

# Using a Hash-Based Method with Transaction Trimming for Mining Association Rules

Jong Soo Park, *Member, IEEE*,  
Ming-Syan Chen, *Senior Member, IEEE*, and Philip S. Yu, *Fellow, IEEE*

**Abstract**—In this paper, we examine the issue of mining association rules among items in a large database of sales transactions. Mining association rules means that, given a database of sales transactions, to discover all associations among items such that the presence of some items in a transaction will imply the presence of other items in the same transaction. The mining of association rules can be mapped into the problem of discovering large itemsets where a large itemset is a group of items that appear in a sufficient number of transactions. The problem of discovering large itemsets can be solved by constructing a candidate set of itemsets first, and then, identifying—within this candidate set—those itemsets that meet the large itemset requirement. Generally, this is done iteratively for each large  $k$ -itemset in increasing order of  $k$ , where a large  $k$ -itemset is a large itemset with  $k$  items. To determine large itemsets from a huge number of candidate sets in early iterations is usually the dominating factor for the overall data-mining performance. To address this issue, we develop an effective algorithm for the candidate set generation. It is a hash-based algorithm and is especially effective for the generation of candidate set for large 2-itemsets. Explicitly, the number of candidate 2-itemsets generated by the proposed algorithm is, in orders of magnitude, smaller than that by previous methods—thus resolving the performance bottleneck. Note that the generation of smaller candidate sets enables us to effectively trim the transaction database size at a much earlier stage of the iterations, thereby reducing the computational cost for later iterations significantly. The advantage of the proposed algorithm also provides us the opportunity of reducing the amount of disk I/O required. Extensive simulation study is conducted to evaluate performance of the proposed algorithm.

**Index Terms**—Data mining, association rules, hashing, performance analysis.

## 1 INTRODUCTION

RECENTLY, mining of databases has attracted a growing amount of attention in database communities due to its wide applicability in retail industry to improving marketing strategy. As pointed out in [5], the progress in bar-code technology has made it possible for retail organizations to collect and store massive amounts of sales data. Catalog companies can also collect sales data from the orders they receive. A record in such data typically consists of the transaction date, the items bought in that transaction, and possibly also customer-ID if such a transaction is made via the use of a credit card or any kind of customer card. It is noted that analysis of past transaction data can provide very valuable information on customer buying behavior, and thus improve the quality of business decisions (such as what to put on sale, which merchandises to be placed on shelves together, how to customize marketing programs, to name a few). It is essential to collect a sufficient amount of sales data (say, over last 30 days) before we can draw any meaningful conclusion from them. As a result, the amount

of these sales data tends to be huge. It is therefore important to devise efficient algorithms to conduct mining on these data. The requirement to process large amount of data distinguishes data mining in the database context from its study in the AI context.

One of the most important data-mining problems is mining association rules [3], [5], [9], [10]. Specifically, given a database of sales transactions, one would like to discover all associations among items such that the presence of some items in a transaction will imply the presence of other items in the same transaction. The problem of mining association rules in the context of database was first explored in [3]. In this pioneering work, it is shown that mining association rules can be decomposed into two subproblems. A set of items is called an itemset. First, we need to identify all itemsets that are contained in a sufficient number of transactions above the minimum (support) requirement. These itemsets are referred to as *large itemsets*. Second, once all large itemsets are obtained, the desired association rules can be generated in a straightforward manner. Subsequent work in the literature followed this approach and focused on the large itemset generations.

Various algorithms have been proposed to discover the large itemsets [3], [5], [10]. Generally speaking, these algorithms first construct a candidate set of large itemsets based on some heuristics, and then discover the subset that indeed contains large itemsets. This process can be done iteratively in the sense that the large itemsets discovered in one iteration will be used as the basis to generate the

• J.S. Park is with the Department of Computer Science, Sungshin Women's University, Seoul, Korea. E-mail: jpark@cs.sungshin.ac.kr.

• M.-S. Chen is with the Electrical Engineering Department, National Taiwan University, Taipei, Taiwan, Republic of China. E-mail: mschen@cc.ee.ntu.edu.tw.

• P.S. Yu is with the IBM Thomas J. Watson Research Center, Yorktown Heights, New York, NY 10598. E-mail: psyu@watson.ibm.com.

Manuscript received 18 Aug. 1995; revised 2 July 1996.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number 104423.

candidate set for the next iteration. For example, in [5], at the  $k$ th iteration, all large itemsets containing  $k$  items, referred to as large  $k$ -itemsets, are generated. In the next iteration, to construct a candidate set of large  $(k + 1)$ -itemsets, a heuristic is used to expand some large  $k$ -itemsets into a  $(k + 1)$ -itemset, if certain constraints are satisfied.<sup>1</sup>

The heuristic to construct the candidate set of large itemsets is crucial to performance. Clearly, in order to be efficient, the heuristic should only generate candidates with high likelihood of being large itemsets because for each candidate, we need to count its appearances in all transactions. The larger the candidate set, the more processing cost required to discover the large itemsets. As previously reported in [5] and also attested by our experiments, the processing in the initial iterations in fact dominates the total execution cost. A performance study was provided in [5] to compare various algorithms of generating large itemsets. It was shown that for these algorithms the candidate set generated during an early iteration is generally, in orders of magnitude, larger than the set of large itemsets it really contains. *Therefore, the initial candidate set generation, especially for the large 2-itemsets, is the key issue to improve the performance of data mining.*

Another performance related issue is on the amount of data that has to be scanned during large itemset discovery. A straightforward implementation would require one pass over the database of all transactions for each iteration. Note that as  $k$  increases, not only is there a smaller number of large  $k$ -itemsets, but also there are fewer transactions containing any large  $k$ -itemsets. Reducing the number of transactions to be scanned and trimming the number of items in each transaction can improve the data-mining efficiency in later stages. In [5], two alternative approaches are considered: Apriori and AprioriTid. In the Apriori algorithm, each iteration requires one pass over the database. In the AprioriTid algorithm, the database is not scanned after the first pass. Rather, the transaction ID and candidate large  $k$ -itemsets present in each transaction are generated in each iteration. This is used to determine the large  $(k + 1)$ -itemsets present in each transaction during the next iteration. It was found that in the initial stages, Apriori is more efficient than AprioriTid, since there are too many candidate  $k$ -itemsets to be tracked during the early stages of the process. However, the reverse is true for later stages. A hybrid algorithm of the two algorithms was also proposed in [5] and shown to lead to better performance in general. The challenge to operate the hybrid algorithm is to determine the switch-over point.

In this paper, we shall develop an algorithm DHP (standing for direct hashing and pruning) for efficient large itemset generation. Specifically, DHP proposed has three major features:

- 1) efficient generation for large itemsets,
- 2) effective reduction on transaction database size, and
- 3) the option of reducing the number of database scans required.

As will be seen later, by utilizing a hash technique, DHP is very efficient for the generation of candidate large itemsets,

in particular for the large 2-itemsets, where the number of candidate large itemsets generated by DHP is, in orders of magnitude, smaller than that by previous methods, thus greatly improving the performance bottleneck of the whole process. In addition, DHP employs effective pruning techniques to progressively reduce the transaction database size. As observed in prior work [5], during the early iterations, tracking the candidate  $k$ -itemsets in each transaction is ineffective, since the cardinality of such  $k$ -itemsets is very large. Note that the generation of smaller candidate sets by DHP enables us to effectively trim the transaction database at a much earlier stage of the iterations, i.e., right after the generation of large 2-itemsets, thereby reducing the computational cost for later iterations significantly. It will be seen that by exploiting some features of association rules, not only the number of transactions, but also the number of items in each transaction can be substantially reduced.

Moreover, the advantage of DHP provides us an opportunity of avoiding database scans in some passes so as to reduce the disk I/O cost involved. Explicitly, due to small candidate sets it deals with, DHP has the option of using a candidate set to generate subsequent candidate sets, and delaying the determination of large itemsets to a later pass when the database is scanned. Since DHP does not scan the database to obtain large itemsets in each pass, some database scans are saved. Extensive experiments are conducted to evaluate the performance of DHP. As shown by our experiments, with a slightly higher cost in the first iteration due to the generation of a hash table, DHP incurs significantly smaller execution times than Apriori<sup>2</sup> in later iterations, not only in the second iteration when a hash table is used by DHP to facilitate the generation of candidate 2-itemsets, but also in later iterations when the same procedure for large itemset generation is employed by both algorithms, showing the advantage of effective database trimming by DHP. Sensitivity analysis for various parameters is conducted.

We mention in passing that the discovery of association rules is also studied in the AI context [15] and there are various other aspects of data mining explored in the literature. Classification is an approach of trying to develop rules to group data tuples together based on certain common characteristics. This has been explored both in the AI domain [16] and in the context of databases [2], [7], [8]. Mining in spatial databases was conducted in [14]. Another source of data mining is on ordered data, such as stock market and point of sales data. Interesting aspects to explore include searching for similar sequences [1], [17], e.g., stocks with similar movement in stock prices, and sequential patterns [6], e.g., grocery items bought over a set of visits in sequence. It is noted that some techniques such as sampling were proposed in [4], [13] to improve mining efficiency. However, the emphasis in those works is to devise methods that can learn knowledge efficiently from sampled subsets, and is thus intrinsically different from the approach in this paper that improves the efficiency of mining

2. A comprehensive study on various algorithms to determine large itemsets is presented in [5], where the Apriori algorithm is shown to provide the best performance during the initial iterations. Hence, Apriori is used as the base algorithm to compare with DHP in this study.

1. A detailed description of the algorithm used in [5] is given in Section 2.

the entire transaction database via candidate itemset reduction and transaction trimming. Additionally, it is worth mentioning that in [5], the hybrid algorithm has the option of switching from Apriori to another algorithm AprioriTid after early passes for better performance. For ease of presentation of this paper, such an option is not adopted here. Nevertheless, the benefit of AprioriTid in later passes is complementary to the focus of DHP on initial passes.

This paper is organized as follows. A detailed problem description is given in Section 2. The algorithm DHP proposed for generating large itemsets is described in Section 3. Performance results are presented in Section 4. Section 5 contains the summary.

## 2 PROBLEM DESCRIPTION

Let  $I = \{i_1, i_2, \dots, i_m\}$  be a set of literals, called items. Let  $D$  be a set of transactions, where each transaction  $T$  is a set of items such that  $T \subseteq I$ . Note that the quantities of items bought in a transaction are not considered, meaning that each item is a binary variable representing if an item was bought. Each transaction is associated with an identifier, called TID. Let  $X$  be a set of items. A transaction  $T$  is said to contain  $X$  if and only if  $X \subseteq T$ . An association rule is an implication of the form  $X \Rightarrow Y$ , where  $X \subset I$ ,  $Y \subset I$ , and  $X \cap Y = \emptyset$ . The rule  $X \Rightarrow Y$  holds in the transaction set  $D$  with *confidence*  $c$  if among those transactions that contain  $X$   $c\%$  of them also contain  $Y$ . The rule  $X \Rightarrow Y$  has *support*  $s$  in the transaction set  $D$  if  $s\%$  of transactions in  $D$  contain  $X \cup Y$ .

As mentioned before, the problem of mining association rules is composed of the following two steps:

- 1) Discover the large itemsets, i.e., all sets of itemsets that have transaction support above a predetermined minimum support  $s$ .
- 2) Use the large itemsets to generate the association rules for the database.

The overall performance of mining association rules is, in fact, determined by the first step. After the large itemsets are identified, the corresponding association rules can be derived in a straightforward manner. In this paper, we shall develop an algorithm to deal with the first step, i.e., discovering large itemsets from the transaction database. Readers interested in more details for the second step are referred to [5]. As a preliminary, we shall describe the method used in the prior work Apriori for discovering the large itemsets from a transaction database given in Fig. 1.<sup>3</sup> As pointed out earlier, a comprehensive study on various algorithms to determine large itemsets is presented in [5], where the Apriori algorithm is shown to provide the best performance during the initial iterations. Hence, Apriori is used as the base algorithm to compare with DHP. In Apriori [5], in each iteration (or each pass) it constructs a candidate set of large itemsets, counts the number of occurrences of each candidate itemset, and then determine large itemsets based on a predetermined minimum support. In the first iteration, Apriori simply scans all the transactions to count the number of occurrences for each item. The set of candidate 1-

Database $D$	
TID	Items
100	A C D
200	B C E
300	A B C E
400	B E

Fig. 1. An example transaction database for data mining.

itemsets,  $C_1$ , obtained is shown in Fig. 2. Assuming that the minimum transaction support required is 2, the set of large 1-itemsets,  $L_1$ , composed of candidate 1-itemsets with the minimum support required, can then be determined.

To discover the set of large 2-itemsets, in view of the fact that any subset of a large itemset must also have minimum support, Apriori uses  $L_1 * L_1$  to generate a candidate set of itemsets  $C_2$  using the a priori candidate generation, where  $*$  is an operation for concatenation.  $C_2$  consists of  $|L_1|(|L_1| - 1)/2$  2-itemsets. Note that when  $|L_1|$  is large,  $|L_1|(|L_1| - 1)/2$  becomes an extremely large number. Next, the four transactions in  $D$  are scanned and the support of each candidate itemset in  $C_2$  is counted. The middle table of the second row in Fig. 2 represents the result from such counting in  $C_2$ . A hash tree is usually used for a fast counting process [11].<sup>4</sup> The set of large 2-itemsets,  $L_2$ , is therefore determined based on the support of each candidate 2-itemset in  $C_2$ .

The set of candidate itemsets,  $C_3$ , is generated from  $L_2$  as follows. From  $L_2$ , two large 2-itemsets with the same first item, such as  $\{BC\}$  and  $\{BE\}$ , are identified first. Then, Apriori tests whether the 2-itemset  $\{CE\}$ , which consists of their second items, constitutes a large 2-itemset or not. Since  $\{CE\}$  is a large itemset by itself, we know that all the subsets of  $\{BCE\}$  are large and then  $\{BCE\}$  becomes a candidate 3-itemset. There is no other candidate 3-itemset from  $L_2$ . Apriori then scans all the transactions and discovers the large 3-itemsets  $L_3$  in Fig. 2. Since there is no candidate 4-itemset to be constituted from  $L_3$ , Apriori ends the process of discovering large itemsets.

## 3 DIRECT HASHING WITH EFFICIENT PRUNING FOR FAST DATA MINING

In this section, we shall describe algorithm DHP (which stands for direct hashing and pruning) for efficient large itemset generation. As pointed out earlier, algorithm DHP proposed has three major features: one is efficient generation for large itemsets, another is effective reduction on transaction database size, and the third is database scan reduction. We shall describe in Section 3.1 the main flow of algorithm DHP, and explain its first feature, i.e., efficient generation for large itemsets. The feature of reducing database size will be described in Section 3.2. For better readability, the scan-reduction technique of DHP is described in Section 4, together with its performance results.

4. The use of a hash tree to count the support of each candidate itemset is a common feature for the two algorithms (i.e., Apriori and DHP) discussed in this paper, and should not be confused with the hash technique used by DHP to generate candidate itemsets.

3. This example database is extracted from [5].

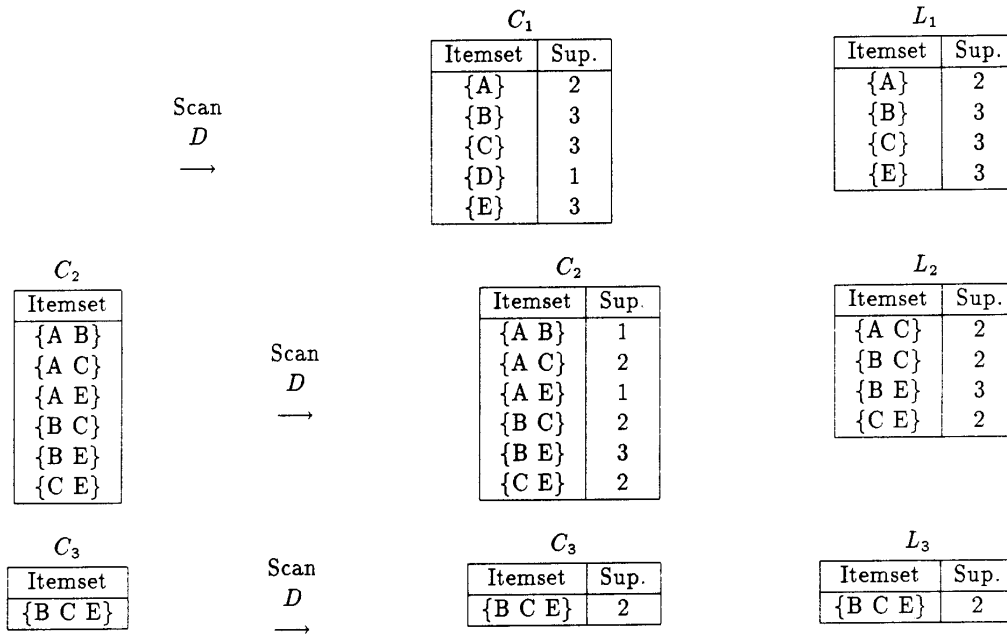


Fig. 2. Generation of candidate itemsets and large itemsets.

### 3.1 Algorithm DHP

As illustrated in Section 2, in each pass we use the set of large itemsets,  $L_i$ , to form the set of candidate large itemsets  $C_{i+1}$  by joining  $L_i$  with  $L_i$  on  $(i-1)$  (denoted by  $L_i * L_i$ ) common items for the next pass. We then scan the database and count the support of each itemset in  $C_{i+1}$  so as to determine  $L_{i+1}$ . In general, the more itemsets in  $C_i$ , the higher the processing cost of determining  $L_i$  will be. Note that

$$|C_2| = \binom{|L_1|}{2}$$

in Apriori. The step of determining  $L_2$  from  $C_2$  by scanning the whole database and testing each transaction against a hash tree built by  $C_2$  is hence very expensive. By constructing a significantly smaller  $C_2$ , DHP can also generate a much smaller  $D_3$  to derive  $C_3$ . Without the smaller  $C_2$ , the database cannot be effectively trimmed. (This is one of the reasons that Apriori-TID in [5] is not effective in trimming the database until later stages.) After this step, the size of  $L_i$  decreases rapidly as  $i$  increases. A smaller  $L_i$  leads to a smaller  $C_{i+1}$ , and thus a smaller corresponding processing cost.

In essence, algorithm DHP presented in Fig. 3 uses the technique of hashing to filter out unnecessary itemsets for next candidate itemset generation. When the support of candidate  $k$ -itemsets is counted by scanning the database, DHP accumulates information about candidate  $(k+1)$ -itemsets in advance in such a way that all possible  $(k+1)$ -itemsets of each transaction after some pruning are hashed to a hash table. Each bucket in the hash table consists of a number to represent how many itemsets have been hashed to this bucket thus far. We note that based on the resulting hash table, a bit vector can be constructed, where the value of one bit is set to be one if the number in the corresponding entry of the hash table is greater than or equal to  $s$ . As can be seen later, such a bit vector can be used to greatly reduce

the number of itemsets in  $C_i$ . This implementation detail is omitted in Fig. 3.

Fig. 3 gives the algorithmic form of DHP which, for ease of presentation, is divided into three parts. Part 1 gets a set of large 1-itemsets and makes a hash table (i.e.,  $H_2$ ) for 2-itemsets. Part 2 generates the set of candidate itemsets  $C_k$  based on the hash table (i.e.,  $H_k$ ) generated in the previous pass, determines the set of large  $k$ -itemsets  $L_k$ , reduces the size of database for the next large  $k$ -itemsets (as will be explained in Section 3.2 later), and makes a hash table for candidate large  $(k+1)$ -itemsets. (Note that in Part 1 and Part 2, the step of building a hash table to be used by the next pass is a unique feature of DHP.) Part 3 is basically same as Part 2 except that it does not employ a hash table. Note that DHP is particularly powerful to determine large itemsets in early stages, thus improving the performance bottleneck. The size of  $C_k$  decreases significantly in later stages, thus rendering little justification its further filtering. This is the very reason that we shall use Part 2 for early iterations, and use Part 3 for later iterations when the number of hash buckets with a count larger than or equal to  $s$  (i.e.,  $|\{x | H_k[x] \geq s\}|$  in Part 2 of Fig. 3) is less than a predefined threshold  $LARGE$ . We note that in Part 3 Procedure `apriori_gen` to generate  $C_{k+1}$  from  $L_k$  is essentially the same as the method used by algorithm Apriori in [5] in determining candidate itemsets, and we hence omit the details on it. Part 3 is included into DHP only for the completeness of our method.

After the setting by Part 1, Part 2 consists of two phases. The first phase is to generate a set of candidate  $k$ -itemsets  $C_k$  based on the hash table  $H_k$ , which is described by Procedure `gen_candidate` in Fig. 4. Same as Apriori, DHP also generates a  $k$ -itemset by  $L_{k-1}$ . However, DHP is unique in that it employs the bit vector, which is built in the previous pass, to test the validity of each  $k$ -itemset. Instead of including all  $k$ -itemsets from  $L_{k-1} * L_{k-1}$  into  $C_k$ , DHP adds a

```

/* Part 1 */
s = a minimum support;
set all the buckets of  $H_2$  to zero; /* hash table */
forall transaction  $t \in D$  do begin
    insert and count 1-items occurrences in a hash tree;
    forall 2-subsets  $x$  of  $t$  do
         $H_2[h_2(x)] ++$ ;
end
 $L_1 = \{c | c.count \geq s, c \text{ exists in the leaf node of the hash tree}\}$ 

/* Part 2 */
k = 2;
 $D_k = D$ ; /* database for large k-itemsets */
while ( $|\{x | H_k[x] \geq s\}| \geq LARGE$ ) { /* make a hash table */
    gen_candidate( $L_{k-1}, H_k, C_k$ );
    set all the buckets of  $H_{k+1}$  to zero;
     $D_{k+1} = \phi$ ;
    forall transactions  $t \in D_k$  do begin
        count_support( $t, C_k, k, \hat{t}$ ); /*  $\hat{t} \subseteq t$  */
        if ( $|\hat{t}| > k$ ) then do begin
            make_hasht( $\hat{t}, H_k, k, H_{k+1}, \tilde{t}$ );
            if ( $|\tilde{t}| > k$ ) then  $D_{k+1} = D_{k+1} \cup \{\tilde{t}\}$ ;
        end
    end
     $L_k = \{c \in C_k | c.count \geq s\}$ ;
    k ++;
}

/* Part 3 */
gen_candidate( $L_{k-1}, H_k, C_k$ );
while ( $|C_k| > 0$ ) {
     $D_{k+1} = \phi$ ;
    forall transactions  $t \in D_k$  do begin
        count_support( $t, C_k, k, \hat{t}$ ); /*  $\hat{t} \subseteq t$  */
        if ( $|\hat{t}| > k$ ) then  $D_{k+1} = D_{k+1} \cup \{\hat{t}\}$ ;
    end
     $L_k = \{c \in C_k | c.count \geq s\}$ ;
    if ( $|D_{k+1}| = 0$ ) then break;
     $C_{k+1} = \text{apriori\_gen}(L_k)$ ; /* refer to [?] */
    k ++;
}

```

Fig. 3. Main program of algorithm DHP.

$k$ -itemset into  $C_k$  only if that  $k$ -itemset passes the hash filtering, i.e., that  $k$ -itemset is hashed into a hash entry whose value is larger than or equal to  $s$ . As can be seen later, such hash filtering can drastically reduce the size of  $C_k$ . Every  $k$ -itemset that passes the hash filtering is included into  $C_k$  and stored into a hash tree [5], [11]. The hash tree built by  $C_k$  is then probed by each transaction later (i.e., in Part 2) when the database is scanned and the support of each candidate itemset is counted. The second phase of Part 2 is to count the support of candidate itemsets and to reduce the size of each transaction, as described by Procedure count\_support in Fig. 4. Same as in [5], a subset function is used to determine all the candidate itemsets contained in

each transaction. As transactions in the database (which is reduced after  $D_2$ ) are scanned one by one,  $k$ -subset of each transaction are obtained and used to count the support of itemsets in  $C_k$ . The methods for trimming a transaction and reducing the number of transactions are described in detail in Section 3.2.

An example of generating candidate itemsets by DHP is given in Fig. 5. For the candidate set of large 1-itemsets, i.e.,  $C_1 = \{\{A\}, \{B\}, \{C\}, \{D\}, \{E\}\}$ , all transactions of the database are scanned to count the support of these 1-items. In this step, a hash tree for  $C_1$  is built on the fly for the purpose of efficient counting. DHP tests whether or not each item exists already in the hash tree. If yes, it increases the count

```

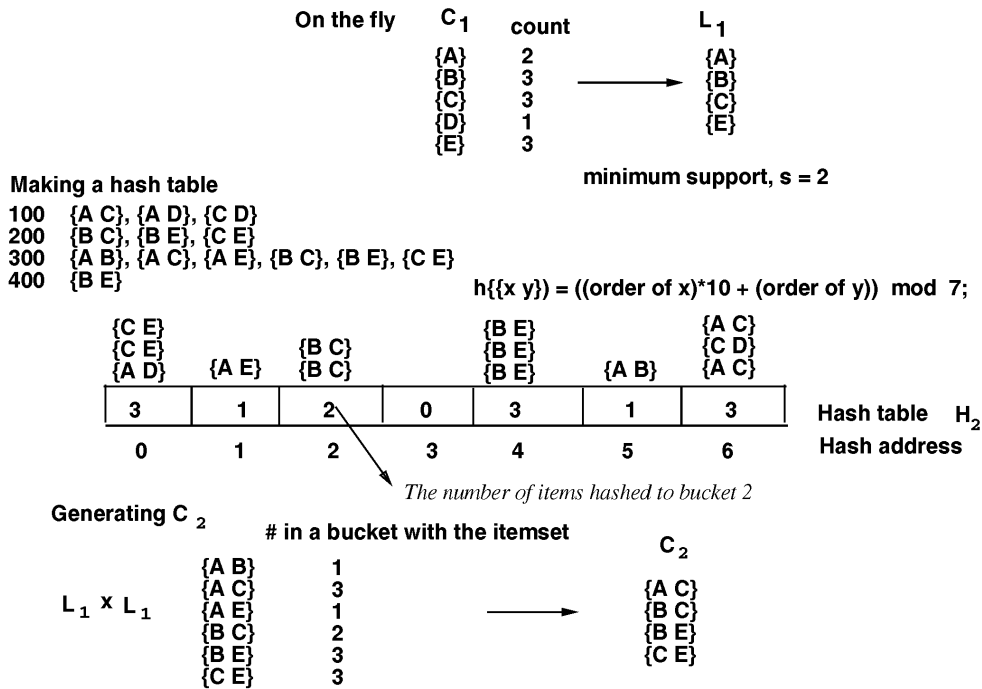
Procedure gen_candidate( $L_{k-1}, H_k, C_k$ )
   $C_k = \phi$ ;
  forall  $c = c_p[1] \cdot \dots \cdot c_p[k-2] \cdot c_p[k-1] \cdot c_q[k-1], c_p, c_q \in L_{k-1}, |c_p \cap c_q| = k-2$  do
    if ( $H_k[h_k(c)] \geq s$ ) then
       $C_k = C_k \cup \{c\}$ ;          /* insert  $c$  into a hash tree */
  end Procedure

Procedure count_support( $t, C_k, k, \hat{t}$ )          /* explained in Section ?? */
  forall  $c$  such that  $c \in C_k$  and  $c (= t_{i_1} \dots t_{i_k}) \in t$  do begin
     $c.count++$ ;
    for ( $j = 1; j \leq k; j++$ )  $a[i_j]++$ ;
  end
  for ( $i = 0, j = 0; i < |\hat{t}|; i++$ )
    if ( $a[i] \geq k$ ) then do begin  $\hat{t}_j = t_{i_i}; j++$ ; end
  end Procedure

Procedure make_hasht( $\hat{t}, H_k, k, H_{k+1}, \hat{t}$ )
  forall  $(k+1)$ -subsets  $x (= \hat{t}_{i_1} \dots \hat{t}_{i_{k+1}})$  of  $\hat{t}$  do
    if (for all  $k$ -subsets  $y$  of  $x, H_k[h_k(y)] \geq s$ ) then do begin
       $H_{k+1}[h_{k+1}(x)]++$ ;
      for ( $j = 1; j \leq k+1; j++$ )  $a[i_j]++$ ;
    end
  for ( $i = 0, j = 0; i < |\hat{t}|; i++$ )
    if ( $a[i] > 0$ ) then do begin  $\hat{t}_j = \hat{t}_{i_i}; j++$ ; end
  end Procedure

```

Fig. 4. Subprocedures for algorithm DHP.

Fig. 5. Example of a hash table and generation of  $C_2$ .

of this item by one. Otherwise, it inserts the item with a count equal to one into the hash tree. For each transaction, after occurrences of all the 1-subsets are counted, all the 2-subsets of this transaction are generated and hashed into

a hash table  $H_2$  in such a way that when a 2-subset is hashed to bucket  $i$ , the value of bucket  $i$  is increased by one. Fig. 5 shows a hash table  $H_2$  for a given database. After the database is scanned, each bucket of the hash table has the

number of 2-itemsets hashed to the bucket. Given the hash table in Fig. 5 and a minimum support equal to 2, we obtain a resulting bit vector  $\langle 1, 0, 1, 0, 1, 0, 1 \rangle$ . Using this bit vector to filter out 2-itemsets from  $L_1 * L_1$ , we have  $C_2 = \{ \{AC\}, \{BC\}, \{BE\}, \{CE\} \}$ , instead of  $C_2 = \{ \{AB\}, \{AC\}, \{AE\}, \{BC\}, \{BE\}, \{CE\} \}$  resulted by Apriori as shown in Fig. 2.

### 3.2 Reducing the Size of the Transaction Database

DHP reduces the database size progressively by not only trimming each individual transaction size but also pruning the number of transactions in the database. Note that, as observed in [5] on mining association rules, any subset of a large itemset must be a large itemset by itself. That is,  $\{B, C, D\} \in L_3$  implies  $\{B, C\} \in L_2$ ,  $\{B, D\} \in L_2$ , and  $\{C, D\} \in L_2$ . This fact suggests that a transaction be used to determine the set of large  $(k + 1)$ -itemsets only if it consists of  $(k + 1)$  large  $k$ -itemsets in the previous pass. In view of this, when  $k$ -subsets of each transaction are counted toward candidate  $k$ -itemsets, we will be able to know if this transaction meets the necessary condition of containing large  $(k + 1)$ -itemsets. This, in turn, means that if we have the number of candidate itemsets close to that of large itemsets when counting  $k$ -subsets, we can efficiently trim transactions and reduce the number of transactions by eliminating items which are found useless for later large itemset generation.

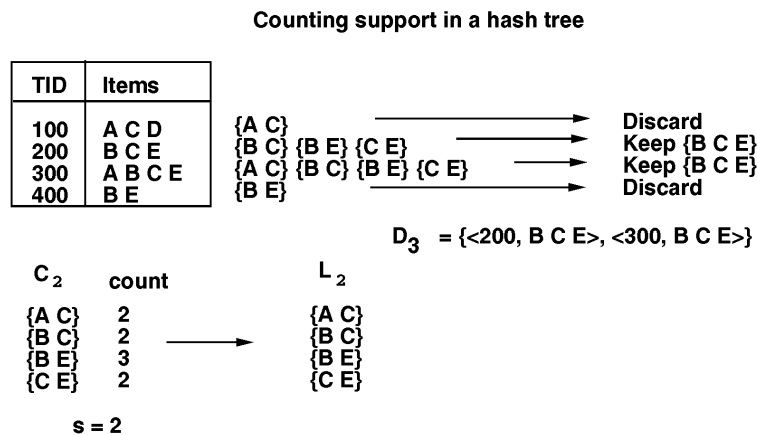
We now take a closer look at how the transaction size is trimmed by DHP. If a transaction contains some large  $(k + 1)$ -itemsets, any item contained in these  $(k + 1)$ -itemsets will appear in at least  $k$  of the candidate  $k$ -itemsets in  $C_k$ . As a result, an item in transaction  $t$  can be trimmed if it does not appear in at least  $k$  of the candidate  $k$ -itemsets in  $t$ . This concept is used in Procedure `count_support` to trim the transaction size. Certainly, the above is only a necessary condition, not a sufficient condition, for an item to appear in a candidate  $(k + 1)$ -itemset. In Procedure `make_hasht`, we further check that each item in a transaction is indeed covered by a  $(k + 1)$ -itemset (of the transaction) with all  $(k + 1)$  of its  $k$ -itemsets contained in  $C_k$ .

An example to trim and reduce transactions is given in Fig. 6. Note that the support of a  $k$ -itemset is increased as long as it is a subset of transaction  $t$  and also a member of  $C_k$ . As described by Procedure `count_support`,  $a[i]$  is used to

keep the occurrence frequency of the  $i$ th item of transaction  $t$ . When a  $k$ -subset containing the  $i$ th item is a member of  $C_k$ , we increase  $a[i]$  by 1, according to the index of each item in the  $k$ -subset (e.g., in transaction 100,  $a[0]$  corresponds to A,  $a[1]$  corresponds to C, and  $a[2]$  corresponds to D). Then, in Procedure `make_hasht`, before hashing of a  $(k + 1)$ -subset of transaction  $\hat{t}$ , we test all the  $k$ -subsets of  $\hat{t}$  by checking the values of the corresponding buckets on the hash table  $H_k$ . To reduce the transaction size, we then check each item  $\hat{t}_i$  in  $\hat{t}$  to see if item  $\hat{t}_i$  is indeed included in one of the  $(k + 1)$ -itemsets eligible for hashing from  $\hat{t}$ . Item  $\hat{t}_i$  is discarded if it does not meet such a requirement.

For example, in Fig. 6, transaction 100 has only a single candidate itemset AC. Then, occurrence frequencies of all the items are:  $a[0] = 1$ ,  $a[1] = 1$ , and  $a[2] = 0$ . Since all the values of  $a[i]$  are less than 2, this transaction is deemed not useful for generating large 3-itemsets and thus discarded. On the other hand, transaction 300 in Fig. 6 has four candidate 2-items and the occurrence frequencies of items are  $a[0] = 1$  (corresponding to A),  $a[1] = 2$  (corresponding to B),  $a[2] = 3$  (corresponding to C), and  $a[3] = 2$  (corresponding to E). Thus, it keeps three items—B, C, E—and discards item A.

For another example, if transaction  $t = ABCDEF$  and five 2-subsets,  $(AC, AE, AF, CD, EF)$ , exist in a hash tree built by  $C_2$ , we get values of array  $a[i]$  as  $a[0] = 3$ ,  $a[2] = 2$ ,  $a[3] = 1$ ,  $a[4] = 2$ , and  $a[5] = 2$ , according to the occurrences of each item. For large 3-itemsets, four items, A, C, E, and F, have a count larger than 2. Thus, we keep the items A, C, E, and F as transaction  $\hat{t}$  in Procedure `count_support` and discard items B and D, since they are useless in later passes. Clearly, not all items in  $\hat{t}$  contribute to the later large itemset generations. C, in fact, does not belong to any large 3-itemset since only AC and CD, but not AD, are large 2-itemsets. It can be seen from Procedure `make_hasht` that spurious items like C are removed from  $\hat{t}$  in the reduced database  $D_3$  for next large itemsets. Consequently, during the transaction scan, many transactions are either trimmed or removed, and only transactions which consist of necessary items for later large itemset generation are kept in  $D_{k+1}$ , thus progressively reducing the transaction database size. The fact that  $D_k$  de-



creases significantly along the pass number  $k$  is the very reason that DHP achieves a shorter execution time than Apriori, even in later iterations when the same procedure for large itemset generation is used by both algorithms. Fig. 6 shows an example of  $L_2$  and  $D_3$ .

## 4 EXPERIMENTAL RESULTS

To assess the performance of DHP, we conducted several experiments on large itemset generations by using an RS/6000 workstation with model 560. As will be shown later, the techniques of using a hash table and progressively reducing the database size enable DHP to generate large itemsets efficiently. The methods used to generate synthetic data are described in Section 4.1. The effect of the hash table size used by DHP is discussed in Section 4.2. Performance of DHP is assessed in Section 4.3. The scan-reduction technique of DHP is presented in Section 4.4.

### 4.1 Generation of Synthetic Data

The method used by this study to generate synthetic transactions is similar to the one used in [5] with some modifications noted below. Table 1 summarizes the meaning of various parameters used in our experiments. Each transaction consists of a series of potentially large itemsets, where those itemsets are chosen from a set of such itemsets  $L$ .  $|L|$  is set to 2,000. The size of each potentially large itemset in  $L$  is determined from a Poisson distribution with mean equal to  $|I|$ . Itemsets in  $L$  are generated as follows. Items in the first itemset are chosen randomly from  $N$  items. In order to have common items in the subsequent  $S_q$  itemsets in each of these  $S_q$  itemsets, some fraction of items are chosen from the first itemset generated, and the other items are picked randomly. Such a fraction, called the correlation level, is chosen from an exponential distribution with mean equal to 0.25 except in Section 4.4. So,  $S_q$  means the number of subsequent itemsets that are correlated to the first itemset, and this value was set to 10 in this study. After the first  $(S_q + 1)$  itemsets are generated, the generation process resumes a new cycle. That is, items in the next itemset are chosen randomly, and the subsequent  $S_q$  itemsets are so determined as to be correlated to that itemset in the way described above. The generation process repeats until  $|L|$  itemsets are generated. It is noted that, in [5], the generation of one itemset is only dependent on the previous itemset. Clearly, a larger  $S_q$  incurs a larger degree of "similarity" among transactions generated. Here, a larger  $S_q$  is used so

as to have more large itemset generation scenarios to observe in later iterations. As such, we are also able to have  $|L|/(S_q + 1)$  groups of transactions to model grouping or clustering in the retailing environment.

Each potentially large itemset in  $L$  has a weight, which is the probability that this itemset will be picked. Each weight is chosen from some given distribution of  $|L|$  potentially large itemsets such as an exponential distribution [5], and then normalized in such a way that the sum of all the weights is equal to one. For these weights, we used a Zipf-like distribution [12] so that we can control the number of large  $k$ -itemsets by varying a parameter of this distribution. The Zipf-like distribution of  $|L|$  potentially large itemsets is generated as follows. The probability  $p_i$  that the potentially large itemset is taken for a particular transaction from  $L$  is  $p_i = c/i^{(1-\theta)}$ , where

$$c = 1 / \sum_{i=1}^{|L|} (1 / i^{(1-\theta)})$$

is a normalization constant. Setting the parameter  $\theta = 0$  corresponds to the pure Zipf distribution, which is highly skewed; whereas  $\theta = 1$  corresponds to the uniform distribution. In this experiment, the value of  $\theta$  was set to between 0.60 and 0.75. Same as [5], we also use a corruption level during the transaction generation to model the phenomenon that all the items in a large itemset are not always bought together. Therefore, we can control easily the number of large  $k$ -itemsets by  $S_q$ ,  $\theta$ , and the correlation level. Each transaction is stored in a file system with the form of <transaction identifier, the number of items, items>.

### 4.2 Effect of the Size of a Hash Table

Note that the hash table size used by DHP affects the cardinality of  $C_2$  generated. In fact, from the process of trimming the database size described in Section 3.2, it can be seen that the size of  $C_2$  subsequently affects the determination of  $D_3$ . Table 2 shows the results from varying the hash table size, where  $|D_1| = 100,000$ ,  $|T| = 10$ ,  $|I| = 4$ ,  $N = 1,000$ ,  $|L| = 2,000$ ,  $\theta = 0.65$ , and  $s = 0.75$  percent. We use Tx.Iy.Dz to mean that  $x = |T|$ ,  $y = |I|$ , and the number of transactions in  $D_1$  is  $z \times 1,000$ . For notational simplicity, we use  $\{H_2 \geq s\}$  in Table 2 to represent the set of buckets which are hashed into by 750 (i.e.,  $100 \times 1,000 \times 0.75$  percent) or more 2-itemsets during the execution of Part 1 in DHP, and  $|\{H_2 \geq s\}|$  to denote the cardinality of  $\{H_2 \geq s\}$ . Also,  $\alpha$  represents the ratio of the number of 2-itemsets hashed into  $\{H_2 \geq s\}$  to the total number of 2-itemsets generated

TABLE 1  
MEANING OF VARIOUS PARAMETERS

$D_k$	Set of transactions for large $k$ -itemsets
$ D_k $	The number of transactions in $D_k$
$H_k$	Hash table containing $ H_k $ buckets for $C_k$
$C_k$	Set of candidate $k$ -itemsets
$L_k$	Set of large $k$ -itemsets
$ T $	Average size of the transactions
$ I $	Average size of the maximal potentially large itemsets
$ L $	Number of maximal potentially large itemsets
$N$	Number of items

TABLE 2  
RESULTS FROM VARYING HASH TABLE SIZES (T10.I4.D100)

$ H_2 $	524,288	262,144	131,072	65,536	32,768
$L_1$	595	595	595	595	595
$ \{H_2 \geq s\} $	119	123	128	148	210
$C_2$	150	220	313	543	1,469
$L_2$	115	115	115	115	115
$\alpha$	0.0273	0.0282	0.0296	0.0336	0.0461
size of $D_3$	326KB	332KB	333KB	351KB	408KB
$ D_3 $	12,450	12,631	12,663	13,393	15,595
total time	6.65	6.53	6.39	6.51	7.47



from the original database  $D_1$ .  $\alpha$  is a factor to represent what fraction of 2-itemsets are involved in candidate 2-itemsets. The total number of 2-itemsets of the experiment at Table 2 is 5,002,249, which is slightly larger than

$$|D_1| \cdot \binom{|T|}{2}.$$

For  $N = 1,000$ , the number of distinct 2-itemsets is  $\binom{N}{2}$ .

Note that, when  $n = \binom{N}{2}$  and  $|H_2|$  is chosen to be the exponent of 2 (which is greater than  $n$ ), we have  $|C_2| / |L_2| = 1.30$  in Table 2. In addition, note that in the second column of Table 2, the database to be used for generating large 3-itemsets, i.e.,  $D_3$ , is very small compared to the original database, indicating a very effective trimming in the database size by DHP. Specifically, the ratio of  $D_3$  to  $D_1$  is 6.91 percent in terms of their sizes, and 12.45 percent in terms of their numbers of transactions. Also, the average number of items in each transaction of  $D_3$  is 4.69 (instead of 10 in the original database). Clearly, a small number of items in each transaction incurs fewer comparisons in a hash tree and leads to a shorter execution time.

In Table 2, the values of  $|H_2|$  in DHP varies from  $\binom{N}{2}$  to  $\binom{N}{2}/16$ . When  $|H_2|$  is approximately  $\binom{N}{2}$  (as in the second column),  $|C_2| / |L_2| = 1.30$ , meaning that a larger  $H_2$  leads to a smaller  $C_2$  at the cost of using more memory. As  $|H_2|$  decreases,  $|C_2|$  and the execution time for  $L_2$  increase. The size of  $D_3$  then increases as well. We found that we can have fairly good overall performance until  $|H_2|$  is a quarter of  $\binom{N}{2}$  (i.e., until the fourth column). However, it is noted that even when the hash table has only  $\binom{N}{2}/16$  buckets, the number of candidate 2-itemsets is still significantly smaller than  $|L_1|(|L_1| - 1)/2$ . Clearly, when either the minimum support is small or the number of total 2-itemsets is large, it is advantageous to use a large  $|H_2|$  for DHP. The reduction ratios of  $|C_2|$  by DHP with respect to  $|L_1|(|L_1| - 1)/2$  for various sizes of  $H_2$  are shown in Fig. 7, where a logarithmic scale is used in y-axis for ease of presentation.

### 4.3 Performance of DHP

#### 4.3.1 Comparison of DHP and Apriori

Table 3 shows the relative performance between Apriori used in [5] and DHP. Here, we use  $|T| = 15$ ; i.e., each transaction has 15 items in average, so as to have more large itemsets in later passes for interest of presentation. The execution times of these two algorithms are shown in Fig. 8. In DHP,  $|H_2|$  is chosen to be the exponent of 2, which is greater than  $\binom{N}{2}$  (i.e., with 524,288 buckets) and parameter  $\theta$  is 0.75 in this subsection. In this experiment,

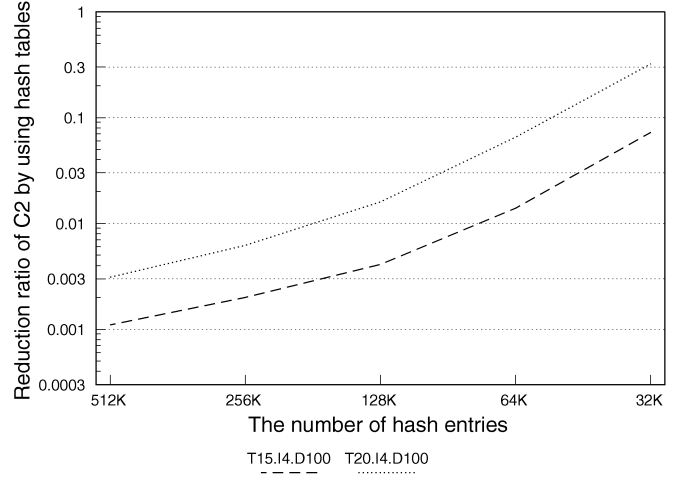


Fig. 7. Reductions ratio of  $|C_2|$  by DHP for different sizes of  $H_2$ .

TABLE 3  
COMPARISON OF EXECUTION TIME (T15.I4.D100)

	Apriori number	DHP		
		number	$D_k$	$ D_k $
$L_1$	760	760	6.54MB	100,000
$C_2$	288,420	318	6.54MB	100,000
$L_2$	211	211		
$C_3$	220	220	0.51MB	20,047
$L_3$	204	204		
$C_4$	229	229	0.25MB	8,343
$L_4$	227	227		
$C_5$	180	180	0.16MB	4,919
$L_5$	180	180		
$C_6$	94	94	0.10MB	2,459
$L_6$	94	94		
$C_7$	29	29	0.06MB	1,254
$L_7$	29	29		
$C_8$	4	4	0.05MB	1,085
$L_8$	4	4		
total time	43.36	13.57		

DHP uses a hash table for the generation of  $C_2$  (i.e., Part 1 of Fig. 3). Starting from the third pass, DHP is the same as Apriori in that the same procedure for generating large itemsets (i.e., Part 3 of Fig. 3) is used by both algorithms, but different from the latter in that a smaller transaction database is scanned by DHP. The last column represents the database size in the  $k$ th pass (i.e.,  $D_k$ ) used by DHP and its cardinality (i.e.,  $|D_k|$ ). More explicitly, Apriori scans the full database  $D_1$  for every pass, whereas DHP only scans the full database for the first two passes and then scans the reduced database  $D_k$  thereafter. As mentioned before, in [5] the hybrid algorithm has the option of switching from Apriori to another algorithm AprioriTid after early passes for better performance, and such an option is not adopted

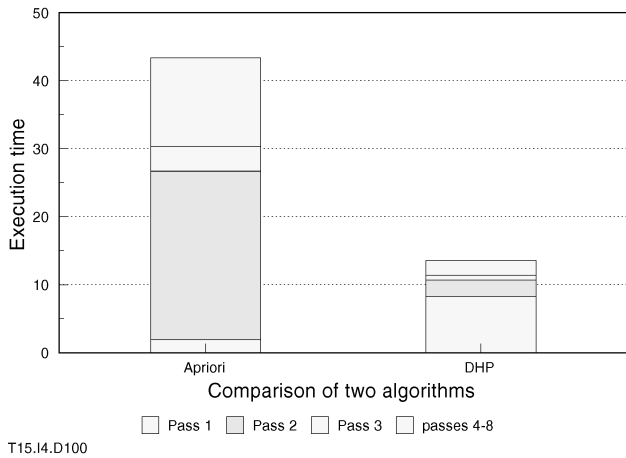


Fig. 8. Execution time of Apriori and DHP.

here. It can, nevertheless, be seen from Fig. 8 that the execution time of the first two passes by Apriori is larger than the total execution time by DHP.<sup>5</sup>

It can be seen from Table 3 that the execution time of the first pass of DHP is slightly larger than that of Apriori due to the extra overhead required for generating  $H_2$ . However, DHP incurs significantly smaller execution times than Apriori in later passes, not only in the second pass when a hash table is used by DHP to facilitate the generation of  $C_2$ , but also in later passes when the same procedure is used, showing the advantage of scanning smaller databases by DHP. Here, the execution time of the first two passes by Apriori is about 62 percent of the total execution time. This is the very motivation of employing DHP for early passes to achieve performance improvement.

Fig. 9 shows the execution time ratios of DHP to Apriori over various minimum supports, ranging from 0.75 percent to 1.25 percent. Fig. 9 indicates that DHP constantly performs well for various minimum supports. Fig. 10 shows the effect of progressively reducing the transaction database by DHP. As pointed out earlier, this very feature of DHP is made feasible in practice due to the early reduction on the size of  $C_2$ , and turns out to be very powerful to facilitate the later itemset generations. Note that DHP is not only reducing the number of transactions but also trimming the items in each transaction. It can be seen that the average number of items is, on one hand, reduced by the latter process (i.e., trimming each transaction size), but on the other hand, is increased by the former process (i.e., reducing the number of transactions) since transactions eliminated are usually small ones. As a result of these two conflicting factors, as shown in Case A of Fig. 10, whose y-axis uses a logarithmic scale, the average number of items in each transaction in  $D_i$  remains approximately the same along the pass number  $i$ . For example, for T20.I4.D100, starting from 20 items, the average number of items in a transaction drops to 6.2, and then increases slightly since several small transactions are eliminated in later passes. To explicitly show the effect of trimming each transaction size,

5. The benefit of AprioriTid in later passes is complementary to the focus of DHP on initial passes. In fact, in Part 3 of DHP, AprioriTid can be used instead of Apriori if desired.

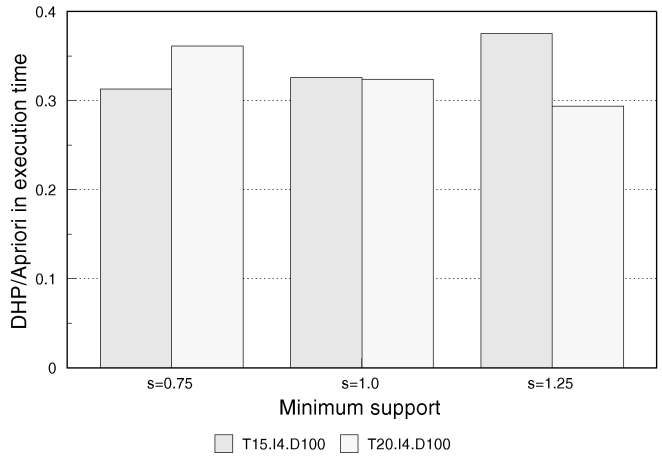


Fig. 9. Execution time comparison between DHP and Apriori for some minimal supports.

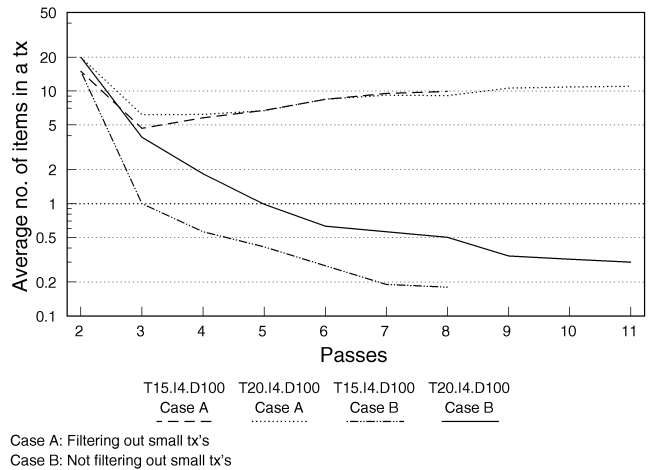


Fig. 10. The average number of items in a transaction in each pass.

we conducted another experiment where transactions are only trimmed but not thrown away along the process. The average number of items in each transaction resulting from this experiment is shown by Case B of Fig. 10, which indicates that the trimming method employed by DHP is very effective.

#### 4.3.2 Scale-Up Experiment for DHP

Fig. 11 shows that the execution time of DHP increases linearly as the database size increases, meaning that DHP possesses the same important feature as Apriori. Also, we examine the performance of DHP as the number of items,  $N$ , increases. Table 4 shows the execution times of DHP when the number of items increases from 1,000 to 10,000 for three data sets: T5.I2.D100, T10.I4.D100, and T20.I6.D100. In the experiments for Table 4, the minimum support is 0.75 percent; the hash table size is the exponent of 2, which is greater than  $\binom{1,000}{2}$ ; and parameter  $\theta$  of the Zipf-like distribution is 0.65. Note that the portion of time on determining  $L_1$  for the case of small transactions (e.g., T5 in Table 4) is relatively larger than that for the case of large transactions (e.g., T20 in Table 4). In other words,

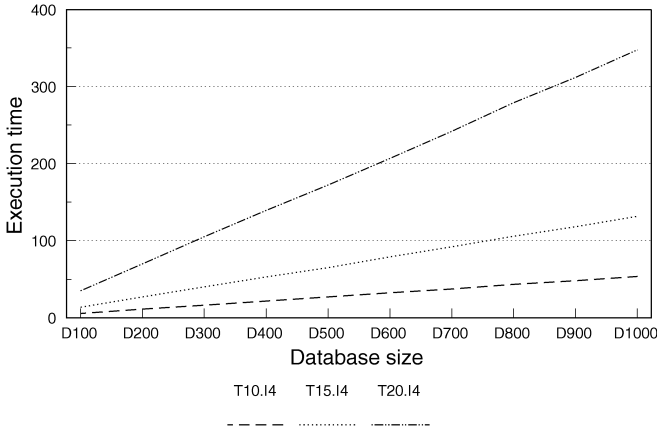


Fig. 11. Performance of DHP when the database size increases.

TABLE 4  
PERFORMANCE OF DHP WHEN  
THE NUMBER OF ITEMS INCREASES

N	T5.I2	T10.I4	T20.I6
1,000	2.26	6.65	45.25
2,500	2.33	6.25	17.86
5,000	2.53	6.04	18.65
7,500	2.58	6.76	19.17
10,000	2.61	6.96	20.10

a large transaction has a larger likelihood of having large itemsets to process than a small transaction. Also, given a fixed minimum support, when the number of items  $N$  increases, the execution time to obtain  $L_1$  increases since the size of  $L_1$  is usually close to  $N$ , but the execution time to obtain larger  $k$ -itemsets decreases since the support for an itemset is averaged out by more items and thus decreases. Consequently, as shown in Table 4, when  $N$  increases, execution times for small transactions increase a little. For large transactions like T20 in Table 4, the execution times decrease a little since the number of large itemsets decreases as we increase the number of items.

#### 4.4 Scan-Reduction Method: Reducing the Amount of Disk I/O Required

Recall that DHP generates a small number of candidate 2-itemsets by using a hashing technique. In fact, this small  $C_2$  can be used to generate the candidate 3-itemsets. Clearly, a  $C'_3$  generated from  $C_2 * C_2$ , instead of from  $L_2 * L_2$ , will have a size greater than  $|C_3|$  where  $C_3$  is generated from  $L_2 * L_2$ . However, if  $|C'_3|$  is not much larger than  $|C_3|$ , and both  $C_2$  and  $C'_3$  can be stored in the main memory, we can find  $L_2$  and  $L_3$  together when the next scan of the database is performed, thereby saving one round of database scan. It can be seen that using this concept, one can determine all  $L_k$ s by as few as two scans of the database (i.e., one initial scan to determine  $L_1$  and a final scan to determine all other large itemsets), assuming that  $C'_k$  for  $k \geq 3$  is generated from  $C'_{k-1}$

and all  $C'_k$ s for  $k > 2$  can be kept in the memory, where  $C'_2 = C_2$ . This technique is called *scan-reduction*.

Table 5 shows results of DHP for cases with and without scan-reduction, where  $|D_1| = 100,000$ ,  $|I| = 4$ , and the parameter  $\theta$  of a Zipf-like distribution is 0.75. It can be seen that the values of  $C_k$  with scan reduction are slightly larger than those without scan-reduction because the former is generated from  $C_{k-1} * C_{k-1}$ . The last column gives the respective execution time: The one in the row of  $L_k$  is CPU time, and that in the row of  $D_k$  is I/O time. To obtain disk I/O time, the disk speed is assumed to be 3 MB/sec and a 1-MB buffer is employed. In the case without scan-reduction, DHP scans the database 11 times to find all the large itemsets, whereas DHP with scan-reduction only involves two database scans. Note that after initial scans, disk I/O involved by both cases will include both disk read and disk write (i.e., writing the trimmed version of the database back to the disk). The CPU and disk I/O times for the case without scan-reduction are 34.91 sec and 24.42 sec, respectively, whereas those for the case with scan reduction are 34.77 sec and 16.82 sec, respectively. Considering both CPU and I/O times, we see a prominent advantage of scan-reduction.

Note that, when the minimum support is relatively small or the number of potentially large itemsets is large,  $C_k$  and  $L_k$  could become large. If  $|C'_{k+1}| > |C'_k|$  for  $k \geq 2$ , then it may cost too much CPU time to generate all subsequent  $C'_j$ ,  $j > k + 1$ , from candidate sets of large itemsets since the size of  $C_j$  may become huge quickly, thus compromising all the benefit from saving disk I/O cost. For the illustrative example in Fig. 2 if  $C_3$  was determined from  $C_2 * C_2$ , instead of from  $L_2 * L_2$ , then  $C_3$  would be  $\{\{ABC\}, \{ABE\}, \{ACE\}, \{BCE\}\}$ . This fact suggests that a timely database scan to determine large itemsets will in fact pay off. After a database scan, one can obtain the sets of large itemsets which are not determined thus far (say, up to  $L_m$ ) and then construct the set of candidate  $(m + 1)$ -itemsets,  $C_{m+1}$ , based on  $L_m$  from that point. According to our experiments, we found that if  $|C'_{k+1}| > |C'_k|$  for some  $k \geq 2$ , it is usually beneficial to have a database scan to obtain  $L_{k+1}$  before the set of candidate itemsets becomes too big. We then derive  $C'_{k+2}$  from  $L_{k+1}$ . After that, we again use  $C'_j$  to derive  $C'_{j+1}$  for  $j \geq k + 2$ . The process continues until the set of candidate  $(j + 1)$ -itemsets becomes empty. An illustrative example for such scenarios is shown in Table 6, where cases with and without scan-reduction are given for comparison purposes. For the case with scan-reduction in Table 6, after skipping the database scan in the second pass, a database scan occurs in the third pass because of the huge size of  $C_3$ . It is important to see that after the database scan in the third pass, the sizes of  $C_k$  are not only reduced significantly (e.g.,  $C_3$  is 26,270 and  $C_4$  is 3,489), but also stay pretty close to those obtained through scanning database every time, showing the very advantage of scanning the database timely.

TABLE 5  
EFFECT OF SCAN-REDUCTION FOR DHP WHEN  $|T| = 20$

$k$	1	2	3	4	5	6	7	8	9	10	11	time (sec)
DHP without scan-reduction												
$C_k$	–	1,150	627	596	618	531	348	167	55	11	1	–
$L_k$	856	675	542	586	618	531	348	167	55	11	1	34.91
$D_k$ (MB)	25.47	25.47	5.61	2.52	1.17	0.63	0.54	0.45	0.21	0.18	0.15	24.42
DHP with scan reduction												
$C_k$	–	1,150	828	761	712	572	361	169	55	11	1	–
$L_k$	856	675	542	586	618	531	348	167	55	11	1	34.77
$D_k$ (MB)	25.47	–	–	–	–	–	–	–	–	–	25.47	16.82

TABLE 6  
THE EFFECT OF SCAN-REDUCTION FOR DHP WHEN  $|T| = 30$

$k$	1	2	3	4	5	6	7	8	9	10	11
DHP without scan-reduction											
$C_k$	–	7,391	8,638	3,489	2,916	2,006	1,085	442	127	23	2
$L_k$	920	3,644	3,375	3,450	2,914	2,006	1,085	442	127	23	2
$D_k$ (MB)	37.23	37.23	25.12	15.48	10.59	7.63	4.79	3.27	2.65	1.37	0.61
DHP with scan-reduction											
$C_k$	–	7,391	26,270	3,489	2,941	2,012	1,086	442	127	23	2
$L_k$	920	3,644	3,375	3,450	2,914	2,006	1,085	442	127	23	2
$D_k$ (MB)	37.23	–	37.23	–	–	–	–	–	–	–	20.95

## 5 CONCLUSIONS

We examined in this paper the issue of mining association rules among items in a large database of sales transactions. The problem of discovering large itemsets was solved by constructing a candidate set of itemsets first and then, identifying, within this candidate set, those itemsets that meet the large itemset requirement. We proposed an effective algorithm DHP for the initial candidate set generation. DHP is a hash-based algorithm and is especially effective for the generation of candidate set for large 2-itemsets, where the number of candidate 2-itemsets generated is, in orders of magnitude, smaller than that by previous methods, thus resolving the performance bottleneck. In addition, the generation of smaller candidate sets enables us to effectively trim the transaction database at a much earlier stage of the iterations, thereby reducing the computational cost for later stages significantly. A scan-reduction technique was also described. Extensive simulation study has been conducted to evaluate performance of the proposed algorithm.

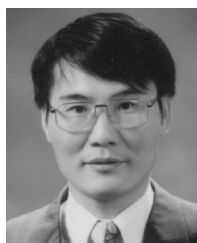
## ACKNOWLEDGMENTS

Jong Soo Park is partially supported by the Korea Science & Engineering Foundation. Ming-Syan Chen is supported, in part, by National Science Council Project No. 87-2213-E-002-009 and 87-2213-E-002-101, Taiwan, Republic of China.

## REFERENCES

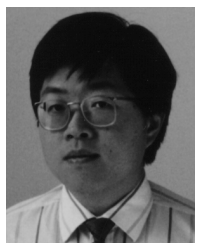
- [1] R. Agrawal, C. Faloutsos, and A. Swami, "Efficient Similarity Search in Sequence Databases," *Proc. Fourth Int'l Conf. Foundations of Data Organization and Algorithms*, Oct. 1993.
- [2] R. Agrawal, S. Ghosh, T. Imielinski, B. Iyer, and A. Swami, "An Interval Classifier for Database Mining Applications," *Proc. 18th Int'l Conf. Very Large Data Bases*, pp. 560–573, Aug. 1992.
- [3] R. Agrawal, T. Imielinski, and A. Swami, "Mining Association Rules between Sets of Items in Large Databases," *Proc. ACM SIGMOD*, pp. 207–216, May 1993.
- [4] R. Agrawal, H. Mannila, R. Srikant, H. Toivonen, and A. Inkeri Verkamo, "Fast Discovery of Association Rules," *Advances in KDDM*, U. Fayyad et al., eds., MIT/AAAI Press, 1995.
- [5] R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules in Large Databases," *Proc. 20th Int'l Conf. Very Large Data Bases*, pp. 478–499, Sept. 1994.
- [6] R. Agrawal and R. Srikant, "Mining Sequential Patterns," *Proc. 11th Int'l Conf. Data Eng.*, pp. 3–14, Mar. 1995.
- [7] T.M. Anwar, H.W. Beck, and S.B. Navathe, "Knowledge Mining by Imprecise Querying: A Classification-Based Approach," *Proc. Eighth Int'l Conf. Data Eng.*, pp. 622–630, Feb. 1992.
- [8] J. Han, Y. Cai, and N. Cercone, "Knowledge Discovery in Databases: An Attribute-Oriented Approach," *Proc. 18th Int'l Conf. Very Large Data Bases*, pp. 547–559, Aug. 1992.
- [9] J. Han and Y. Fu, "Discovery of Multiple-Level Association Rules from Large Databases," *Proc. 21th Int'l Conf. Very Large Data Bases*, pp. 420–431, Sept. 1995.
- [10] M. Houtsma and A. Swami, "Set-Oriented Mining for Association Rules in Relational Databases," *Proc. 11th Int'l Conf. Data Eng.*, pp. 25–33, Mar. 1995.
- [11] E.G. Coffman Jr. and J. Eve, "File Structures Using Hashing Functions," *Comm. ACM*, vol. 13, no. 7, pp. 427–432 and 436, July 1970.
- [12] D. Knuth, *The Art of Computer Programming, Volume 3: Sorting and Searching*. Addison-Wesley, 1973.
- [13] H. Mannila, H. Toivonen, and A. Inkeri Verkamo, "Efficient Algorithms for Discovering Association Rules," *Proc. AAAI Workshop Knowledge Discovery in Databases*, pp. 181–192, July 1994.

- [14] R.T. Ng and J. Han, "Efficient and Effective Clustering Methods for Spatial Data Mining," *Proc. 18th Int'l Conf. Very Large Data Bases*, pp. 144–155, Sept. 1994.
- [15] G. Piatetsky-Shapiro, "Discovery, Analysis and Presentation of Strong Rules," *Knowledge Discovery in Databases*, pp. 229–248, 1991.
- [16] J.R. Quinlan, "Induction of Decision Trees," *Machine Learning*, vol. 1, no. 81–106, 1986.
- [17] J.T.-L. Wang, G.-W. Chirn, T.G. Marr, B. Shapiro, D. Shasha, and K. Zhang, "Combinatorial Pattern Discovery for Scientific Data: Some Preliminary Results," *Proc. ACM SIGMOD*, Minneapolis, pp. 115–125, May 1994.



**Jong Soo Park** received the BS degree in electrical engineering with honors from Pusan National University, Pusan, Korea, in 1981; and the MS and PhD degrees in electrical engineering from the Korea Advanced Institute of Science and Technology, Seoul, Korea, in 1983 and 1990, respectively. From 1983 to 1986, he served as an engineer at the Korean Ministry of National Defense. He was a visiting researcher at the IBM Thomas J. Watson Research Center, Yorktown Heights, New York, from July 1994 to

July 1995. He is currently an associate professor in the Department of Computer Science at the Sungshin Women's University in Seoul. His research interests include geographic information systems and data mining. He is a member of the ACM and the IEEE.



**Ming-Syan Chen** (S'87-M'88-SM'93) received the BS degree in electrical engineering from the National Taiwan University, Taipei, Taiwan, Republic of China, in 1982; and the MS and PhD degrees in computer information and control engineering from the University of Michigan, Ann Arbor, in 1985 and 1988, respectively. Dr. Chen is a professor in the Electrical Engineering Department of the National Taiwan University. His research interests include database systems, multimedia technologies, and internet software.

He had been a research staff member at the IBM Thomas J. Watson Research Center, Yorktown Heights, New York, from 1988 to 1996, primarily involved in projects related to parallel databases, multimedia systems, and data mining.

Dr. Chen is serving as an editor of *IEEE Transactions on Knowledge and Data Engineering* and is leading the distance learning activities of the National Taiwan University, Taiwan. He has published more than 70 refereed papers in international journals and the proceedings of international conferences in his research areas, and more than 30 of the papers have appeared in ACM and IEEE journals and transactions. In 1996, he served as co-guest editor of a special issue of *IEEE Transactions on Knowledge and Data Engineering* on data mining. He is an inventor of many international patents in the areas of interactive video payout, video server design, interconnection network, and concurrency and coherency control protocols. He has received numerous awards for his inventions and patent applications, including the Outstanding Innovation Award from IBM in 1994 for his contribution to parallel transaction design for a major database product. Dr. Chen is a senior member of the IEEE and a member of the Association for Computing Machinery.



**Philip S. Yu** (S'76-M'78-SM'87-F'93) received the BS degree in electrical engineering from the National Taiwan University, Taipei, Taiwan, Republic of China, in 1972; the MS and PhD degrees in electrical engineering from Stanford University in 1976 and 1978, respectively; and the MBA degree from New York University in 1982. Since 1978, he has been with the IBM Thomas J. Watson Research Center, Yorktown Heights, New York. Currently, he is manager of the Software Tool and Technique group there.

His current research interests include database systems, Internet applications, data mining, multimedia systems, transaction and query processing, parallel and distributed systems, disk arrays, computer architecture, performance modeling, and workload analysis. He has published more than 200 papers in refereed journals and conference proceedings, more than 140 research reports, and 90 invention disclosures. He holds, or has applied for, 50 U.S. patents.

Dr. Yu was an editor of *IEEE Transactions on Knowledge and Data Engineering*. In addition to serving as a program committee member for various conferences, he was program chair of the Second International Workshop on Research Issues on Data Engineering: Transaction and Query Processing, and program co-chair of the 11th International Conference on Data Engineering. He has received several IBM and external honors, including Best Paper Awards, IBM Outstanding Innovation Awards, Outstanding Technical Achievement Awards, Research Division Awards, and 15 Invention Achievement Awards. He is a fellow of the IEEE and the ACM.