

# Efficient and Effective Metasearch for a Large Number of Text Databases

Clement Yu<sup>1</sup>, Weiyi Meng<sup>2</sup>, King-Lup Liu<sup>1</sup>, Wensheng Wu<sup>1</sup>, Naphtali Risse<sup>3</sup>

<sup>1</sup> Dept. of EECS, University of Illinois at Chicago, Chicago, IL 60607

<sup>2</sup> Dept. of Computer Science, SUNY – Binghamton, Binghamton, NY 13902

<sup>3</sup> School of Computer Science, Florida International University, Miami, FL 33199

## Abstract

Metasearch engines can be used to facilitate ordinary users for retrieving information from multiple local sources (text databases). In a metasearch engine, the contents of each local database is represented by a representative. Each user query is evaluated against the set of representatives of all databases in order to determine the appropriate databases to search. When the number of databases is very large, say in the order of tens of thousands or more, then a traditional metasearch engine may become inefficient as each query needs to be evaluated against too many database representatives. Furthermore, the storage requirement on the site containing the metasearch engine can be very large. In this paper, we propose to use a hierarchy of database representatives to improve the efficiency. We provide an algorithm to search the hierarchy. We show that the retrieval effectiveness of our algorithm is the same as that of evaluating the user query against all database representatives. We also show that our algorithm is efficient. In addition, we propose an alternative way of allocating representatives to sites so that the storage burden on the site containing the metasearch engine is much reduced.

## 1 Introduction

The Internet has become a vast information resource in recent years. To help ordinary users find desired data in this environment, many *search engines* have been created. Each search engine has a *text database* that is defined by the set of documents that can be searched by the search engine. Usually, an inverted file index for all documents in the database is created and stored in the search engine. For each *term* which can represent

a significant word or a combination of several (usually adjacent) significant words, this index can identify the documents that contain the term quickly.

Frequently, the information needed by a user is stored in the databases of multiple search engines. Consider the case when a user wants to find research papers in some subject area. It is likely that the desired papers are scattered in a number of publishers' and/or universities' databases. Substantial effort would be needed for the user to search each database and identify useful papers from the retrieved papers. A solution to this problem is to implement a *metasearch engine* on top of many local search engines. A metasearch engine is just an interface. It does not maintain its own index on documents. However, a sophisticated metasearch engine may maintain information about the contents of each of its underlying search engines (will be called the *representative* of the search engine/database) in order to provide better service. When a metasearch engine receives a user query, it first passes the query to the appropriate local search engines, and then collects (sometimes, reorganizes) the results from its local search engines. With such a solution, only one query is needed from the above user to invoke multiple search engines.

A closer examination of the metasearch approach reveals the following problems. If the number of local search engines in a metasearch engine is large, then it is likely that for a given query, only a small percentage of all search engines may contain sufficiently useful documents to the query. In order to avoid or reduce the possibility of invoking useless search engines for a query, we should first identify those search engines that are most likely to provide useful results to the query and then pass the query to only the identified search engines. Examples of systems that employ this approach include WAIS [7], ALIWEB [9], gGLOSS [5], Savvy-Search [6], ProFusion [4, 3], and D-WISE [22]. The problem of identifying potentially useful databases to search is known as the *database selection problem*. Database selection is done by comparing each query with all database representatives. If a user only wants the  $m$  most similar documents across all local databases, for some positive integer  $m$ , then the  $m$  documents to be retrieved from the identified databases need to be specified and retrieved. This is the *collection fusion problem*.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. CIKM '99 11/99 Kansas City, MO, USA  
© 1999 ACM 1-58113-146-1/99/0010...\$5.00

In this paper, we present an integrated solution to the database selection problem and the collection fusion problem. When the number of databases is not too large, say no more than 100, we employ a two-level hierarchy of database representatives, where the root node represents the representative for the virtual global database that logically contains the documents in all local databases and each leaf node is a representative for a local database. A necessary and sufficient condition to rank the databases optimally is given. Experimental results are provided to show the superior performance of our approach. The experimental results essentially demonstrate that retrieval using our method in distributed environment will yield the same result as if all data were stored in one site. When the number of local databases is large, the two-level hierarchy is not efficient. In this case, we organize the database representatives into a hierarchy of more than 2 levels. We provide an efficient *best-first search* algorithm to search the hierarchy. We show that the search yields exactly the same retrieval results as if all database representatives were placed in a two-level hierarchy. Since the retrieval performance for the two-level hierarchy is very good, this indicates that our best-first search algorithm is guaranteed to yield equally good retrieval effectiveness. We also show that our method is efficient.

The rest of the paper is organized as follows. In Section 2, we present our solution to the database selection problem and the collection fusion problem when the database representatives are organized into a two-level hierarchy. In Section 3, our solution is extended to a hierarchy of more than two levels. The effectiveness and efficiency issues of our algorithm are addressed. We conclude the paper in Section 4. Due to the space limitation, we will not be able to provide the proofs of some results here. The proofs and more details can be found in [21] which is accessible on the Web.

## 1.1 Related Works

In the last several years, a large number of papers on issues related to metasearch engine or distributed collections have been published. In this section, only the research results most closely related to ours and their differences are identified. Please see [13] for a more comprehensive review of other work in this area.

1. The gGLOSS project [5] is similar in the sense that it ranks databases according to some measure. However, there is no necessary and sufficient condition for optimal ranking of databases; there is no precise algorithm for determining which documents from which databases are to be returned. Organizing representatives in a hierarchy and searching them were mentioned in [5], but a precise algorithm for searching was lacking and there was no evidence that searching a hierarchy of more than two levels would yield the same retrieval effectiveness of searching a two-level hierarchy.
2. A necessary and sufficient condition for ranking databases optimally was given in [8]. However, [8]

considered only the databases and queries that are for structured data. In contrast, unstructured text data is considered in this paper. In [1], a theoretical framework was provided for achieving *optimal* results in a distributed environment. However, such an approach depended on parameters which may be difficult to estimate in practice and no experimental results were provided.

3. In [17], experimental results were given to demonstrate that it was possible to retrieve documents in distributed environments with essentially the same effectiveness as if all data were in one site. However, the results depended on the existence of a *versatile collection* where terms related to any given query needed to be extracted and then added to the query before an actual retrieval of documents take place. In the Internet environment where data are highly heterogeneous, it is unclear whether such a versatile collection can in fact be constructed. Even if such a collection can be constructed, the storage penalty could be very high in order to accommodate the heterogeneity. Furthermore, no algorithm for searching a hierarchy of representatives for large number of databases was given.

## 2 Retrieval from a Two-Level Hierarchy

In Section 2.1, the two-level metasearch architecture approach is reviewed. This architecture is suitable for a moderate number of search engines. Using this architecture, a representative which indicates the contents of a database is constructed for each local database and a representative for the global database is also formed. Each user query is compared against these database representatives. The databases are ranked such that databases which are most likely to contain the most similar documents to the user query are placed ahead of other databases. In Section 2.2, we provide a necessary and sufficient condition for ranking the databases optimally. The condition is simple: The databases should be arranged in descending order of the similarity of the most similar document in each database. In Section 2.3, we describe how the similarity of the most similar document in a database can be estimated from the user query, its database representative and the representative for the global database. Having the databases ranked, a strategy for retrieving documents from individual databases, for stopping searching additional databases and for merging the retrieved documents from the searched databases is sketched in Section 2.4. Experimental results are provided in Section 2.5. For ease of presentation, we assume that the similarities of documents to a given query are distinct.

A query in this paper is simply a set of words submitted by a user. It is transformed into a vector of *terms* with *weights* [14], where a term is essentially a content word and the dimension of the vector is the number of all distinct terms. When a term appears in a query, the component of the query vector corresponding to the term, which is the

*term weight*, is positive; if it is absent, the corresponding term weight is zero. The weight of a term usually depends on the number of occurrences of the term in the query (relative to the total number of occurrences of all terms in the query) [14, 20]. This is the term frequency weight. The weight of a term may also depend on the number of documents having the term relative to the total number of documents in the database, which is the inverse document frequency weight. A document is similarly transformed into a vector with weights. The similarity between a query and a document can be measured by the dot product of their respective vectors. Often, the dot product is divided by the product of the norms of the two vectors, where the norm of a vector  $(x_1, x_2, \dots, x_n)$  is  $\sqrt{\sum_{i=1}^n x_i^2}$ . This is to normalize the similarity between 0 and 1. The similarity function with such a normalization is known as the *Cosine* function [14, 20]. For the sake of concreteness, we shall use in this paper the version of the *Cosine* function [16] where the query uses the inverse document frequency weight and the term frequency weight and the document uses the term frequency weight.

## 2.1 Two-Level Architecture of Metasearch

In this architecture, the highest level (the root node) contains the representative for the global database. There is only one additional level. This level contains a representative for each local database. When a user query is submitted, it is processed first by the metasearch engine against all these database representatives. The databases are then ranked. Finally, the metasearch engine invokes a subset of the search engines and co-ordinates the retrieval and merging of documents from individual search engines.

## 2.2 Optimal Ranking of Databases

We assume that a user is interested in retrieving the  $m$  most similar documents to his/her query for a given  $m$ . Databases  $[D_1, D_2, \dots, D_p]$  are ranked optimally if there exists a  $k$  such that the first  $k$  databases  $D_1, D_2, \dots, D_k$  contains the  $m$  most similar documents and each of these databases contains at least one such document. A necessary and sufficient condition to rank databases optimally [19] is as follows.

- Databases  $[D_1, D_2, \dots, D_p]$  are ranked optimally if the similarity of the most similar document in  $D_i$  is larger than that of the most similar document in  $D_j$ , if  $i < j$ .

**Example 2.1** Suppose there are 4 databases. For a query  $q$ , suppose the similarities of the most similar documents in databases  $D_1, D_2, D_3$  and  $D_4$  are 0.8, 0.5, 0.7 and 0.2, respectively. Then the databases should be ranked  $[D_1, D_3, D_2, D_4]$ . ■

This result applies to any similarity function as well as any function which assigns degrees of relevance to documents.

## 2.3 Estimation of the Similarity of the Most Similar Document

We now estimate the similarity of the most similar document in a database using its representative, the representative of the global database and the query. The representative of a database consists of all terms in the database. For each term, two quantities are kept. They are the *maximum normalized weight* and the *average normalized weight* of the term. For term  $t_i$ , they are denoted by  $mnu_i$  and  $awi_i$ , respectively.  $mnu_i$  is  $\max\{d_i/|d|\}$ , where  $|d|$  is the norm (length) of the document  $d = (d_1, \dots, d_i, \dots, d_n)$  and the maximum is over all documents in the database.  $d_i/|d|$  is the normalized weight of the  $i$ th term.  $awi_i$  is simply the average of such values over all documents in the database, including documents not having the term. For the representative of all databases, it also consists of all terms appearing in any of the databases. For each term, only one quantity, namely the document frequency (the number of documents having the term in the global database) of the term, is kept. This generates a weight which is known as the inverse document frequency weight of the term and is denoted by  $idf_i$  for term  $t_i$ .

The similarity of the most similar document of a database  $D$  with respect to a query  $q = (q_1 * idf_1, \dots, q_k * idf_k)$ , where  $q_i$  is the term frequency of the  $i$ -th query term, is estimated by

$$est(D) = \max_{1 \leq i \leq k} \left\{ \sum_{\substack{j=1 \\ j \neq i}}^k q_j * idf_j * aw_j + q_i * idf_i * mnu_i \right\} / |q|$$

The intuition for having this estimate is that the most similar document in a database is likely to have the maximum normalized weight of the  $i$ th query term, for some  $i$ . This yields the second half of the above expression within the braces. For each of the other query terms, the document takes the average normalized value. This yields the first half. Then, the maximum is taken over all  $i$ , since the most similar document may have the maximum normalized weight of any one of the  $k$  query terms. Normalization by the query norm,  $|q|$ , yields a value less than or equal to 1. We shall drop  $|q|$  for ease of presentation. This will not have any impact as the relative similarity values of the most similar documents of the different databases are not changed. The estimate given here is slightly different from our earlier estimate in [19]. It is easy to see that  $est(D)$  can be computed in time linear to the number of terms in the query. The estimate can be further improved by combining two terms which appear in adjacent positions in some previously processed queries. This is used in our experiments described in Section 2.5.

## 2.4 Coordinate the Retrieval from Different Search Engines

Suppose the databases have been ranked in the order  $[D_1, \dots, D_p]$ , we now briefly review an algorithm [18] which

determines (1) the value of  $k$  such that the first  $k$  databases are searched, and (2) which documents from these  $k$  databases should be used to form the list of  $m$  documents to be returned to the user. Suppose the first  $s$  databases have been invoked from the metasearch engine. Each of these search engines returns the similarity of the most similar document to the metasearch engine which computes the minimum of these  $s$  values. Each of the  $s$  search engines returns documents to the metasearch engine whose similarities are greater than or equal to the minimum value. If  $m$  or more documents have been returned, then they are sorted in descending order of similarity and the first  $m$  documents are returned to the user. Otherwise, the next database will be invoked and the process is repeated until  $m$  documents are returned to the user. It has been shown that if the databases have been ranked optimally (with the databases containing the desired documents ahead of other databases) for a given query, then this algorithm will retrieve all the  $m$  most similar documents with respect to the query.

## 2.5 Experimental Results

15 databases are used in our experiments. These databases are formed from articles posed to 52 different newsgroups in the Internet. These articles were collected at Stanford University [5]. Each newsgroup that contains more than 500 articles forms a separate database. Smaller newsgroups are merged to produce larger databases.

There are altogether 6,597 queries submitted by real users. Both the data and the queries were used in the gGLOSS project [5]. From these 6,597 queries, we obtain two subsets of queries. The first subset consists of the first 1,000 queries, each having no more than 6 terms. They will be referred later as *short queries*. The second subset consists all queries having 7 or more terms. There are 363 *long queries*.

The following tables provide a comparison of our method with the high-correlation method in gGLOSS [5], where each number in the table represents the percentage of the  $m$  most similar documents retrieved by a method.

$m$	our method	HC method
5	98.41%	69.73%
10	99.29%	75.64%
20	99.58%	82.87%
30	99.70%	87.68%

(a) Short Queries

$m$	our method	HC method
5	90.22%	61.44%
10	93.58%	67.64%
20	97.09%	75.82%
30	98.54%	81.82%

(b) Long Queries

For short queries, the improvements of our method over the high-correlation method range from 41.1% to 13.7%.

For long queries, the improvements range from 46.8% to 20.4%. Furthermore, for short queries, our method retrieves on the average 98.4% to 99.7% of the most similar documents. Since Internet queries are typically very short (the average number of terms per Internet query is about 2.2), our result indicates that it is possible to retrieve the most similar documents from multiple databases with essentially the same effectiveness as if all documents were stored in one database. The number of databases searched by our algorithm is only 7.5% to 14.0% more than the number of databases that contain the  $m$  most similar documents and the number of documents transmitted is about 11.2% to 24.2% beyond  $m$ .

## 3 Searching a General Hierarchy of Database Representatives

### 3.1 Definition of Representatives of Superdatabases

In the last section, we had a two-level architecture for database representatives. This architecture is suitable for a moderate number of search engines (or databases — we shall use search engine or database interchangeably since logically each search engine retrieves documents from a logical database), say 100 databases. However, when the number of databases is very large, say thousands or tens of thousands, then there will be storage and efficiency problems. First, the amount of storage to contain all database representatives could be enormous. Second, estimating the similarity of the most similar document for each database could be time consuming for large number of databases. For these reasons, we introduce a general hierarchy of database representatives and a search algorithm for such a hierarchy so that the number of estimations can be significantly reduced. We will also address the storage issue.

We first define this hierarchy of database representatives. As in the two-level architecture, the lowest level contains all representatives of individual local databases. Local databases can be logically grouped into *superdatabases*. For example, if superdatabase  $S_1$  contains databases  $D_1, D_2$  and  $D_3$ , then all documents in these databases are logically contained in  $S_1$ . Physically, superdatabases do not exist. The next level of the hierarchy contains representatives of superdatabases formed from local database representatives directly. The representative of a superdatabase will be called a *super-representative*. Each super-representative is constructed from the representatives that are one level below them and is kept physically. The super-representative,  $RS$ , which is constructed from a set,  $R$ , of representatives, consists of two quantities for each term that appears in any of the representatives in  $R$ . The two quantities are the *maximum normalized weight* and the *maximum averaged normalized weight* of the term. Let the maximum normalized weight of term  $t_i$  in the  $j$ -th representative in  $R$  be  $mnw_{ij}$ . Then, the maximum normalized weight of term  $t_i$  in the super-



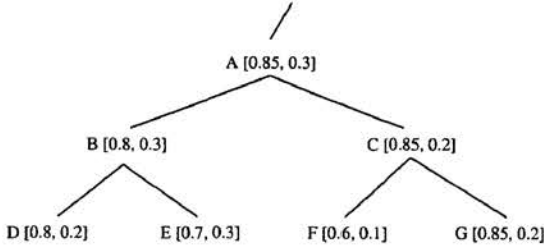


Figure 1: Illustrating the computation of super-representative for one term.

representative  $RS$  is  $\max_{\text{all } j \text{ in } R} mnw_{ij}$ . Similarly, the maximum average normalized weight of term  $t_i$  in the super-representative  $RS$  is obtained by taking the maximum of the corresponding quantities in the component representatives. If the number of super-representatives at a given level is still too large, then they can be grouped into higher level super-representatives by repeating the above process. The root node representative contains the same information as that in the root node of the two-level hierarchy case discussed in Section 2, i.e., the global document frequency of each term which appears in any local database is kept.

Figure 1 illustrates the process of deriving the quantities for a single term in a hierarchy. The first quantity associated with each node is the maximum normalized weight and the second quantity is the maximum average normalized weight.

### 3.2 An Algorithm for Searching the Hierarchy

Let  $RS$  be the super-representative for superdatabase  $S$  that logically contains local databases  $D_1, \dots, D_r$  at the next lower level of the hierarchy (i.e.,  $S$  is the parent of  $D_1, \dots, D_r$ ). The similarity of the most similar document with the query  $q = (q_1 * idf_1, \dots, q_k * idf_k)$  in  $S$  can be estimated using  $RS$  and the root representative as follows. Let the estimate be

$$est(S) = \max_{1 \leq i \leq k} \left\{ mnw'_i * idf_i * q_i + \sum_{\substack{j=1 \\ j \neq i}}^k maw'_j * idf_j * q_j \right\}$$

where  $mnw'_i = \max\{mnw_i(1), \dots, mnw_i(r)\}$ ,  $mnw_i(e)$  is the maximum normalized weight of term  $t_i$  in database  $D_e$ ,  $1 \leq e \leq r$ ;  $maw'_j = \max\{aw_j(1), \dots, aw_j(r)\}$ ,  $aw_j(e)$  is the average normalized weight of term  $t_j$  in the database  $D_e$ ,  $1 \leq e \leq r$ . This estimation formula can also be used to estimate the similarity of the most similar document in a superdatabase whose children are lower level superdatabases. The only change is to replace each  $aw_j(e)$  by  $maw_j(e)$ , where  $maw_j(e)$  is the maximum average normalized weight of term  $t_j$  in the  $e$ th child super-representative.

We now provide a best-first search algorithm to search

this hierarchy of representatives. The main idea of this algorithm is as follows. For a given query, from the children of the root, it selects the representative which yields the largest estimated similarity. If it is the representative of a local database, then the corresponding search engine is invoked and documents are retrieved in the way as described in Section 2.4. If it is a super-representative of an intermediate node in the hierarchy, then the similarity of the most similar document in each of its child representatives is estimated. These child representative nodes are arranged in descending order of similarity and merged with the current list of nodes in descending order of similarity to form a list of representative nodes in the same order. In this list, the estimated similarities can be due to the representatives of local databases or from the non-leaf super-representatives. But, we always take the largest value. If the largest value is from the representative of a local database, then the corresponding search engine is invoked and documents are retrieved according to Section 2.4; otherwise, the best-first search process is executed on the node with the largest estimated similarity value. The details are as follows.

**Search( $m, q, \text{Root}$ )** /\*  $m$  is the number of documents to be retrieved;  $q$  is the query;  $\text{Root}$  is the root node of the hierarchy \*/

1. Initialization:  $min\text{-}sim := 1$ ; /\* the minimum of the similarities of the retrieved documents from previously searched databases is initially set to be 1 (the highest possible similarity) \*/
  2. The similarity of the most similar document in each child of the root node,  $\text{Root}$ , is estimated. These child nodes are arranged in a list  $L$  in descending order of the estimated similarities.
  3. The first node, say  $N$ , is removed from  $L$ .  
If it is the representative of a local database  $D$ , then Step (3.1): {
    - (a) local database  $D$  is searched;
    - (b) the most similar document and its similarity with  $q$ ,  $csim$ , are returned to the result-merger;
    - (c) If  $csim > min\text{-}sim$  then {
      - i. send from database  $D$  all documents with similarity  $\geq min\text{-}sim$  to the result-merger;
      - ii. if  $m$  or more documents have been received by the result-merger, then take the  $m$  most similar documents and stop;
- } else {for each local database  $D'$  which has been searched, do
- i. send all documents from  $D'$  which have similarities  $\geq csim$  (but have not been transmitted) to the result-merger;
  - ii.  $min\text{-}sim := csim$ ;

```

    iii. if  $m$  or more documents have been received
        by the result-merger, then take the  $m$  most
        similar documents and stop;
  } } /* end of step (3.1) */
else { Step (3.2):
  (a) the similarity of the most similar document in
      each child of  $N$  is estimated;
  (b) these child nodes are arranged in a list  $L_1$  in de-
      scending order of estimated similarity;
  (c)  $L$  and  $L_1$  are merged to form  $L$  in descending
      order of estimated similarity;
}
4. repeat Step (3);

```

The *result merger* mentioned in the algorithm is a software component of the metasearch engine which collects the transmitted documents from the searched databases to form the list of  $m$  documents to present to the user.

### 3.3 Effectiveness of Algorithm Search( $m, q, \text{Root}$ )

In Algorithm *Search*( $m, q, \text{Root}$ ), the result-merger gathers the documents retrieved from various local databases until  $m$  or more documents have been received. Step (3.1) is exactly the same process for deciding which documents from the selected databases will form the  $m$  documents to present to the user as that described in Section 2.4. In the two-level hierarchy, the databases (the leaf-nodes) are arranged in a list in descending order of estimated similarity and then Step (3.1) is executed to those databases in the ordered list. In Algorithm *Search*( $m, q, \text{Root}$ ), the leaf-nodes and the intermediate nodes are interleaved in the list of nodes in descending order of estimated similarity. Whenever a leaf-node, i.e., a local database is reached, Step (3.1) is executed to determine the documents to retrieve from the local database. Whenever a non-leaf node is encountered, it is replaced by its children. If the leaf-nodes executed by Step (3.1) are exactly in the same order as if databases were arranged in descending order of estimated similarity, then this algorithm will give the same retrieval performance of a two-level hierarchy as described in Section 2. The following proposition establishes this fact.

**Proposition 3.1** Consider any two local databases  $D_i$  and  $D_j$ . Suppose the estimated similarity of the most similar document in  $D_i$  is higher than that of the most similar document in  $D_j$ , i.e.,  $est(D_i) > est(D_j)$ . Then Step (3.1) of Algorithm *Search*( $m, q, \text{Root}$ ) will execute on database  $D_i$  before it executes on database  $D_j$ .

**Observation:** This proposition guarantees that the databases will be searched in descending order of estimated similarity using our method of estimating similarity of the most similar document in a database or superdatabase. The

same result holds for any other estimation method as long as the estimation method is a non-decreasing function of the two parameters, namely the maximum normalized weight and the maximum average normalized weight (the average normalized weight in the case of a local database).

### 3.4 Efficiency of Algorithm Search( $m, q, \text{Root}$ )

We first give a simple example to illustrate the efficiency of searching a hierarchy in comparison to having a two-level architecture.

Suppose there are 900 disjoint databases. Every 30 databases are grouped into a superdatabase. Then, there will be three levels. The root is at the highest level; the second level has 30 superdatabases (actually their super-representatives); the third level consists of the representatives of the 900 local databases. If the two-level architecture were used, all 900 estimation computations would have been carried out for each query. In the three-level architecture, 30 estimation computations are carried out for the 30 superdatabases. For the superdatabase having the largest estimated similarity, another 30 estimation computations are done for the local databases. Suppose the estimated similarity of a rather highly ranked local database is less than the estimated similarity of the most similar document in another superdatabase, then an additional 30 estimation computations will be carried out. Suppose the  $m$  most similar documents are found from these 60 local databases that are searched, then altogether there are 90 estimation computations. In comparison to the 900 estimation comparisons, one order of magnitude of savings is achieved.

We now establish a result to show that our best-first search algorithm is reasonably efficient. Let  $s$  be the number of databases containing the  $m$  most similar documents. Note that  $s$  is bounded by  $m$ . Let  $h$  be the height of the hierarchy and  $r$  be the number of children of each non-leaf node in the hierarchy. The result is:

**Proposition 3.2** For any single term query  $q$ , the number of estimation computations is no more than  $(s+1)*r*(h-1)$ . (There is no assumption required for this result.)

In the situation that a query has multiple terms, the above result is true if the databases are properly clustered together in the sense that similar databases belong to the same logical superdatabases. For a database or a superdatabase  $S$ , let  $mnw(t \text{ in } S)$  denote the maximum normalized weight of term  $t$  in  $S$ .  $aw(t \text{ in } S)$  and  $maw(t \text{ in } S)$  can be defined similarly.

**Definition 1** A set of databases and superdatabases are suitably clustered in the hierarchy  $H$  with respect to a query  $q$  if for each term  $t$  of the query, the following condition holds. For any given  $m$ , where  $m$  is the number of most similar documents to the query to be retrieved, let  $D_i$  be a database containing at least one of the  $m$  most similar documents and  $C$  be a superdatabase that does not contain any of

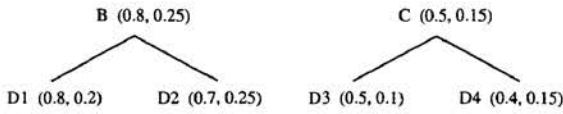


Figure 2: Illustrating Suitably Clustered Databases

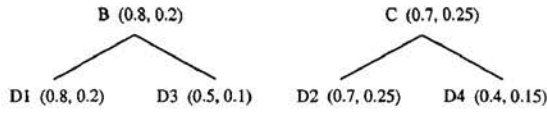


Figure 3: Illustrating Not Suitably Clustered Databases

the  $m$  documents, then  $mnw(t \text{ in } D_i) > mnw(t \text{ in } D)$  and  $aw(t \text{ in } D_i) \geq aw(t \text{ in } D)$  for any database  $D$  contained in  $C$ .

We motivate this definition by the following example. There are two superdatabases in Figure 2. The statistics  $(a, b)$  of a term in two different databases and superdatabases are shown. For example,  $D_1(0.8, 0.2)$  means the maximum normalized term weight and the average normalized term weight of the term in  $D_1$  are respectively 0.8 and 0.2. Suppose in Figure 2, database  $D_1$  contains all of the  $m$  most similar documents, but databases  $D_3$  and  $D_4$  do not contain any of them. The above definition says that if the databases are *suitably clustered* with respect to query  $q$ , then  $mnw(t \text{ in } D_1) > mnw(t \text{ in } D_3)$ ,  $aw(t \text{ in } D_1) \geq aw(t \text{ in } D_3)$ ,  $mnw(t \text{ in } D_1) > mnw(t \text{ in } D_4)$  and  $aw(t \text{ in } D_1) \geq aw(t \text{ in } D_4)$ . The reason why these inequalities are reasonable is that databases with high values of the statistics of certain terms should be within a cluster and similarly databases with low values of the terms should also be within a cluster. In practice, this condition may not be true for all queries but are likely to hold for most queries if the databases are clustered properly. Figure 3 shows a violation of this definition as it is not true that  $aw(D_1) \geq aw(D_2)$ . Here, the high values of the statistics of the term in  $D_2$  are mixed with the low values of the statistics of the term in  $D_4$ , which implies that the databases are not suitably clustered.

**Proposition 3.3** *If the databases are suitably clustered with respect to the query  $q$ , then the number of similarity computations is bounded by  $(s + 1) * (h - 1) * r$ , where  $s$  is the number of databases containing the  $m$  most similar documents to  $q$ ,  $h$  is the height of the hierarchy and  $r$  is the number of children per non-leaf node.*

### 3.5 Placement of Non-Leaf Nodes of the Hierarchy

It was assumed in the above discussion that the hierarchy of representatives are placed completely in the site containing the metasearch engine. This is clearly a reasonable choice. However, there is a potential problem with this approach.

The amount of storage required to store the representatives can be rather large. The problem and its potential remedy can be illustrated by the following example.

Suppose the size of a database representative is 2% of the size of the local database. In the case of a non-root superdatabase  $S$ , let  $S$  be composed from a set of local databases or superdatabases which are one level below  $S$  in the hierarchy. Assume that the size of the representative of  $S$  is on the average twice the size of one of its component representatives. Recall that there are two quantities for each term in a database representative. Thus, the above assumption is based on the premise that the component databases or superdatabases of  $S$  have a lot of terms in common and therefore the number of distinct terms in  $S$  is only twice of the number of distinct terms in a component database or superdatabase. We now use the previous example involving 900 databases to illustrate the storage requirement. Let each database have the same size  $s$ . Then the total size needed to store the 900 database representatives is  $900 * 2\% * s = 18 * s$ . The representatives of the 30 superdatabases take  $30 * 4\% * s = 1.2 * s$ . Thus, the storage requirement is  $19.2 * s$ . In general, if there are  $n$  local databases and each superdatabase is made up of  $r$  databases or superdatabases which are one level below it, with  $r$  reasonably large, say  $r \geq 10$ , then the amount of storage for all representatives and super-representatives are approximately  $n * 2\% * s$  and  $n * 4\% * s/r$ , respectively, as the number of non-leaf nodes which are two or more levels above the leaf nodes are much less than the number of nodes at the bottom two levels. When  $n$  is large, say in the order of thousands, it will impose a rather heavy burden on the site  $S_1$  containing the metasearch engine, in terms of storage and processing power, if all these database representatives are placed in that site. If  $S_1$  is powerful enough, then searching the hierarchy of database representatives can be carried out as described; otherwise, we propose that all super-representatives be placed at the site  $S_1$  and each database representative be placed at the site where the database resides. Thus, the storage requirement at  $S_1$  is approximately  $1.2 * s$  for our example.

We now describe the process to execute a query  $q$  when local representatives are stored at their local database site. The query is submitted to the site of the metasearch engine. The global document frequency of each term in  $q$  is kept in the root representative, as before. These document frequencies together with the query  $q$  are used to estimate the similarity of the most similar document in each superdatabase and executes the best-first algorithm until the parent  $P$  of some local database representatives is reached. At that point, query  $q$  as well as the global document frequencies of the terms in  $q$  will be sent to each child (site) of  $P$ . Each of these child sites then estimates the similarity of the most similar document in its own database and forwards the estimate to the global site  $S_1$ , which will continue to execute the best-first search algorithm  $Search(m, q, Root)$ . In other words, for a given query, the algorithm is executed distributively in multiple sites consisting of the sites  $S_1$  and

some of the sites containing the local databases. The trade-off between this way of storing the database representatives versus the traditional way in which all database representatives are stored in one site is that the former way will incur higher communication cost but less processing cost, since the estimations in local sites can be executed in parallel. Clearly, the former method imposes much less storage burden on the site  $S_1$ , since the representatives of the local databases (which take much more storage than the representatives of the superdatabases) are kept in their own sites.

## 4 Conclusion

We have shown that for a distributed environment involving reasonable number of databases a two-level hierarchy of database representatives together with an appropriate estimation algorithm can yield nearly the same effectiveness as if all documents were in one site. When the number of databases is very large, we propose a hierarchy of representatives with the number of levels  $> 2$ . We provide an algorithm to search the hierarchy. It is shown that the search will produce the same effectiveness as the corresponding two-level hierarchy. We also show that the search of the hierarchy is efficient for single term queries which are submitted frequently in the Internet environment. Under the assumption that databases are clustered properly, we show that the search for multi-term queries is also efficient.

We plan to implement the hierarchy search process for number of levels  $> 2$  in the future and perform experiments with it. We will also incorporate information (such as co-occurrence of terms and linkages between web pages) that affects relevance of documents into the solutions to the database selection problem and the collection fusion problem.

**Acknowledgement:** This research is supported by the following organizations: NSF (IIS-9902872, IIS-9902792, CCR-9816633, CCR-9803974, CDA-9711582, HRD-9707076), NASA (NAGW-4080, NAG5-5095) and ARO (NAAH04-96-1-0049, DAAH04-96-1-0278). We are grateful to L. Gravano and H. Garcia-Molina for providing us with the collection of documents and queries used in our experiments.

## References

- [1] C. Baumgarten. *A Probabilistic Model for Distributed Information Retrieval*. ACM SIGIR Conference, 1997.
- [2] J. Callan, Z. Lu, and W. Bruce Croft. *Searching Distributed Collections with Inference Networks*. ACM SIGIR Conference, 1995.
- [3] Y. Fan, and S. Gauch. *Adaptive Agents for Information Gathering from Multiple, Distributed Information Sources*. 1999 AAAI Symposium on Intelligent Agents in Cyberspace, March 1999.
- [4] S. Gauch, G. Wang, and M. Gomez. *ProFusion: Intelligent Fusion from Multiple, Distributed Search Engines*. Journal of Universal Computer Science, 2(9), 1996.
- [5] L. Gravano, and H. Garcia-Molina. *Generalizing GLOSS to Vector-Space databases and Broker Hierarchies*. VLDB Conference, 1995.
- [6] A. Howe, and D. Dreilinger. *SavvySearch: A Meta-Search Engine that Learns Which Search Engines to Query*. AI Magazine, 18(2), 1997.
- [7] B. Kahle, and A. Medlar. *An Information System for Corporate Users: Wide Area information Servers*. Technical Report TMC199, Thinking Machine Corporation, April 1991.
- [8] T. Kirk, A. Levy, Y. Sagiv, and D. Srivastava. *The Information Manifold*. AAAI Spring Symposium on Information Gathering in Distributed Heterogeneous Environments. 1995.
- [9] M. Koster. *ALIWEB: Archie-Like Indexing in the Web*. Computer Networks and ISDN Systems, 27:2, 1994.
- [10] U. Manber, and P. Bigot. *The Search Broker*. USENIX Symposium on Internet Technologies and Systems (NSITS'97), 1997.
- [11] W. Meng, K-L. Liu, C. Yu, X. Wang, Y. Chang, and N. Rishe. *Determining Text Databases to Search in the Internet*. VLDB Conference, 1998.
- [12] W. Meng, K. Liu, C. Yu, W. Wu, and N. Rishe. *Estimating the Usefulness of Search Engines*. IEEE Data Engineering Conference, 1999.
- [13] W. Meng, C. Yu, and K. Liu. *Challenges and Solutions for Building an Efficient and Effective Metasearch Engine*. Technical Report, Dept. of CS, SUNY at Binghamton, 1999.
- [14] G. Salton and M. McGill. *Introduction to Modern Information Retrieval*. New York: McGraw-Hill, 1983.
- [15] E. Selberg, and O. Etzioni. *The MetaCrawler Architecture for Resource Aggregation on the Web*. IEEE Expert, 1997.
- [16] E. Voorhees, N. Gupta, and B. Johnson-Laird. *Learning Collection Fusion Strategies*. ACM SIGIR Conference, 1995.
- [17] J. Xu, and J. Callan. *Effective Retrieval with Distributed Collections*. ACM SIGIR Conference, 1998.
- [18] C. Yu, K. Liu, W. Wu, W. Meng and N. Rishe. *Finding the Most Similar Documents across Multiple Text Databases*. IEEE Conference on Advances in Digital Libraries, May 1999.
- [19] C. Yu, K. Liu, W. Wu, W. Meng and N. Rishe. *A Methodology to Retrieve Text Documents from Multiple Databases*. Tech. report, U. of Illinois at Chicago, 1999.
- [20] C. Yu, and W. Meng. *Principles of Database Query Processing for Advanced Applications*. Morgan Kaufmann, San Francisco, 1998.
- [21] C. Yu, W. Meng, K. Liu, W. Wu, and N. Rishe. *Efficient and Effective Metasearch for a Large Number of Text Databases*. Tech. report, U. of Illinois at Chicago, 1999. (<http://panda.cs.binghamton.edu/~meng/cikm99.ps>)
- [22] B. Yuwono, and D. Lee. *Server Ranking for Distributed Text Resource Systems on the Internet*. Int'l Conf. On Database Systems For Advanced Applications, 1997.