

# Document Ranking on Weight-Partitioned Signature Files

DIK LUN LEE and LIMING REN  
The Ohio State University

---

A signature file organization, called the weight-partitioned signature file, for supporting document ranking is proposed. It employs multiple signature files, each of which corresponds to one term frequency, to represent terms with different term frequencies. Words with the same term frequency in a document are grouped together and hashed into the signature file corresponding to that term frequency. This eliminates the need to record the term frequency explicitly for each word. We investigate the effect of false drops on retrieval effectiveness if they are not eliminated in the search process. We have shown that false drops introduce insignificant degradation on precision and recall when the false-drop probability is below a certain threshold. This is an important result since false-drop elimination could become the bottleneck in systems using fast signature file search techniques. We perform an analytical study on the performance of the weight-partitioned signature file under different search strategies and configurations. An optimal formula is obtained to determine for a fixed total storage overhead the storage to be allocated to each partition in order to minimize the effect of false drops on document ranks. Experiments were performed using a document collection to support the analytical results.

*Categories and Subject Descriptors:* H.2.2 [Database Management]: Physical Design—access methods; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—retrieval models; H.3.6 [Information Storage and Retrieval]: Library Automation

*General Terms:* Design, Experimentation, Performance

*Additional Key Words and Phrases:* Access method, document retrieval, information retrieval, signature file, superimposed coding, text retrieval

---

## 1. INTRODUCTION

The signature file approach has been investigated quite extensively in the last decade. It has been applied to a large variety of applications, including textual databases such as news databases [Stanfill and Kahle 1986],

---

Part of this research was done while the first author was on sabbatical at the Online Computer Library Center (OCLC), Dublin, Ohio, from September, 1992, to May, 1993.

Authors' address: Department of Computer and Information Science, The Ohio State University, 2036 Neil Avenue, Columbus, OH 43210-1277; email: {dlee; ren}@cis.ohio-state.edu.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 1996 ACM 1046-8188/96/0400-0109 \$03.50

text block			... text ..... database ...
word signatures:			
	text	001 000 110 010	
	database	000 010 101 001	
block signature (v)		001 010 111 011	
Queries	Query Signatures	Results	
1) retrieval	010 001 000 011	← no match	
2) database	000 010 101 001	← match	
3) information	001 000 111 000	← false drop	

Fig. 1. Signature generation and comparison based on superimposed coding.

multiattribute retrieval in relational databases [Aho and Ullman 1979; Lee and Leng 1990], multimedia office filing [Christodoulakis et al. 1986], and chemical databases for DNA matching [Lipman and Pearson 1985].

Superimposed coding is the most common way for generating signatures from a text. In superimposed coding, a text is divided into text blocks containing the same number of unique, nontrivial words. Each word in a text block is hashed into a word signature. A block signature is generated by superimposing all word signatures generated from the block. In a query, the query terms are hashed and superimposed into a query signature in a similar way. Then the query signature is matched against each signature in the signature file. Figure 1 is an example showing the generation of the block signature from a text block.

The signature file is a filtering mechanism that will eliminate most, but not all, of the text blocks that will not match the query. The first case shown in Figure 1 illustrates this point. The query signature does not “match” with the text signature in that some of the bits in the text signature are zero while the corresponding bits in the query signature are set to one. If the query term “retrieval” is indeed in the text, the query signature would be one of the word signatures forming the text signature; thus, every bit set in the query signature will be set in the text signature. The second case shows a match, where for each bit in the query signature set to one the corresponding bit in the block signature is also set to one. The third case is a *false drop*. False drops are text blocks that the signature file identifies as containing the query terms (i.e., a match), but indeed they do not. They can be eliminated by further comparing the query terms with the text blocks.

Compared to inverted files, the signature file approach has two major advantages:

- (1) Its storage overhead can be controlled easily and, in general, is very low. Therefore, if only a limited space for indexes is available, the signature length can be set so that the whole signature file will fit into the given space.
- (2) Its structure is very simple. Thus, it has low processing overhead on insertion. Since signatures are conceptually organized as a sequential

file, it is easy to partition a signature file for parallel processing [Stanfill and Kahle 1986].

The signature file approach also has the following drawbacks:

- (1) Its search speed is low compared to inverted files, since the signature file has to be searched exhaustively if signatures are organized in a single sequential file.
- (2) A signature encodes the presence of terms in a text. Therefore, it can easily support the Boolean retrieval model. However, in order to support document ranking, term frequencies must be stored; this is difficult to do on signatures, which are simple bit vectors.
- (3) Signature files introduce false drops that are expensive to eliminate.

Recent research on signature files has addressed the improvement of search speed, as exemplified by the partitioned signature file approach [Lee and Leng 1989; 1990; Zezula et al. 1991]. This article addresses the last two drawbacks of signature files listed above. We first present a technique, called the weight-partitioned signature file, for implementing document ranking with signature files. The technique we propose does not represent term frequency values explicitly. Furthermore, it also avoids exhaustive search on the signature file. An analytical study is performed on various performance issues of the weight-partitioned signature file. We also present an extensive evaluation on the degradation of retrieval effectiveness for the case when false drops are not eliminated from the signature file.

Section 2 discusses the problem of false drops in signature files and the motivation of this study. Section 3 defines the term weights and the similarity formula under the vector-space model. The design of the signature file for supporting document ranking is outlined in Section 4. Section 5 presents an analytical study and the experimental results on the performance of the weight-partitioned signature file under different search strategies and storage configurations. Section 6 presents the conclusions and outlines some problems for future research.

### 1.1 Problems of False Drops

False drops are an inherent property of signature files. The *false-drop probability*, defined as the ratio of false drops to the number of unqualified signatures, is dependent on the size of the signatures and, thus, on the storage overhead. The false-drop probability can be reduced arbitrarily by using a longer signature length while keeping the number of distinct keywords in a text block the same. However, to eliminate false drops completely, it is necessary to retrieve the text and to perform pattern matching between the query term(s) and the text after matching the signature file. This process imposes a great penalty on the speed of signature files. The cost of false-drop elimination is particularly acute for large collections, since the number of false drops is proportional to the number of signatures and, thus, to the size of the collection. Furthermore, when queries are not highly selective, a large number of matches will be

returned, and each of the matches has to be verified with pattern matching. The high cost raises a serious concern that it will become the bottleneck in retrieval, especially in systems employing fast signature search techniques [Lee et al. 1994]. In other words, the search speed will be slow no matter how fast the signature file itself can be searched.

This article investigates the feasibility of retaining the false drops generated by the signature file. For Boolean retrieval, it is problematic, since users are accustomed to the semantics of Boolean retrieval where each document retrieved will satisfy the Boolean condition that the user specifies. Compared to the Boolean model, retrieval in statistical models is less predictable in that documents are retrieved based on their relative ranking rather than on their absolute scores. Therefore, given a query, it is difficult to tell why a document is retrieved without examining the other documents. This makes the effect of false drops less observable to the user, even if the document ranks are slightly perturbed by the false drops. This article is based on the conjecture that, since document ranking is not a precise process (i.e., only some of the retrieved documents are relevant to the users), the addition of some "noise" or perturbation, due to false drops, to the resulting document set should have a small impact on retrieval effectiveness. If the conjecture is true, then the signature file approach will become a very attractive technique for implementing document ranking.

## 2. TERM WEIGHTS AND SIMILARITY MEASURE

Term weights are determined by the  $tf \times idf$  strategy [Salton and Buckley 1988]. (Table I summarizes the important symbols in this article.) Precisely, the weight of term  $j$  in document  $i$  is defined as

$$w_{i,j} = tf_{i,j} \times idf_j,$$

where  $tf_{i,j}$  is the frequency of occurrence of term  $j$  in document  $i$ , and  $idf_j$  is the inverse document frequency of term  $j$  in the entire collection. The similarity between a document and a query is based on the cosine measure,

$$\Theta(D_i, Q) = \frac{\sum_{j=1}^V w_{Q,j} \times w_{i,j}}{\sqrt{\sum_{j=1}^V w_{Q,j}^2 \sum_{j=1}^V w_{i,j}^2}}, \quad (1)$$

where  $V$  is the vocabulary size, and  $w_{Q,j}$  is the weight of term  $j$  in the query, defined in the same way as document term weights.

It is clear from Eq. (1) that the factor  $\sum_{j=1}^V w_{Q,j}^2$  in the denominator is a constant for each document, so it only affects the absolute scores but not the ranks of the documents. Therefore, this factor is not used in our study. The factor  $\sum_{j=1}^V w_{i,j}^2$  is very expensive to compute, because the weight of every term in a document has to be computed, not just for the terms specified in the query. Note that this factor cannot be precomputed for each document because the weight depends on  $idf$ , which changes as documents are inserted into and deleted from a collection. Therefore, we approximate

Table I. Major Symbols

$C$	Total storage constraint
$D$	A document
$D_i$	$i$ th document in the database
$d$	Number of terms in document $D$
$d_i$	Number of terms in document $D_i$
$I_i$	Average number of terms in a document with term frequency $i$
$I_i(D)$	Number of terms in $D$ with term frequency $i$
$idf_j$	Inverse document frequency of term $j$ in the database
$l_i$	Total number of signatures in $SF_i$
$l_i(D)$	Number of signatures in $SF_i$ for document $D$
$m$	Signature length
$m_i$	Signature length in $SF_i$
$\bar{m}$	Abbreviation for $m_1, \dots, m_n$
$N$	Total number of documents
$n$	Number of $tf$ groups
$p(m, w, s)$	False-drop probability given $m$ , $w$ , and $s$
$p_i$	False-drop probability of a signature of file $tf_i$
$Q$	Query
$q_i$	False-drop probability of a signature against file $tf_i$
$s$	Number of words hashed into a signature
$s_i$	Number of words hashed into a signature in $SF_i$
$\bar{s}$	Abbreviation for $s_1, \dots, s_n$
$SF_i$	Signature file corresponding to term frequency $i$
$\Theta(D_i, Q)$	Similarity between document $D_i$ and query $Q$
$tf_{i,j}$	Term frequency of term $j$ in document $i$
$V$	Size of the document vocabulary
$w$	Weight of a word signature (number of bits set to 1)
$w_i$	Weight of a word signature in $tf_i$ (number of bits set to 1)
$\bar{w}$	Abbreviation for $w_1, \dots, w_n$
$w_{i,j}$	Weight of term $j$ in document $i$
$w_{Q,j}$	Weight of term $j$ in query $Q$
$W$	General term ranking weight

the normalization factor with the square root of the number of terms in a document (i.e., the term weights are set to one when the normalization factor is computed). The effect of the approximation has been shown to be negligible compared to the exact normalization factor [Lee and Chuang 1994]. In summary, the similarity is computed as

$$\Theta(D_i, Q) = \frac{\sum_{j=1}^V w_{Q,j} \times w_{i,j}}{\sqrt{d_i}}, \quad (2)$$

where  $d_i$  is the number of terms in document  $D_i$ .

### 3. WEIGHT-PARTITIONED SIGNATURE FILE

One problem of signature files is the difficulty in encoding term frequency information. Croft and Savino [1988] evaluated several approaches, which range from not representing term frequencies at all to approximating term frequencies with the number of text blocks containing the term. None of the approaches allow exact representation of term

frequencies. Wong and Lee [1990] proposed a method to represent the exact values of term frequencies by using a separate signature file for terms with the same term frequency. It is done by first grouping terms with the same term frequency in a document. Then, signatures are generated from each group and stored in the signature file corresponding to that term frequency. In this article we follow a similar approach. The main difference is that our approach uses superimposed coding, whereas Wong and Lee use an encoding scheme that, unlike superimposed codes, does not introduce any false drops.

Figure 2 is an example of the signature generation process, and Figure 3 illustrates the general design of the retrieval system. A document first goes through a preprocessing stage where common words are removed and keywords stemmed. Then, the keywords of a document are grouped according to term frequency. Signatures are generated from keywords in each group as in a traditional signature file, but signatures generated from different groups are stored in different signature files. The group corresponding to term frequency  $i$  is denoted by  $tf_i$ , and the corresponding signature file is denoted by  $SF_i$ . During a search, the query signature is compared to every  $SF_i$ . If a match is found in  $SF_i$ , the term frequency of the term producing the match has a term frequency equal to  $i$ , which can be used in computing the document score.

The algorithms for generating and searching the signature files are given in Algorithms 1 and 2, where  $n$  and  $N$  denote, respectively, the number of  $tf$ -groups and the number of documents in the collection. They are shown at the conceptual level for the purpose of illustration only and are not necessarily implemented as described in real implementations:

*Algorithm 1. Signature File Generation.*

1. sort keywords and tabulate term frequencies of keywords;
2. for  $i = 1$  to  $n$  do
3.     for every  $s$  keywords in  $tf$ -group  $i$  do
4.         generate signature and insert into  $SF_i$ ;
5.     end {for}
6. end {for}

*Algorithm 2. Signature File Search.*

1. for all term  $j$  in  $Q$  do
2.     obtain  $idf$  of term  $j$ ;
3.     generate signature for term  $j$ ;
4.     end {for}
5. for  $i = 1$  to  $N$  do
6.     for  $j = 1$  to  $n$  do
7.         for all term  $k$  in  $Q$  do
8.             if term  $k$ 's signature matches  $D_i$ 's signature in the  $j$ th partition then
9.                 increment  $D_i$ 's score by  $i$  by  $j \times idf_k \times w_{Q,k}$
10.                 break;
11.             end {if}
12.         end {for}
13.     end {for}
14. end {for}

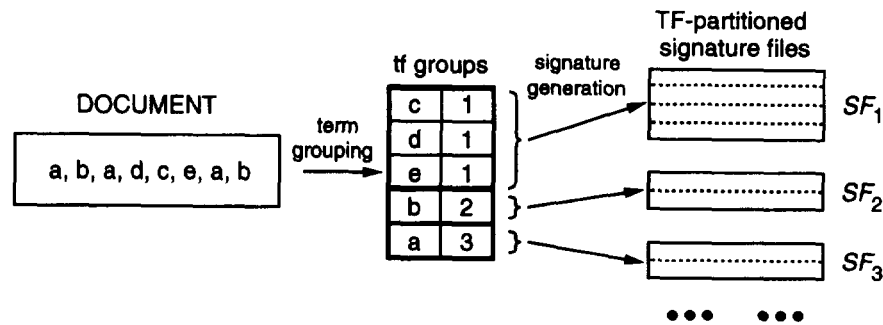


Fig. 2. Example of weight-partitioned signature files.

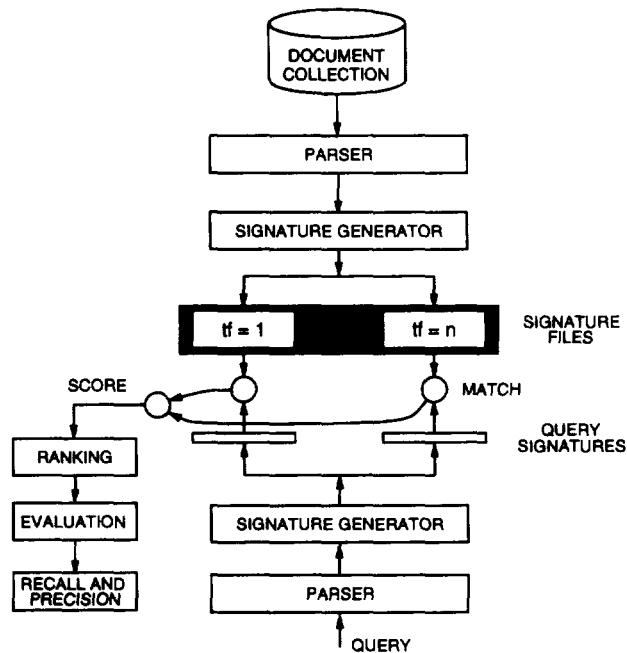


Fig. 3. Retrieval system based on weight-partitioned signature files.

In the weight-partitioned signature file, multiple occurrences of the same keyword in a document will generate only one word signature and, thus, are represented in the signature file only once. This is different from the traditional method where a keyword will be represented as many times as the number of text blocks containing it (although multiple occurrences of a keyword within the same text block will be collapsed). This collapsing of multiple occurrences into a single representation in the signature file significantly reduces the storage overhead for long documents; on the other hand, the distribution of a document's signatures into several signature files causes additional overhead. Furthermore, this method does not always search all of the signatures in the signature file; whenever a match is

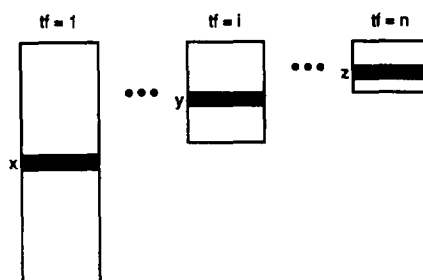


Fig. 4. Effects of search strategies.

found, the search can proceed to the next document, since a query signature will find *at most* one match (ignoring false drops) within a document. This explains the break statement in step 10 of Algorithm 2.

Since the signature file is partitioned, the partitions can be searched in different orders and may have different signature lengths. The effects of different choices are explained in Sections 3.1 and 3.2.

### 3.1 Search Strategies

A search can start from the signature file with the highest term frequency and move toward the signature file with the lowest term frequency, or vice versa. We refer to these two approaches, respectively, as the HL (high to low) and LH (low to high) methods. There are two contradicting factors affecting the merits of these two search orders, as shown in Figure 4. Suppose signatures  $x$  and  $z$  are false drops matching the query term being searched. If the search starts from the highest term frequency,  $z$  will be identified as a match. In this case, the effect of the false drop on the document score will be very high since  $z$  has a high term frequency. On the contrary, if the search starts from the lowest term frequency, the false drop  $x$  will have a small effect on the document score, since the term frequency is low. This suggests that LH is better than HL. However, there is a contradicting factor due to the fact that the number of signatures at high term frequencies is much less than that at low term frequencies, owing to the Zipfian distribution of terms among term frequencies. If we further assume that the probability that a true match occurs at  $SF_i$  is the same for all  $i$ , then, on the average, the chance of encountering a false drop before a true match is found is smaller in the HL approach than in the LH approach. In Figure 4, if  $y$  is a true match, then there will be much fewer signatures to the "right" than to the "left" of  $y$ . Consequently, the probability that the search will arrive at  $y$  without hitting any false drops is higher if the search starts from the "right" (i.e., high term frequency). Furthermore, the search speed is improved, since once a match is found the search can proceed immediately to the next query term. This factor is magnified if the user has a tendency to specify query terms that have high term frequencies (i.e., true matches tend to skew toward high term frequencies). This suggests that the search should start from the highest term frequency. The performance of these two search strategies is studied in Section 4.



### 3.2 Optimal Signature File Configuration

When the signatures are partitioned among several files, it is possible to use different signature configurations for different partitions. Since a false match at high term frequencies affects the document scores more than those at low term frequencies, it is desirable to reduce the false-drop probabilities at high term frequencies. On the other hand, if we set lower false-drop probabilities for the signatures with low term frequencies, the absolute number of false drops will be large, since there are more signatures at the low-frequency end. Therefore, there are conflicting factors affecting the determination of the signature lengths for the partitions. For a given storage overhead (i.e., a fixed total size for all of the signature partitions), we want to obtain an optimal way of allocating the storage to the partitions so as to minimize the effect of false drops on precision and recall. This is investigated in Section 4.

## 4. ANALYTICAL RESULTS

In superimposed coding, there are four major parameters related by the following formula [Roberts 1979]:

$$p = (1 - (1 - w/m)^s)^w,$$

where  $p$  is the false-drop probability;  $w$  is the weight of a word signature (number of bits set by a word);  $m$  is the signature length; and  $s$  is the number of words hashed into a block signature. It has been shown that, in order to minimize the false-drop probability for a given storage overhead, the number of bits set in a block signature must be equal to half of the signature length, and

$$w = (1/\ln 2)^2 s (1 - 2^{-1/s}) \ln(1/p), \quad (3)$$

$$m = (1/\ln 2)^2 s \ln(1/p). \quad (4)$$

Formula (3) can be approximated by

$$w = (1/\ln 2) \log_e(1/p), \quad (5)$$

when  $s \geq 5$ . These formulas are used later in this article to determine the values of  $m$ ,  $s$ , and  $w$  when two of these parameters are given.

To improve clarity, we sometimes use  $p(m, w, s)$  instead of  $p$  to indicate the relationship between  $p$ ,  $w$ ,  $m$ , and  $s$ . With this notation, the false-drop probability for signatures in partition  $SF_i$  is given by  $p(m_i, w_i, s_i)$ , or  $p_i$ .

In the experiment,  $p$  is a variable that is varied for each experiment run. The signature length,  $m$ , is a design parameter that is fixed for each test collection. Choosing a length that fits into the native word length of the machine is most efficient, but the native word length is short and, thus, will limit the achievable minimum false-drop probability. Furthermore, a low false-drop probability with a short signature length entails an ex-

tremely small number of words to be hashed in a signature, resulting in a large number of signatures and disk accesses (if buffering is not used). On the other hand, a long signature length is less efficient to manipulate and will generate underutilized signatures for short documents.

#### 4.1 HL versus LH Search Strategies

In this section we analyze the performance of the HL and LH search strategies. The performance of these two strategies depends on the distributions of false and true drops. The distribution of false drops further depends on the signature lengths of the partitions. In order to compare these two methods, we assume that the distributions and the signature lengths are uniform.

In the search algorithm, a signature file  $SF_i$  is treated as a single unit. A match with any signature in  $SF_i$  results in  $tf$  equal to  $i$ . We need to find the probability of a false match when we compare a random signature with all of the signatures in  $SF_i$ . We denote the false-drop probability against  $SF_i$  as  $q_i$  to distinguish it from the usual false-drop probability, which is the probability for two signatures to match combinatorially.

LEMMA 4.1.1. *The false-drop probability against  $tf_i$  is given by*

$$q_i = 1 - (1 - p(m_i, w_i, s_i))^{l_i} \approx l_i p(m_i, w_i, s_i),$$

where  $l_i$  is the number of signatures in  $SF_i$ .

PROOF. This is clear.  $\square$

We may use different  $m_i$ ,  $w_i$ , and  $s_i$  for different  $tf$ -groups. As long as

$$1 \geq p_1 \geq \dots \geq p_n \geq 0,$$

we have

$$0 \leq 1 - p_1 \leq \dots \leq 1 - p_n \leq 1.$$

From Zipf's Law, we know that  $l_1 \geq l_2 \geq \dots \geq l_{n-1}$  and  $l_1 > l_n$ . This implies that

$$(1 - p_1)^{l_1} \leq \dots \leq (1 - p_{n-1})^{l_{n-1}} \quad \text{and} \quad (1 - p_1)^{l_1} \leq (1 - p_n)^{l_n},$$

or, equivalently,  $q_1 \geq \dots \geq q_{n-1}$  and  $q_1 \geq q_n$ . These relations among the  $q_i$  imply that false-drop probabilities are higher at the low  $tf$  end than at the high  $tf$  end. This turns out to be very important in our analysis. One thing we want to point out is that  $p_1 \geq \dots \geq p_n$  does not necessarily mean that  $m_1 \geq \dots \geq m_n$  (neither does the other direction hold), since there are other parameters such as  $s_i$  and  $w_i$  involved.

LEMMA 4.1.2. *Given a random word signature, the expected  $tf$  at which a match is found can be obtained by  $S_n^{(n)}$  for HL and  $T_1^{(n)}$  for LH, where the*

definitions of  $S_i^{(n)}$  and  $T_i^{(n)}$  are given as follows:

- (1) For HL:  $S_n^{(i)} = q_i i + (1 - q_i) S_{i-1}^{(n)}$  for  $i = n, \dots, 2$ ,  
 $S_1^{(n)} = q_1$ .
- (2) For LH:  $T_i^{(n)} = i q_1 + (1 - q_i) T_{i+1}^{(n)}$  for  $i = 1, \dots, n - 1$ ,  
 $T_n^{(n)} = n q_n$ .

Furthermore,  $S_n^{(n)} \geq T_1^{(n)}$ .

PROOF. In the HL method, the probability that there is a match in  $tf_n$  is  $q_n$ . If there is a match in  $tf_n$ , the expected  $tf$  value is  $n$ . Otherwise, we have to check the next lower  $tf$  group. Therefore, we get the following recurrence relation:

$$S_n^{(n)} = q_n n + (1 - q_n) S_{n-1}^{(n)}.$$

For  $tf_1$ , there is no lower  $tf$  group; the expected  $tf$  value is given by  $q_1$ .  $T_n^{(n)}$  for the LH method can be derived in a similar way. That  $S_n^{(n)} \geq T_1^{(n)}$  is clear from the definitions of  $S_n^{(n)}$  and  $T_1^{(n)}$ .  $\square$

LEMMA 4.1.3. Suppose a word appears in a document with  $tf = j$ ; then HL yields the following expected  $tf$ :

$$j + ((n - j)q_n + (n - j - 1)q_{n-1} + \dots + q_{j+1});$$

and LH yields

$$j - ((j - 1)q_1 + (j - 2)q_2 + \dots + q_{j-1}),$$

when high-order terms such as  $q_i q_j$  and  $q_i q_j q_k$  are dropped.

PROOF. We will use the formulas in Lemma 4.1.2. The key in Lemma 4.1.3 is that the search will stop at  $tf_k$ , where  $k \geq j$  for HL and  $k \leq j$  for LH.

- (1) For HL:  $S_n^{(i)} = q_i i + (1 - q_i) S_{i-1}^{(n)}$  for  $i = n, \dots, 2$ ,  
 $S_1^{(n)} = q_1$ .

By eliminating  $S_i^{(n)}$  for  $i = n - 1, \dots, j$  from the above equations, we have

$$S_n^{(n)} = q_n n + (n - 1)q_{n-1}(1 - q_n) + \dots + (j + 1)q_{j+1}(1 - q_n) \dots (1 - q_{j+2}) \\ + j(1 - q_n)(1 - q_{n-1}) \dots (1 - q_{j+1}).$$

Multiplying out everything and dropping small high-order terms, we get

$$S_n^{(n)} \approx n q_n + \dots + (j + 1)q_{j+1} + j - j(q_n + \dots + q_{j+1}) \\ = j + ((n - j)q_n + (n - j - 1)q_{n-1} + \dots + q_{j+1}).$$

(2) For LH: As in HL, we have

$$T_1^{(n)} \approx j - ((j-1)q_1 + (j-2)q_2 + \dots + q_{j-1}).$$

Since false drops exist, it is clear that in the above two expressions  $S_n^{(n)} \geq j$  is an overestimation and that  $T_1^{(n)} \leq j$  is an underestimation.  $\square$

Since we know that the exact  $tf$  value is  $j$ , we can check the errors introduced by these two methods:

$$|S_n^{(n)} - j| \approx (n-j)q_n + (n-j-1)q_{n-1} + \dots + q_{j+1},$$

$$|T_1^{(n)} - j| \approx (j-1)q_1 + (j-2)q_2 + \dots + q_{j-1}.$$

It is easy to see that  $|S_n^{(n)} - j| \leq |T_1^{(n)} - j|$ , provided  $n-j \leq j-1$ . Therefore,  $S_n^{(n)}$  is a better approximation of the true  $tf$  value if  $j \geq (n+1)/2$ .

The number  $(n+1)/2$  is a very conservative estimation. In fact, since  $q_1$  (the false-drop probability for  $SF_1$ ) is typically much larger than  $q_n$  (the false-drop probability for  $SF_n$ ), the  $tf$  value above which the HL approach yields better results than the LH method is much smaller than  $(n+1)/2$ , owing to the Zipfian distribution of terms among term frequencies. In the following discussion,  $I_i(D)$  is used to denote the number of terms in  $D$  with term frequency  $i$ , and  $l_i(D)$  is used to denote the number of signatures in  $SF_i$  for document  $D$ . Given a document  $D$  and the parameters  $m_i$ ,  $s_i$ , and  $w_i$  for each  $tf_i$ , we have the relation  $l_i(D) \approx I_i(D)/s_i$ . Zipf's Law can be mathematically approximated as [Salton 1989]

$$I_i(D) = \frac{2}{i(i+1)} I_1(D) \quad \text{for } i \geq 1. \quad (6)$$

From formula (6), we can get

$$I_{\geq n}(D) \approx \frac{2}{n} I_1(D). \quad (7)$$

For simplicity, we will use  $I_n$  to denote  $I_{\geq n}$ .

**LEMMA 4.1.4.** *Suppose Zipf's Law holds for our document collection,  $n=30$  and*

$$\frac{p_1}{s_1} \geq \dots \geq \frac{p_n}{s_n}. \quad (8)$$

*Then for any term occurring in the document with  $tf \geq 4$ ,  $S_n^{(n)}$  is a better approximation of the  $tf$  value than  $T_1^{(n)}$ .*

**PROOF.** Without loss of generality, we assume that the average number of terms in  $tf_1$  for our document collection is  $I_1$  and that  $tf = j$ . Then, using (6)

and (7), we have

$$\begin{aligned}
|S_n^{(n)} - j| &\approx (n-j)l_n p_n + (n-j-1)l_{n-1} p_{n-1} + \dots + l_{j+1} p_{j+1} \\
&\approx \frac{2(n-j)I_1}{n} \frac{p_n}{s_n} + \dots + \frac{2I_1}{(j+1)(j+2)} \frac{p_{j+1}}{s_{j+1}} \\
&\leq I_1 \frac{p_j}{s_j} \left( \frac{2(n-j)}{n} + \frac{2(n-j-1)}{(n-1)n} + \dots + \frac{2}{(j+1)(j+2)} \right), \\
|T_1^{(n)} - j| &\approx (j-1)q_1 + (j-2)q_2 + \dots + q_{j-1} \\
&\geq I_1 \frac{p_j}{s_j} \left( (j-1) + \frac{2(j-2)}{2 \times 3} + \dots + \frac{2}{(j-1)j} \right).
\end{aligned}$$

Note that  $I_1 p_j / s_j$  is a common factor in both formulas. A simple computation will show that  $|S_n^{(n)} - j| \leq |T_1^{(n)} - j|$  for any  $j \geq 4$ .  $\square$

Condition (8) is not restrictive. Clearly, it is true when the signature lengths are uniform. We will show later that this condition also holds for our optimal signature configurations. From the proof, we can see that  $tf \geq 4$  is still a conservative result. In our experiment with the TREC subset, HL is better than LH for  $j \geq 2$  using the configuration  $n = 30$ ,  $m = 29$ ,  $s = 2$ , and  $w = 10$ . Lemma 4.1.3 can be used to find the exact  $j$  when the statistical data about the document collection are available.

Lemma 4.1.3 considers the situation where the query term indeed exists in the documents. We now consider the case when there is no such word in the document (i.e., the match is a false drop).

**LEMMA 4.1.5.** *In the case of a false drop, the expected  $tf$  value given by  $S_n^{(n)}$  and  $T_1^{(n)}$  all have the form*

$$\sum_{i=1}^n i q_i + \text{high-order terms}.$$

**PROOF.** Eliminating all  $S_i^{(n)}$  for  $i = n-1, \dots, 1$  from part (1) of Lemma 4.1.2, we get

$$S_n^{(n)} = \sum_{i=1}^n i q_i + \text{high-order terms}.$$

Similarly, we can obtain the expression for  $T_1^{(n)}$ .  $\square$

**THEOREM 4.1.6.** *If the false-drop probabilities  $p_i$  and terms per signature  $s_i$  satisfy*

$$\frac{p_1}{s_1} \geq \dots \geq \frac{p_n}{s_n},$$

*then the HL search strategy is superior than the LH method.*

**PROOF.** From Lemma 4.1.4, we can see that, for true matches, HL produces  $tf$  values that are closer to the true values than LH does on the average.

On the other hand, Lemma 4.1.5 indicates that the expected errors created by these two methods on a false match all have the form  $\sum_{i=1}^n iq_i +$  high-order terms. The difference between the errors introduced by the two methods is negligible, considering that the normal value of  $q_i$  is typically very small. Combining these two facts, we can claim that HL is better than LH.  $\square$

We will see in Section 4.2 that, in the optimal signature configuration, (1) the signature lengths of the partitions are increasing and (2) the false-drop probabilities are decreasing from low to high term frequencies. Under this condition, it is clear that HL is even more desirable than LH.

#### 4.2 Optimal Signature Configuration

We have implicitly assumed so far that every signature partition uses the same signature length  $m$ , the same  $w$ , and the same  $s$ . However, for the weight-partitioned signature file, different partitions may be assigned different signature lengths, different signature weights, etc. This implies that, for a given total storage overhead, we can control the amount of storage allocated to each partition by changing the configuration of that partition. Thus, for a fixed total storage overhead, there is an optimal way of assigning signature lengths (or storage) to each partition to minimize the impact of false drops on precision and recall. This section derives the optimal signature configurations.

Since different  $SF_i$  may have different signature configurations and, hence, different false-drop probabilities, the weighted average of the false-drop probabilities of all  $SF_i$  is used to characterize the false-drop probability of the whole signature file. The effective false-drop probability can be written as

$$\begin{aligned} p(\vec{m}, \vec{w}, \vec{s}) &= \sum_{i=1}^n \frac{l_i(D)}{\sum_{j=1}^n l_j(D)} p(m_i, w_i, s_i) \\ &= \sum_{i=1}^n \frac{I_i(D)}{\sum_{j=1}^n I_j(D)} p(m_i, w_i, s_i), \end{aligned} \tag{9}$$

where  $\vec{m}$  denotes  $m_1, m_2, \dots, m_n$ . In this formula,  $l_i(D)$  and  $I_i(D)$ , for  $i = 1, \dots, n$ , are document-dependent variables. We will use Zipf's Law to eliminate them. Applying (6) and (7) to (9), we get a formula for the document-independent, effective false-drop probability:

$$p = \sum_{i=1}^{n-1} \frac{1}{i(i+1)} p(m_i, w_i, s_i) + \frac{1}{n} p(m_n, w_n, s_n). \quad (10)$$

When  $m_1 = \dots = m_n$ ,  $s_1 = \dots = s_n$ , and  $w_1 = \dots = w_n$  (i.e., all  $SF_i$  use the same signature configuration), we can easily check that

$$p(\vec{m}, \vec{w}, \vec{s}) = p(m_1, w_1, s_1) = \dots = p(m_n, w_n, s_n).$$

In other words, the uniform signature structure is just a special case of the nonuniform signature structure. We now present a result about the relation between the uniform signature length and the nonuniform signature length approaches.

**LEMMA 4.2.1.** *For a signature file with parameters  $m$ ,  $s$ , and  $w$ , reducing  $s$  by a factor of  $\lambda$  has the same effect on false-drop probability as increasing  $m$  by a factor of  $\lambda$ , where  $\lambda \geq 1$ . Furthermore, increasing or decreasing  $m$  and  $s$  by the same factor does not change the false-drop probability.*

**PROOF.** We will use the false-drop probability formula

$$p(m, w, s) = (1 - (1 - w/m)^s)^k.$$

In fact, for small  $w$  and  $s$  (relative to  $m$ ), we have

$$\left(1 - \frac{w}{m}\right)^s = \left[\left(1 - \frac{1}{m/w}\right)^{m/w}\right]^{ws/m} \approx e^{-ws/m}.$$

Thus, we have  $p(m, w, s) \approx (1 - e^{-ws/m})^w$ . Now increasing  $m$  by a factor of  $\lambda$ , we get

$$p(\lambda m, w, s) \approx (1 - e^{-ws/\lambda m})^w.$$

Reducing  $s$  by a factor of  $\lambda$ , we get

$$p(m, w, s/\lambda) \approx (1 - e^{-w(s/\lambda)/m})^s = (1 - e^{-ws/\lambda m})^w.$$

We can reach the second conclusion in a similar way.  $\square$

To obtain formulas to determine the optimal signature configurations for the partitions, we first obtain a formula for the general case, where each

partition is associated with a general ranking weight  $W$ .<sup>1</sup> Then, we can apply the formula to the case where  $tf$  is used as the ranking weight (we called this the TF-partitioned signature file). Theorem 4.2.4 determines the optimal signature configurations for the signature partitions to minimize the effect of false drops on retrieval effectiveness. The condition under which the theorem applies is similar to that for the false-drop probability formula [Roberts 1979] (i.e., when  $w_i$ ,  $s_i$ , and  $l_i$  are small relative to  $m_i$ , which is true for most applications).

It is well known that for the signature file approach the space needed is proportional to  $\ln(1/p)$ . Given the space constraint, we first compute the optimal false-drop probabilities for all  $SF_i$ . From the false-drop probabilities, we can compute the other parameters, such as signature lengths, bits per word, etc.

We assume that document terms are partitioned into different groups, say,  $G_1, G_2, \dots, G_n$ , and that each term in  $G_i$  carries a weight  $W_i$ . Without loss of generality, we can assume that  $W_1 \leq \dots \leq W_n$ . The corresponding signature files are denoted as  $SF_1, SF_2, \dots, SF_n$ , as before. Under these assumptions, the storage overhead for a document is  $\sum_{i=1}^n l_i m_i$ , where  $l_i$  is the number of signatures in  $SF_i$ , and  $m_i$  is the length of signatures in  $SF_i$ .

**LEMMA 4.2.2.** *For false matches, the expected weights given by  $S_n^{(n)}$  and  $T_1^{(n)}$  all have the form*

$$\sum_{i=1}^n W_i q_i + \text{high-order terms}.$$

**PROOF.** Similar to Lemma 4.1.5.  $\square$

**LEMMA 4.2.3.** *Suppose the storage constraint is  $\sum_{i=1}^n l_i m_i = C$ . In order to minimize the effect of false drops, the relations between different false-drop probabilities  $p_i$  are*

$$\frac{W_i}{s_i} p(m_i, w_i, s_i) = \frac{W_j}{s_j} p(m_j, w_j, s_j) \quad \text{for } 1 \leq i \leq j \leq n,$$

$$p(m_n, w_n, s_n) = \exp\left(-\frac{C(\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln(W_n s_k / W_k s_n)}{\sum_{k=1}^n l_k s_k}\right).$$

**PROOF.** We know from the previous result that the expected  $tf$  value is

$$\sum_{i=1}^n W_i q_i$$

<sup>1</sup> In order to distinguish this weight, which is used to determine the importance of a match in the vector-space model, from the signature weight, which is the number of ones in a signature, the former is hereafter called *ranking weight*.



for a false drop. We will minimize the function  $\sum_{i=1}^k W_i q_i$  under the storage constraint  $\sum_{i=1}^n l_i m_i = C$ .

The existence of a minimum value is clear, since (1) the constraint defines a compact set in space  $R^n$  and (2) the function is continuous. It is also not difficult to see that the minimum value cannot be achieved on boundaries. First, from Lemma 4.1.1, we have

$$\sum_{i=1}^n W_i q_i \approx \sum_{i=1}^n W_i l_i p_i.$$

There are too many variables to handle. We first eliminate all  $m_i$ . From formula (4), we have

$$m_i = (1/\ln 2)^2 s_i \log_e(1/p_i).$$

The constraint can be rewritten as

$$C = \sum_{i=1}^n l_i m_i = \sum_{i=1}^n (1/\ln 2)^2 s_i l_i \log_e(1/p_i).$$

Lagrange's method can be used to obtain the minimum. Let

$$F(\lambda, \tilde{m}, \tilde{w}, \tilde{s}) = \lambda \left( \sum_{i=1}^n (1/\ln 2)^2 s_i l_i \log_e(1/p_i) - C \right) + \sum_{i=1}^n W_i l_i p_i.$$

The partial derivatives with respect to  $p_i$  and  $\lambda$  are

$$\frac{\partial F}{\partial p_i} = \lambda \frac{l_i s_i}{p_i} + W_i l_i \quad i = 1, \dots, n,$$

$$\frac{\partial F}{\partial \lambda} = \sum_{i=1}^n (1/\ln 2)^2 s_i l_i \log_e(1/p_i) - C.$$

Now set all of the partial derivatives to zero, and cancel  $l_i$  from the  $i$ th equation for  $i = 1, \dots, n$ . We have

$$\frac{W_i}{s_i} p_i + \lambda = 0 \quad i = 1, \dots, n, \quad (11)$$

$$\sum_{i=1}^n (1/\ln 2)^2 s_i l_i \log_e(p_i) + C = 0 \quad (12)$$

Eliminating  $\lambda$  from the  $i$ th and  $j$ th equations, we get

$$\frac{W_i}{s_i} p(m_i, w_i, s_i) = \frac{W_j}{s_j} p(m_j, w_j, s_j)$$

for  $1 \leq i \leq j \leq n$ . Setting  $j = n$ , we get

$$p_i = \frac{W_n s_i}{W_i s_n} p_n.$$

Substituting these equations into Eq. (12), we have

$$\sum_{i=1}^n s_i l_i \ln \frac{W_n s_i}{W_i s_n} + \left( \sum_{i=1}^n s_i l_i \right) \ln p_n + (\ln 2)^2 C = 0.$$

Solving for  $p_n$ , we get

$$p_n = \exp \left( - \frac{C (\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right). \quad \square$$

There are several things worth pointing out. First, since  $W_1 \leq \dots \leq W_n$ , we have

$$\frac{p_1}{s_1} \geq \dots \geq \frac{p_n}{s_n}.$$

These relations ensure that HL is better than LH from Lemma 4.1.4. Second, it is a little surprising that the signature file size  $l_i$  disappeared in

$$\frac{W_i}{s_i} p(m_i, w_i, s_i) = \frac{W_j}{s_j} p(m_j, w_j, s_j).$$

This can be explained intuitively by the fact that there are two conflicting factors: large  $l_1$  means more false drops, thus requiring more storage to be assigned to  $tf_1$ ; on the other hand, with a fixed space, we cannot assign too much storage to  $tf_1$ . The optimal method we are presenting will give us an optimal way of allocating storage among all  $tf$ -groups.

**THEOREM 4.2.4.** *Under the same assumption as in the previous lemma, to minimize the effect of false drops, the signature length  $m_i$ , bits per word  $w_i$ ,*

and words per signature  $s_i$  must satisfy

$$m_i = \frac{1}{\ln^2 2} s_i \left( \ln \left( \frac{W_i s_n}{W_n s_i} \right) + \frac{C(\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right), \quad (13)$$

$$w_i = \frac{1}{\ln 2} s_i (1 - 2^{-1/s_i}) \left( \ln \left( \frac{W_i s_n}{W_n s_i} \right) + \frac{C(\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right), \quad (14)$$

where  $i = 1, \dots, n$ .

PROOF. From Lemma 4.2.3, we know for  $1 \leq i \leq j \leq n$  that

$$\frac{W_i}{s_i} p_i = \frac{W_j}{s_j} p_j, \quad p_n = \exp \left( - \frac{C(\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right).$$

If we let  $j = n$ , we get  $p_i = (W_n s_i / W_i s_n) p_n$ . This implies that

$$p_i = \frac{W_n s_i}{W_i s_n} \exp \left( - \frac{C(\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right).$$

From formulas (3), (4), and (15), we arrive at

$$m_i = \frac{1}{\ln^2 2} s_i \ln \frac{1}{p_i} = \frac{1}{\ln^2 2} s_i \left( \ln \left( \frac{W_i s_n}{W_n s_i} \right) + \frac{C(\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right),$$

$$w_i = \frac{1}{\ln 2} s_i (1 - 2^{-1/s_i}) \ln \frac{1}{p_i}$$

$$= \frac{1}{\ln 2} s_i (1 - 2^{-1/s_i}) \left( \ln \left( \frac{W_i s_n}{W_n s_i} \right) + \frac{C(\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right)$$

with some simple algebraic operations.  $\square$

In order to use Theorem 4.2.4, we need to decide on  $C$ ,  $s_i$ , and  $l_i$ . Note that only  $C$  and  $s_i$  are design parameters, since  $l_i = \lceil I_i / s_i \rceil$ , which is determined by the document collection.  $C$  can be determined from the storage overhead desired. When  $C$  is fixed, the other parameters, namely,  $s_i$ , need to be determined. The selection of  $s_i$  needs some consideration. Suppose in a document collection that each document has on the average  $I_i$  terms with term frequency  $i$ . Then, a factor of  $I_1$  should be used as  $s_1$  to make sure we hash approximately the same number of words into a signature. For example, the TREC subset has on the average 105 terms in  $I_1$ . Then,  $s_1 = 5$  is a good choice since each signature in  $SF_1$  will contain 5 terms. On the other hand,  $s_1 = 8$  will be a bad choice since one of the signatures will contain only one word. After we select  $s_1 = s$ ,  $s_i = \lceil \min(s, I_i) \rceil$  is the logical choice for  $SF_i$ , where  $i = 2, \dots, n$ . Finally, since we have  $s_i l_i = \lceil I_i \rceil$ , the optimal formulas (13) and (14) can be simplified as

$$m_i = \frac{1}{\ln^2 2} s_i \left( \ln \left( \frac{W_i s_n}{W_n s_i} \right) + \frac{C (\ln 2)^2 + \sum_{k=1}^n \lceil I_k \rceil \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n \lceil I_k \rceil} \right), \quad (15)$$

$$w_i = \frac{1}{\ln 2} s_i (1 - 2^{-1/s_i}) \left( \ln \left( \frac{W_i s_n}{W_n s_i} \right) + \frac{C (\ln 2)^2 + \sum_{k=1}^n \lceil I_k \rceil \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n \lceil I_k \rceil} \right). \quad (16)$$

Although it is natural to use different signature lengths in the weight-partitioned signature file method, it is difficult in practice (e.g., in terms of programming effort) to use different signature lengths for different  $SF_i$  because the signature lengths must be memorized for signature generation. A better alternative is to use the same signature length for all signature partitions, but with different  $s_i$  values. The next corollary gives an optimal signature file configuration in which all signatures have the same length.

**COROLLARY 4.2.5.** *Suppose the storage constraint is  $C$ . Then the following signature file configuration will minimize the effect of false drops:*

$$m = \min\{\bar{m}_1, \dots, \bar{m}_n\}, \quad s_i = \left\lceil \frac{m}{\bar{m}_i} * \bar{s}_i \right\rceil \quad \text{for } i = 1, \dots, n,$$

$$w_i = \frac{1}{\ln 2} s_i (1 - 2^{-1/s_i})$$

$$\cdot \left( \ln \left( \frac{W_i s_n}{W_n s_i} \right) + \frac{C (\ln 2)^2 + \sum_{k=1}^n s_k l_k \ln \left( \frac{W_n s_k}{W_k s_n} \right)}{\sum_{k=1}^n l_k s_k} \right) \quad \text{for } i = 1, \dots, n,$$

where  $l_i = \lceil I_i/s_i \rceil$  and where  $w_i$ ,  $\bar{m}_i$ , and  $\bar{s}_i$  satisfy Theorem 4.2.4.

PROOF. We know from Lemma 4.2.1 that increasing or decreasing  $m$  and  $s$  by the same factor simultaneously will not change the false-drop probability. We first use Theorem 4.2.4 to compute all  $\bar{m}_i$  and  $\bar{s}_i$ . Then, we select the shortest signature length (or whichever we prefer) as our new uniform signature length  $m$  and obtain  $s_i$  by scaling down  $\bar{s}_i$  by the factor  $m/\bar{m}_i$ . The weight of each term,  $w_i$ , is the same as in Theorem 4.2.4.  $\square$

We can see that Theorem 4.2.4 and Corollary 4.2.5 both depend on  $I_1, \dots, I_n$ , the sizes of the term frequency groups. These values can be determined by taking the averages of the term frequency group sizes over the document collection. Alternatively, we can use Zipf's Law to get an approximation of  $I_2, \dots, I_n$  from  $I_1$ . From (6) and (7), we get

$$s_i = \lceil \min(s, I_i) \rceil = \left\lceil \min\left(s, \frac{I_1}{i(i+1)}\right) \right\rceil.$$

Then, we can use Theorem 4.2.4 to determine  $m_i$  and  $w_i$  and to generate the signature files.

The general result can now be applied to the TF-partitioned signature file using the term frequencies as the ranking weights. That is,  $W_i = i$  in Eqs. (13)–(16). This amounts to assigning weights to a match in proportion to the term frequencies. Since false drops occurring at high term frequencies will have a large negative impact on the retrieval effectiveness, we want to allocate more storage to high-term-frequency partitions by increasing their weights.

#### 4.3 Effect of False Drops on Document Ranks

Another application of Lemma 4.2.3 is on false-drop elimination. False-drop elimination is expensive to perform. When false-drop elimination is not performed, we want to verify analytically that reducing the number of false drops beyond a certain threshold yields no gain in precision and recall. Furthermore, if  $p$  denotes the exact value of a parameter, then  $p'$  denotes the approximate value of the parameter derived from the weight-partitioned signature file (i.e., with false drops). For example,  $tf'(t)$  denotes the  $tf$  value of term  $t$  obtained in the search process, and  $tf(t)$  denotes the true  $tf$  value of  $t$  in the document;  $\Theta'(D, Q)$  denotes the similarity obtained from the weight-partitioned signature file, and  $\Theta(D, Q)$  denotes the actual similarity if no false drops exist.  $\Theta'(D, Q)$  is also used to denote the expected value of  $\Theta(D, Q)$ , for simplicity. The expected  $\Theta'(D, Q)$  is a theoretical representation of  $\Theta'(D, Q)$ , on the average.

LEMMA 4.3.1. *Given a query  $Q$ , the difference between the true similarity and the expected similarity of a document  $D$ , obtained using the weight-partitioned signature file, is a linear combination of  $q_i$ .*

PROOF. We divide the query terms into two groups:  $Q = Q_1 \cup Q_2$ , where  $Q_1 = D \cap Q$  and  $Q_2 = Q \setminus D$ . The true similarity between document  $D$  and query  $Q$  is given by

$$\Theta(D, Q) = \frac{\sum_{t \in Q_1} tf(t) \times idf(t)}{\sqrt{d}},$$

and the expected similarity obtained by the weight-partitioned signature file method is

$$\Theta'(D, Q) = \frac{\sum_{t \in Q_1} tf'(t) \times idf(t) + \sum_{t \in Q_2} tf'(t) \times idf(t)}{\sqrt{d}}.$$

Due to errors caused by false drops, the  $tf$ -value obtained by the HL method may be larger than the true  $tf$ -value, although the difference is typically small. The formula for the expected  $tf$  values of the terms in  $Q_1$  is given by Lemma 4.1.3, that is,

$$tf'(t) - tf(t) \approx ((n - tf(t))q_n + (n - tf(t) - 1)q_{n-1} + \dots + q_{tf(t)+1}).$$

And the expected  $tf$ -value for elements of  $Q_2$  is given by

$$tf'(t) = \sum_{i=1}^n iq_i.$$

The right-hand sides of these two expectations are linear combinations of  $q_i$  for a fixed query  $Q$  and document  $D$ . From the above formulas, we have

$$\begin{aligned} \Theta'(D, Q) - \Theta(D, Q) &= \frac{\sum_{t \in Q_1} (tf'(t) - tf(t)) \times idf(t)}{\sqrt{d}} + \frac{\sum_{t \in Q_2} (tf'(t)) \times idf(t)}{\sqrt{d}} \\ &= \sum_{t \in Q_1} \frac{idf(t)}{\sqrt{d}} (tf'(t) - tf(t)) + \sum_{t \in Q_2} \frac{idf(t)}{\sqrt{d}} (tf'(t)), \end{aligned}$$

which is clearly a linear combination of  $q_i$  since  $tf'(t) - tf(t)$  for  $t \in Q_1$  and  $tf'(t)$  for  $t \in Q_2$  are for every term  $t$  in  $Q$ .  $\square$

From Lemma 4.3.1, we can see that the expected  $\Theta'(D, Q)$  approaches  $\Theta(D, Q)$  as  $q_i$  approach zero. For query  $Q$ , the vector-space model assigns a score  $\Theta(D, Q)$  to each document. For simplicity, assume that no two documents receive the same score. In this way, a complete ranking of all the documents in  $\mathcal{D}$  is produced. Under this assumption, the following theorem can be obtained:

**THEOREM 4.3.2.** *For a document collection  $\mathcal{D}$  and a set of queries  $\mathcal{Q}$ , there is a constant  $\epsilon$  such that, when the false-drop probabilities for all  $SF_i$  are smaller than  $\epsilon$ , the expected precision and recall resulting from searching the weight-partitioned signature file for  $Q \in \mathcal{Q}$  will be the same as that without false drops. In other words, reducing the false-drop probabilities beyond this  $\epsilon$  yields no gain in precision and recall.*

**PROOF.** For each query  $Q \in \mathcal{Q}$ , the vector-space model assigns a unique score  $\Theta(D, Q)$  to each document  $D$ . Sorting all documents by decreasing order of their scores, we obtain a complete ranking. Use  $\delta_Q$  to denote the minimum difference between two adjacent documents. In other words,

$$\delta_Q = \min\{\Theta(D_i, Q) - \Theta(D_{i+1}, Q) \mid i \geq 1\}.$$

Finally, let  $\delta = \min\{\delta_Q \mid Q \in \mathcal{Q}\}$ . Since no two documents receive the same score,  $\delta$  is a positive number. From Lemma 4.3.1, for each document  $D$  and query  $Q$ , we have

$$\lim_{q_i \rightarrow 0} (\Theta'(D, Q) - \Theta(D, Q)) = 0.$$

Thus, there must be a positive number  $\epsilon$  such that, when all of the false-drop probabilities are less than  $\epsilon$ ,

$$0 \leq \Theta'(D, Q) - \Theta(D, Q) < \delta.$$

This means that the expected ranking order given by  $\Theta'(D, Q)$  will be identical to the ranking order given by  $\Theta(D, Q)$ . Reducing false-drop probabilities beyond this  $\epsilon$  will not alter the precision and recall values.  $\square$

It is not always the case that no two documents receive the same score, but we can modify Theorem 4.3.2 to imply there is a number  $\epsilon$  such that, if all false-drop probabilities are less than  $\epsilon$ ,

$$\Theta(D_1, Q) > \Theta(D_2, Q) \Rightarrow \Theta'(D_1, Q) > \Theta'(D_2, Q).$$

Theorem 4.3.2 does not give any method to compute the threshold  $\epsilon$ . The main difficulty is that we do not have a prior collection of queries. Even with this knowledge, obtaining this number is not simple. Constant update of the collection makes it even more complicated to predict. Sampling or experimentation may be a practical method to obtain this threshold  $\epsilon$ .

## 5. EXPERIMENTAL RESULTS

The criteria for retrieval effectiveness are precision and recall. The goal of the experiments is

- (1) to observe the extent to which false drops can be tolerated without causing unacceptable degradation to recall and precision,

Table II. Assignment of Signature Lengths

Term frequency	Multiplying factor
1	1
2	2
3–29	4
30	8

- (2) to confirm the analytical results on the effectiveness of the LH and HL search strategies, and
- (3) to observe the performance of various signature length assignment methods.

Experiments were performed with a document subset extracted from the TREC collection [Harman 1993a]. The full TREC collection supplied by NIST consists of about 2GB worth of articles. Since the full collection exceeded the resource of our computing environment, a subset, hereafter called the TREC subset, of 10,000 *Wall Street Journal* articles from 1987, with a raw size of about 27MB, was used.<sup>2</sup>

In the experiments, a set of 25 standard queries was used; each of the queries came with a standard relevance judgment. The precision and recall values are obtained as follows. First, the standard queries are run on the collections. The ranked outputs of the queries are then compared to their respective list of relevant documents. Finally, the precision and recall values are calculated and averaged over all of the queries. In the computation, the recall values are rounded to the nearest 0.05 increment, and interpolation is used to obtain the precision values for all 21 recall points (from 0 to 1, at 0.05 increments). Note that, owing to interpolation, the precision at the zero recall point is actually the precision of the recall point that is nearest to zero. The 21 pairs of values can then be plotted in a precision and recall graph. This procedure is commonly used in evaluating text retrieval methods [Harman 1993b; Salton and McGill 1983]. For clarity and conciseness, we further average the 21 precision values to obtain the *average precision* for each set of parameters used in the experiments. In order to observe the effects of false drops on precision, recall, and storage overhead, we used a very wide range of false-drop probabilities in the experiments. In practice, the false-drop probability is typically below 1%.

Three different signature length assignment methods are tested in the experiments. The first two are derived from the optimal formulas (Eqs. (13) and (14)) using (1) term frequencies as ranking weights of the signature partitions (denoted as TF in Figures 5–7) and (2) uniform weights (denoted as U). The third method assigns signature lengths based on the data in Table II.

<sup>2</sup> The subset was extracted at the Online Computer Library Center (OCLC), Dublin, Ohio.



This method is based on the intuition that longer signature lengths should be used for higher term frequencies to reduce false drops occurring at high TF partitions. However, the multiplying factors are arbitrarily chosen. The inclusion of this method provides a data point with which the other two methods can compare. This method is denoted as EXP in Figures 5–7.

Note that, in all experiments, the maximum term frequency is set to 30. For large collections (used in experiments described later), there may be many terms with term frequencies greater than 30 and thus are mapped to  $tf = 30$ . This may result in a larger number of terms for  $tf = 30$  (or whatever ceiling was chosen).

The storage overhead is used as the main controlling parameter in the experiments. Storage overhead is measured by the size of the signature file divided by the total size of the documents *after stemming and common-word removal*. This measures the true storage overhead attributable to the signature file mechanism (as opposed to the space reduction achieved by stemming and common-word removal). When the raw collection size is used, the storage overhead is much lower. For instance, in the TREC subset, the storage overhead calculated using the raw collection size is typically only half of the storage overhead calculated using our method. For a fixed storage overhead, the signature lengths and the other related parameters are computed based on the three methods described above.

### 5.1 Performance of HL and LH Strategies

Figure 5 compares the average precisions of the HL and LH search strategies using the TREC subset. It is clear that the HL method is consistently better than the LH method, as predicted by our analysis. The difference is large for high false-drop probabilities (i.e., low storage overheads), but narrows down as the false-drop probability decreases. In general, the gap between the HL and LH methods increases when the number of false drops increases. This explains the large gap in the EXP method, because the EXP method produces the largest number of false drops given a fixed storage overhead. Since the HL method is much better than the LH method, we will only show the performance of the HL method in subsequent discussion.

### 5.2 Comparison of Signature Length Assignment Methods

The performance of the three signature file generation methods for the HL search strategies is shown in Figure 6. We can see from the figures that the EXP method is consistently the worst in all of the test runs, whereas the uniform and optimal methods give much better performance. Furthermore, the optimal method is the best everywhere. In fact, the LH approach in the optimal method is as good as the HL method in the uniform method. This confirms our previous discussion that the optimal method utilizes the space resource more effectively.

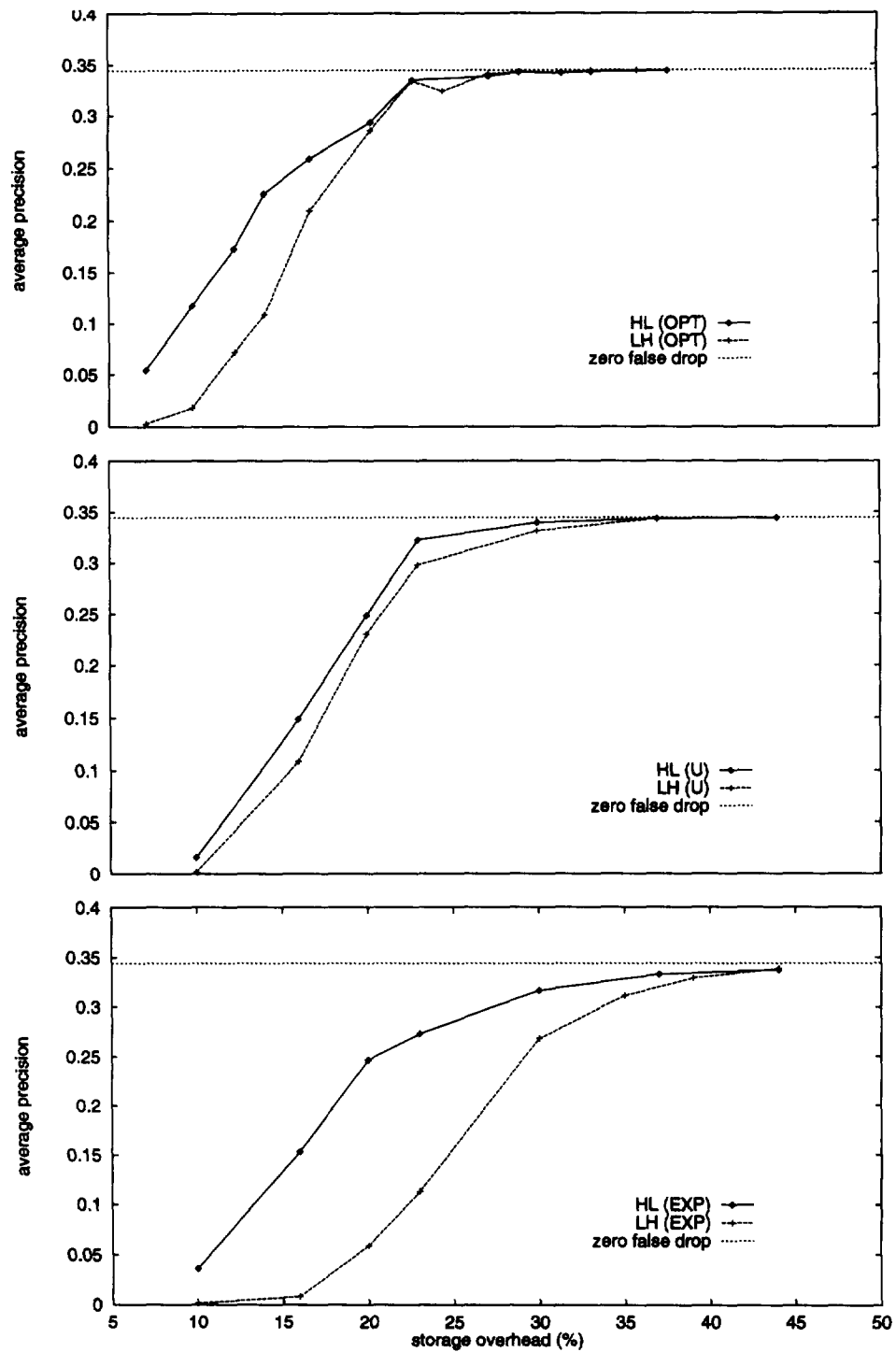


Fig. 5. Performance of LH versus HL search strategies.

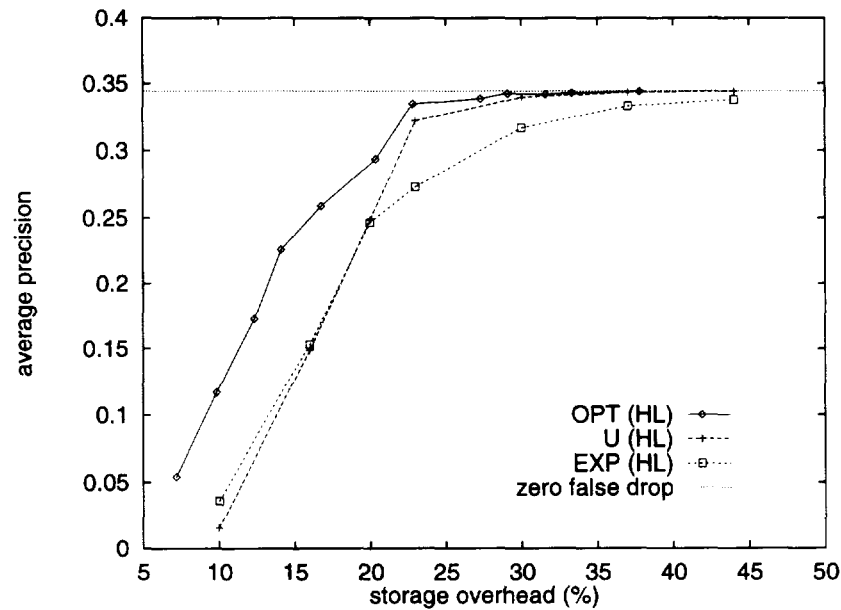


Fig. 6. Comparison of signature-length-assignment methods (HL).

### 5.3 Effect of False Drops on Retrieval Effectiveness

We can see from the figures that, although the average precision in general improves as false-drop probability decreases, there is no observable improvement on retrieval effectiveness beyond a certain threshold. For the TREC subset, the average precision reaches the perfect level (when there are no false drops) when the storage overhead is about 37%. According to Figure 7, this corresponds to roughly 0.01% false-drop probability and about 150 false drops generated from a query.

### 5.4 Number of False Drops and False-Drop Probability

Figure 7 shows the average number of false drops produced by a query and the weighted false-drop probability. The former is obtained by counting the total number of false drops in each run and then dividing it by the number of queries in each run. The latter is computed from the analytical formula (Eq. (9)). We can see that the analytical and experimental data are very consistent with each other. The uniform and TF methods are almost identical in the number of false drops, as well as in false drop probability. The EXP method is much worse than the other two methods.

## 6. CONCLUSION AND FUTURE WORK

This article presents a general technique for implementing the vector-space statistical model with signature files. In this technique, terms in a document are grouped according to their term frequencies and hashed as a group into the signature file with the corresponding term frequency. Thus,

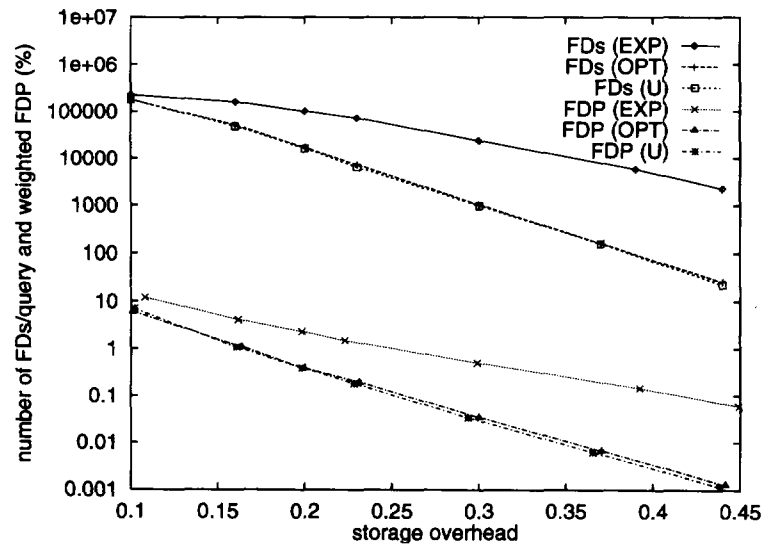


Fig. 7. Number of false drops (FDs), and weighted false-drop probability (FDP) (HL).

when a match is found, the signature file in which the match occurs will indicate the term frequency of the matching term. This avoids storing the term frequency of each term explicitly.

We have presented an analytical study of the performance of two search strategies, namely, the LH and HL methods, and have obtained an optimal way of assigning signature lengths, and thus storage, to the signature partitions to minimize the effect of false drops on precision and recall. Experiments were performed to support our analytical study.

The work presented in this article is a significant step toward a viable way of supporting document ranking with signature files. The signature file technique will become much more competitive in speed if it does not have to perform false-drop elimination. Another significant observation from the experiments is that a storage overhead of around 25% gives a precision that is quite comparable to the ideal case where there is no false drops. There is little gain in retrieval effectiveness if a larger storage overhead is used. However, there are still many interesting problems to be investigated. For instance, we are investigating search heuristics to improve search performance further and to reduce the effect of false drops [Wong and Lee 1993]. We are also investigating methods that use the signature file to produce a coarse ranking and to apply exact ranking (without false drops) only to the top-ranked documents, thus obtaining low storage overhead and high precision [Knaus and Schäube 1993].

#### ACKNOWLEDGMENT

The authors are indebted to the anonymous reviewers for their comments, which led to many improvements to this article.

## REFERENCES

- AHO, A. V. AND ULLMAN, J. D. 1979. Optimal partial-match retrieval when fields are independently specified. *ACM Trans. Database Syst.* 4, 2 (June), 168–179.
- CHRISTODOULAKIS, S., THEODORIDOU, M., HO, F., PAPA, M., AND PATHRIA, A. 1986. Multimedia document presentation, information extraction and document formation in MINOS: A model and a system. *ACM Trans. Off. Inf. Syst.* 4, 4 (Oct.), 345–383.
- CROFT, W. B. AND SAVINO, P. 1988. Implementing ranking strategies using text signatures. *ACM Trans. Off. Inf. Syst.* 6, 1 (Jan.), 42–62.
- HARMAN, D. 1993a. Overview of the first Text Retrieval Conference. In *Proceedings of the 16th Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval* (Pittsburgh, Pa, June). ACM, New York, 36–47.
- HARMAN, D. 1993b. Overview of the second Text Retrieval Conference (TREC-2). In *Proceedings of the 2nd Text Retrieval Conference (TREC-2)* (Gaithersburg, Md., Aug.), 1–20.
- KNAUS, D. AND SCHAUBE, P. 1993. Effective and efficient retrieval from large and dynamic document collections. In *Proceedings of the 2nd Text Retrieval Conference (TREC-2)* (Gaithersburg, Md., Aug.), 163–170.
- LEE, D. L. AND CHUANG, H. 1994. Performance of document ranking and relevance feedback. Submitted for publication.
- LEE, D. L. AND LENG, C.-W. 1989. Partitioned signature files: Design issues and performance evaluation. *ACM Trans. Inf. Syst.* 7, 2 (Apr.), 158–180.
- LEE, D. L. AND LENG, C.-W. 1990. A partitioned signature file structure for multiattribute and text retrieval. In *Proceedings of the 6th International Conference on Data Engineering* (Los Angeles, Calif., Feb.), 389–397.
- LEE, D. L., KIM, Y. M., AND PATEL, G. 1994. Efficient signature file methods for text retrieval. *IEEE Trans. Data Knowl. Eng.* To be published.
- LIPMAN, D. J. AND PEARSON, W. R. 1985. Rapid and sensitive protein similarity searches. *Science* 227 (Mar.), 1435–1441.
- ROBERTS, C. S. 1979. Partial-match retrieval via method of superimposed codes. *Proc. IEEE* 67, 12 (Dec.), 1624–1642.
- SALTON, G. 1989. *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer*. Addison-Wesley, Reading, Mass.
- SALTON, G. AND BUCKLEY, C. 1988. Term-weighting approaches in automatic text retrieval. *Inf. Process. Manage.* 24, 5, 513–523.
- SALTON, G. AND MCGILL, M. J. 1983. *Introduction to Modern Information Retrieval*. McGraw-Hill, New York.
- STANFILL, C. AND KAHLE, B. 1986. Parallel free-text search on the connection machine system. *Commun. ACM* 29, 12 (Dec.), 1229–1239.
- WONG, W. Y. P. AND LEE, D. L. 1990. Signature file methods for implementing a ranking strategy. *Inf. Process. Manage.* 26, 5 (Sept.), 641–653.
- WONG, W. Y. P. AND LEE, D. L. 1993. Implementations of partial document ranking using inverted files. *Inf. Process. Manage.* 29, 5 (Sept.), 647–699.
- ZEZULA, P., RABITTI, F., AND TIBERIO, P. 1991. Dynamic partitioning of signature files. *ACM Trans. Inf. Syst.* 9, 4 (Oct.), 336–339.

Received November 1993; revised February 1994; accepted November 1994