A    An algorithm of Phase 1.

| next | present | operation |
|------|---------|-----------|
| (    | free    | delimiter(a):=delimiter(a)+1;<br>p:=p+1; |
| SET  | free    | setid(a):=SET;   a:=a+1;<br>p:=p+1; |
| •    | SET     | match(a-1):=v;   check(a):=v;<br>v:=v+1;   p:=p+1; |
| •    | )       | w:=a;<br>L1:w:=w-1;<br>    if match(w)=0   then    match(w):=v;<br>    if delimiter(w)=0  then  go to L1;<br>    delimiter(w):=delimiter(w)-1;<br>    check(a):=match(a-1);<br>    v:=v+1;   p:=p+1; |
| +    | SET     | w:=a;<br>L2:w:=w-1;<br>LL:if  w=1   then   L3:begin<br>                         check(a):=check(w);<br>                         p:=p+1;<br>                         end<br>             else<br>             if  delimiter(w)=0<br>             then  go to L2  else  go to L3; |
| +    | )       | w:=a;<br>L4:w:=w-1;<br>    if  delimiter(w)=0   then  go to L4;<br>    delimiter(w):=delimiter(w)-1;<br>    go to LL; |

*(continued)*

| | | |
|---|---|---|
| ) | SET | `p:=p+1;` |
| | ) | `w:=a;`<br>`L5:w:=w-1;`<br>`    if delimiter(w)=0  then  go to L5;`<br>`    delimiter(w):=delimiter(w)-1;`<br>`    p:=p+1;` |
| ∇ | free | `w:=a;`<br>`L6:w:=w-1;`<br>`    if match(w)≠0  then  go to L7;`<br>`    match(w):=v;`<br>`L7:if w=1  then  go to END`<br>`              else  go to L6;` |

$$(S_1+S_2 \cdot S_3) \cdot (S_4+S_5 \cdot (S_6+S_7))+S_8$$

⇓

p= | 0 | 1 | 2 | 3 | 4 | 5 | 6 | · · · · | 20 | 21 | 22 |

| ∇ | ( | $S_1$ | + | $S_2$ | · | $S_3$ | ) | · | ( | $S_4$ | + | $S_5$ | · | ( | $S_6$ | + | $S_7$ | ) | ) | + | $S_8$ | ∇ |

| address | set id. | match | check | delimiter | |
|---------|---------|-------|-------|-----------|---|
| 10 | | | | | |
| 9 | | | | | |
| 8 | $S_8$ | 4 | 0 | | |
| 7 | $S_7$ | 4 | 3 | | |
| 6 | $S_6$ | 4 | 3 | 1 | 0 |
| 5 | $S_5$ | 3 | 2 | | |
| 4 | $S_4$ | 4 | 2 | 1 | 0 |
| 3 | $S_3$ | 2 | 1 | | |
| 2 | $S_2$ | 1 | 0 | | |
| 1 | $S_1$ | 2 | 0 | 1 | 0 |

Fig.B   An example of preprocessing (Phase 1).

```
1C EXECUTING SET FUNCTIONS BY USING HASHING TECHNIQUES
2C FEBRUARY 1976     BY  S. NISHIHARA
3      COMMON ITAB(3,2000),IOVF(3,1000),IS1(1800),IS2(900),IS3(360)
4      DIMENSION ICOUNT(21)
5C PHASE 0 ****************************************************************
6C INPUT PARAMETERS, CARDINAL NUMBER OF EACH SET
7C AND INCREMENT SIZES.
8      READ(5,1000) N1,N2,N3,N12,N23,N31,N123
9      READ(5,1000) INC1,INC2,INC3
10 1000 FORMAT(7I5)
11C PHASE 1 ****************************************************************
12C GENERATE RANDOM NUMBERS USED AS KEYS.
13      IPOSSN=0
14      IY=1471541918
15      KURI=1
16  602 CONTINUE
17      DO 100 IW=1,N1
18      CALL RANDOM2(YFL,IY)
19      IS1(IW)=IY
20  100 CONTINUE
21      I1=N12+N123
22      DO 101 IW=1,I1
23      IS2(IW)=IS1(IW)
24  101 CONTINUE
25      I2=I1+1
26      DO 102 IW=I2,N2
27      CALL RANDOM2(YFL,IY)
28      IS2(IW)=IY
29  102 CONTINUE
30      DO 103 IW=1,N123
31      IS3(IW)=IS2(IW)
32  103 CONTINUE
33      DO 104 IW=1,N23
34      I3=N123+IW
35      I4=I1+IW
36      IS3(I3)=IS2(I4)
37  104 CONTINUE
38      I5=N123+N23
39      DO 105 IW=1,N31
40      I6=I5+IW
41      I7=I1+IW
42      IS3(I6)=IS1(I7)
43  105 CONTINUE
44      I8=I5+N31+1
45      DO 106 IW=I8,N3
46      CALL RANDOM2(YFL,IY)
47      IS3(IW)=IY
48  106 CONTINUE
49CC
50C PHASE 2 ****************************************************************
51C CALCULATE THE STARTING ADDRESS OF EACH SET IS1, IS2 AND IS3.
52      CALL RANDOM2(YFL,IY)
53      IP1=IY-(IY/N1)*N1
54      CALL RANDOM2(YFL,IY)
55      IP2=IY-(IY/N2)*N2
```

```
56   550 CALL RANDOM2(YFL,IY)
57       IP3=IY-(IY/N3)*N3
58       WRITE(6,1010) IY
59  1010 FORMAT(1H ,25HCURRENT RANDOM NUMBER****,I15)
60 C PHASE 3 ****************************************************************
61 C STORE ALL ELEMENTS IN SET IS1, AND COUNT
62 C THE COLLISIONS FOR X**2 TEST.
63 CHAINING METHOD
64       CALL CLEAR(IPOVF)
65       IPROB=0
66       DO 200 I=1,N1
67       KP=IP1+1
68       KEY=IS1(KP)
69 C STORE THE KEY
70       I1=KEY/3
71       IAD=I1-(I1/2000)*2000+1
72 C** PROBING **   ACCESS THE FIRST KEY
73       IPROB=IPROB+1
74       IF(ITAB(3,IAD) .NE. 0) GO TO 201
75 C** PROBING **   STORE
76       IPROB=IPROB+1
77       ITAB(1,IAD)=1
78       ITAB(3,IAD)=KEY
79       GO TO 202
80 C
81   201 IF(ITAB(2,IAD) .NE. 0) GO TO 203
82 C** PROBING **   POINTER SET
83       IPROB=IPROB+1
84       ITAB(2,IAD)=IPOVF
85 C** PROBING **   STORE
86       IPROB=IPROB+1
87       IOVF(1,IPOVF)=1
88       IOVF(3,IPOVF)=KEY
89       GO TO 204
90 C
91   203 I2=ITAB(2,IAD)
92 C** PROBING **   ACCESS NEXT KEY
93   206 IPROB=IPROB+1
94       IF(IOVF(2,I2) .EQ. 0) GO TO 205
95       I2=IOVF(2,I2)
96       GO TO 206
97 C
98 C** PROBING **   POINTER SET
99   205 IPROB=IPROB+1
100      IOVF(2,I2)=IPOVF
101 C** PROBING **   STORE
102      IPROB=IPROB+1
103      IOVF(1,IPOVF)=1
104      IOVF(3,IPOVF)=KEY
105 C
106   204 IPOVF=IPOVF+1
107      IF(IPOVF .GT. 1000) STOP 9999
108 C
109   202 IP1=IP1+INC1
110      IF(IP1 .GE. N1) IP1=IP1-N1
```

```
111   200 CONTINUE
112       WRITE(6,1001) IPROB
113  1001 FORMAT(1H ,//35H**NUMBER OF PROBES TO STORE SET S1=,I8)
114C
115C PHASE 7 ***********************************************************
116       DO 500 I=1,21
117       ICOUNT(I)=0
118   500 CONTINUE
119C
120       DO 501 I=1,2000
121       LEN=1
122       IF(ITAB(3,I) .EQ. 0) GO TO 502
123       LEN=LEN+1
124       IF(ITAB(2,I) .EQ. 0) GO TO 502
125       LEN=LEN+1
126       J=ITAB(2,I)
127   503 IF(IOVF(2,J) .EQ. 0) GO TO 502
128       LEN=LEN+1
129       J=IOVF(2,J)
130       GO TO 503
131C
132   502 IF(LEN .GT. 21) LEN=21
133       ICOUNT(LEN)=ICOUNT(LEN)+1
134   501 CONTINUE
135C
136       XX=N1
137       XX=XX/2000.0
138       WRITE(6,1500) XX,N1
139  1500 FORMAT(1H ,12H   X**2-TEST,5X,12HLOAD FACTOR=,
140      1F6.3,5X,3HN1=,I6)
141       WRITE(6,1501) ICOUNT(1)
142       WRITE(6,1502) (ICOUNT(I),I=2,11)
143       WRITE(6,1502) (ICOUNT(I),I=12,21)
144  1501 FORMAT(1H ,5X,I7,7H(BLANK))
145  1502 FORMAT(1H ,5X,10I7)
146C****************************************************************
147C PHASE 4 ***********************************************************
148   600 CONTINUE
149       DO 300 I=1,N2
150       KP=IP2+1
151       KEY=IS2(KP)
152C SEARCH THE KEY
153       I1=KEY/3
154       IAD=I1-(I1/2000)*2000+1
155C** PROBING **   ACCESS THE FIRST KEY
156       IPROB=IPROB+1
157       IF(ITAB(3,IAD) .EQ. KEY) GO TO 301
158       IF(ITAB(2,IAD) .EQ. 0) GO TO 302
159       IAD=ITAB(2,IAD)
160C** PROBING **   ACCESS NEXT KEY
161   304 IPROB=IPROB+1
162       IF(IOVF(3,IAD) .EQ. KEY) GO TO 303
163       IF(IOVF(2,IAD) .EQ. 0) GO TO 302
164       IAD=IOVF(2,IAD)
165       GO TO 304
```

```
166C
167C** PROBING **  SET FLAG 2
168  303 IPROB=IPROB+1
169      IOVF(1,IAD)=2
170      GO TO 302
171C
172C** PROBING **  SET FLAG 2
173  301 IPROB=IPROB+1
174      ITAB(1,IAD)=2
175C
176  302 IP2=IP2+INC2
177      IF(IP2 .GE. N2) IP2=IP2-N2
178  300 CONTINUE
179C SET IS2 PROCESSING  COMPLETED
180C
181C PHASE 5 *****************************************************************
182      KOSU=0
183      DO 400 I=1,N3
184      KP=IP3+1
185      KEY=IS3(KP)
186C SEARCH THE KEY
187      I1=KEY/3
188      IAD=I1-(I1/2000)*2000+1
189C** PROBING **  ACCESS THE FIRST KEY
190      IPROB=IPROB+1
191      IF(ITAB(3,IAD) .EQ. KEY) GO TO 401
192      IF(ITAB(2,IAD) .EQ. 0) GO TO 402
193      IAD=ITAB(2,IAD)
194C** PROBING **  ACCESS NEXT KEY
195  404 IPROB=IPROB+1
196      IF(IOVF(3,IAD) .EQ. KEY) GO TO 403
197      IF(IOVF(2,IAD) .EQ. 0) GO TO 402
198      IAD=IOVF(2,IAD)
199      GO TO 404
200C
201  403 IF(IOVF(1,IAD) .NE. 2) GO TO 402
202C** PROBING **  UPDATE FLAG TO 3
203      IPROB=IPROB+1
204      IOVF(1,IAD)=3
205      KOSU=KOSU+1
206      GO TO 402
207C
208  401 IF(ITAB(1,IAD) .NE. 2) GO TO 402
209C** PROBING **  UPDATE FLAG TO 3
210      IPROB=IPROB+1
211      ITAB(1,IAD)=3
212      KOSU=KOSU+1
213C
214  402 IP3=IP3+INC3
215      IF(IP3 .GE. N3) IP3=IP3-N3
216  400 CONTINUE
217C SET IS3 PROCESSING  COMPLETED
218C PHASE 6 *****************************************************************
219      WRITE(6,1002) IPROB,KOSU
220 1002 FORMAT(1H ,//24H TOTAL NUMBER OF PROBES=,
```

```
221        1I8,13H      *RESULT=,I8)
222    601 WRITE(6,1003) KURI
223   1003 FORMAT(1H ,10H**********,I2,10H**********,//)
224        IF(IPOSSN .EQ. 0) KURI=KURI+1
225        IF(KURI .LT. 31) GO TO 602
226        STOP
227        END
```

```
10         SUBROUTINE CLEAR(IPOVF)
20         COMMON ITAB(3,2000),IOVF(3,1000)
30         DO 10 I=1,2000
40         ITAB(1,I)=0
50         ITAB(2,I)=0
60         ITAB(3,I)=0
70      10 CONTINUE
80         DO 11 I=1,1000
90         IOVF(1,I)=0
100        IOVF(2,I)=0
110        IOVF(3,I)=0
120     11 CONTINUE
130        IPOVF=0
140        RETURN
150        END
```

| REPORT DOCUMENTATION PAGE | REPORT NUMBER |
| --- | --- |
| | ISE-TR-77-7 |

TITLE

### PERFORMING SET OPERATIONS BY USING HASHING TECHNIQUES

AUTHOR(s)

Seiichi Nishihara (Institute of Information Sciences
and Electronics, University of
Tsukuba)

Hiroshi Hagiwara  (Department of Information Science,
Kyoto University)

| REPORT DATE | NUMBER OF PAGES |
| --- | --- |
| September 12, 1977 | 24 |

| MAIN CATEGORY | CR CATEGORIES |
| --- | --- |
| Data Management | 4.33, 4.34, 3.73 |

KEY WORDS

set processing, hashing, scatter storage, database,
data manipulation, information retrieval

ABSTRACT

Performing set operations is  one of the basic techniques in
the fields of information retrieval, data structure and data
base management.
    In this paper,  it is shown that hashing techniques can
effectively be applied  to performing set operations,  where
each set is a set of keys.  Each entry of a hash table con-
tains a key field,  a pointer field  and a match level field.
The last field is  used  to  indicate  how  well  the  key
satisfies the set formula under consideration.
    Some algorithms  to process  set formulas containing no
complementary set  are given  and  the efficiency  is proved
by some experiments.

SUPPLEMENTARY NOTES