



PERGAMON

Available at
www.ComputerScienceWeb.com
POWERED BY SCIENCE @ DIRECT®

Information Systems 29 (2004) 697–717



www.elsevier.com/locate/infosys

Efficiency and effectiveness of query processing in cluster-based retrieval[☆]

Fazli Can^{*,1,2}, Ismail Sengör Altingövde, Engin Demir

Computer Engineering Department, Bilkent University, Bilkent, Ankara 06533, Turkey

Received 9 August 2002; accepted 3 June 2003

Abstract

Our research shows that for large databases, without considerable additional storage overhead, cluster-based retrieval (CBR) can compete with the time efficiency and effectiveness of the inverted index-based full search (FS). The proposed CBR method employs a storage structure that blends the cluster membership information into the inverted file posting lists. This approach significantly reduces the cost of similarity calculations for document ranking during query processing and improves efficiency. For example, in terms of in-memory computations, our new approach can reduce query processing time to 39% of FS. The experiments confirm that the approach is scalable and system performance improves with increasing database size. In the experiments, we use the cover coefficient-based clustering methodology (C³M), and the Financial Times database of TREC containing 210 158 documents of size 564 MB defined by 229 748 terms with total of 29 545 234 inverted index elements. This study provides CBR efficiency and effectiveness experiments using the largest corpus in an environment that employs no user interaction or user behavior assumption for clustering.

© 2003 Elsevier Ltd. All rights reserved.

Keywords: Clustering; Cluster-based retrieval; Information retrieval; Performance; Query processing

1. Introduction

The well-known *clustering hypothesis* states that “closely associated documents tend to be relevant

to the same request.” It is this hypothesis that motivates clustering of documents in a database [1]. In the IR research, clustering has been originally introduced with the expectation of increasing the efficiency and effectiveness of the retrieval process [2,3].

In best-match cluster-based retrieval (CBR), it is assumed that there is a flat (one-level) clustering structure. In this environment, the queries are first compared with the clusters, or more accurately with the cluster representatives called centroids. Detailed query by document comparison is performed only within the selected clusters. In hierarchical (multi-level) clustering structures, it

[☆] Recommended by Ricardo Baeza-Yates.

*Corresponding author.

E-mail addresses: canf@muohio.edu (F. Can),
ismaila@cs.bilkent.edu.tr (I.S. Altingövde),
endemir@cs.bilkent.edu.tr (E. Demir).

URL: <http://www.users.muohio.edu/canf/>

¹ Present address: Computer Science and Systems Analysis Department, Miami University, Oxford, OH 45056, USA.

² The majority of this work has been completed when the first author was on sabbatical leave at Bilkent University.

is possible to implement top-down or bottom-up cluster-based search strategies [1,4,5].

As with any such algorithm, the efficiency of CBR is important. In addition to being efficient, CBR should be effective in the sense of meeting user needs. Users may employ the clustering structure in exploring the document space to locate items close to some known documents. This is called *browsing*. Documents within a cluster can also be stored in close proximity to each other within a disk medium to minimize I/O delays [5, pp. 222–227]. However, it is shown that the operational effectiveness of many clustering algorithms is low [6]. It is also observed that the efficiency of CBR is less than the efficiency of full search, FS, i.e., inverted index search of all documents [4,7–9].

In this paper we study the efficiency and effectiveness of query processing in large clustered document collections using the best-match strategy. In the experiments we use the cover coefficient-based clustering methodology, C³M, which we have introduced in our previous work. It is known that C³M has superior performance with respect to other algorithms of the literature [10] and can be used in dynamic environments in an incremental manner for cluster maintenance [11]. The contributions of this study are the following.

- It shows that for large databases it is possible to perform CBR with an efficiency and effectiveness level, which is comparable with that of FS without considerable additional storage overhead. The CBR method proposed in this paper employs a simple and novel storage structure that blends the cluster membership information with the inverted file posting lists. As will be shown this approach significantly reduces the cost of similarity calculations during query processing. Our new storage structure is generic and clearly applicable with other clustering algorithms.
- It provides CBR experiments using the largest corpus in an environment with no user interaction or user behavior assumption. In the experiments we use the Financial Times database of the TREC Disk 4 containing 210 158 documents of text size 564 MB defined by

229 748 terms with total of 29 545 234 inverted index elements. For example, the studies reported in [12,13] provide experiments with larger corpora; however, in their evaluations they assume that the user picks the optimal cluster or try to generate refined clusters with user interaction from an existing global clustering structure. In our work, all decisions are made automatically using similarity measures based on a global clustering structure with no user interaction.

There are various optimization techniques used for inverted index searches [14–17]. They aim to use only the most informative parts of inverted list and try to increase efficiency of query processing without deteriorating retrieval effectiveness. Such techniques can be employed during query processing to further improve the query optimization provided by CBR; however, search optimization (or pruning) with techniques other than clustering is beyond the scope of this study.

One of the advantages of cluster-based IR is that users can *browse* the documents of the best-matching clusters. This may give them the opportunity of seeing additional relevant documents not ranked highly and that may or may not have common term(s) with the submitted query (intuitively documents with no common term with the query but still relevant would have better chance of being a member the best-matching clusters; however, this needs further investigation). Our approach provides efficient access to the clustered documents based on the common terms with the queries and browsing *can* also be supported as an inherent advantage of clustering. On the other hand, the best-matching clusters can be too large to browse and users may prefer to see a ranked document list; in such an environment browsing using a method such as scatter/gather may be practical [12]. This option has not been considered in this study.

The rest of the paper is organized as follows. Section 2 explains our clustering algorithm C³M and the file structures that can be used for the implementation of CBR including our new approach. Section 3 describes the experimental

environment in terms of document database and queries. Section 4 covers the experimental results and their discussion. Section 5 reviews the previous and related work. The conclusions and pointers for future research are given in Section 6. A list of frequently used symbols and acronyms is provided in Table 1.

2. Clustering algorithm and file structures for CBR implementation

In this section we briefly explain our clustering algorithm, the file structures that can be used for CBR, and our new CBR file structure that blends the clustering information with the traditional inverted file.

2.1. Clustering algorithm: C³M

As we indicated earlier, in the experiments we use the C³M algorithm, which is known to have good information retrieval performance. The C³M algorithm assumes that the operational environment is based on the vector space model. Using this model, a document collection can be abstracted by a document matrix, *D*, of size *m* by *n* whose individual entries, *d_{ij}* ($1 < i < m, 1 < j < n$), indicate the number of occurrences of term *j*(*t_j*) in document *i*(*d_i*).

Determining the number of clusters in a collection is a difficult problem [18]. In other clustering

algorithms, if it is required, the number of clusters, *n_c*, is usually a user specified parameter; in C³M it is determined by using the cover-coefficient (CC) concept [10; 19, pp. 376–377]. In C³M, some of the documents are selected as cluster seeds and non-seed documents are assigned to one of the clusters initiated by the seed documents. According to CC, for an *m* by *n* document matrix the value range of *n_c* and the average cluster size (*d_c*) are as follows:

$$1 \leq n_c \leq \min(m, n); \max(1, m/n) \leq d_c \leq m.$$

In C³M, the document matrix *D* is mapped into an *m* by *m* cover-coefficient (*C*) matrix using a double-stage probability experiment. This asymmetric *C* matrix shows the relationships among the documents of a database. Note, however, that the implementation of C³M does not require the complete *C* matrix. The diagonal entries of *C* are used to find the number of clusters, *n_c*, and the selection of cluster seeds. During the construction of clusters, the relationships between a non-seed document (*d_i*) and a seed document (*d_j*) is determined by calculating the *c_{ij}* entry of *C*, where *c_{ij}* indicates the extent with which *d_i* is covered by *d_j*. Therefore, the whole clustering process implies the calculation of $(m + (m - n_c) \times n_c)$ entries of the total *m²* entries of *C*. This is a small fraction of *m²*, since *n_c* \ll *m* (for some database examples please refer to Table 2). A thorough discussion and complexity analysis of C³M are available in [10].

Table 1
Expanded form and meaning of frequently used acronyms and symbols

Acronym	Expanded form	Symbol	Meaning
C ³ M	Cover coefficient-based clustering methodology	<i>d_c</i>	Average no. of documents per cluster
CBR	Cluster-based retrieval	<i>d_s</i>	No. of documents selected during information retrieval
CVDV ^a	Centroid vector document vector	<i>m</i>	No. of documents
CVIIS ^a	Centroid vector inv. index search	<i>n</i>	No. of terms
ICDV ^a	Inverted centroid document vector	<i>n_c</i>	No. of clusters
ICIIS ^a	Inverted centroid inverted index search	<i>n_s</i>	No. of selected best matching clusters
ICsIIS ^a	Inverted centroid skip inv. index search	<i>n_t</i>	Average no. of target clusters per query
FS	Full search	<i>n_{tr}</i>	Average no. of target clusters per query in random clustering
FT	Financial times database (TREC Disk 4)	<i>t_g</i>	Term generality (avg. no. of documents per term)
FTs, FTm	FT small and FT medium size versions	<i>x_d</i>	Depth of indexing (avg. no. of terms per document)

^aCBR implementation method.

Table 2
Characteristics of the FT (TEC Disk 4) and some other databases

Database	m , no. of documents	n , no. of terms	x_d , avg. no. of distinct terms/doc.	n_c , no. of clusters	d_c , avg. No. of docs./clust.
BLISS-1 ^a	152 850	166 216	25.7	6468	25
MARIAN	42 815	59 536	11.2	5218	8
INSPEC	12 684	14 573	32.5	475	27
NPL	11 429	7 491	20.0	359	32
FT	210 158	229 748	140.6	1640	128

^a Approximate n_c value is calculated using the cover-coefficient-based formula: $n_c = n/x_d$.

The CC concept reveals the relationships between indexing and clustering [10]. The CC-based indexing-clustering relationships are formulated as follows:

$$n_c = t/(x_d \times t_g) = (m \times n)/t = m/t_g = n/x_d,$$

and

$$d_c = m/n_c = t_g.$$

In these formulas, the meanings of the variables not used in the text so far are as follows:

t : the total number of non-zero entries in D matrix, t_g : t/n the average number of different documents a term appears (term generality), and x_d : t/m the average number of distinct terms per document (depth of indexing).

These relationships can be used to predict the clustering structure that would be generated by the algorithm.

It is shown that the algorithm can be used in a dynamic environment in an incremental fashion and such an approach saves clustering time and generates a clustering structure comparable to that of cluster regeneration by C³M [11,20].

C³M is a non-overlapping (partitioning) type clustering algorithm. In this paper, we introduce a new version of C³M that creates overlapping clusters to test the effects of overlapping on efficiency and effectiveness of CBR. In the overlapping version a document can be assigned to more than one cluster. For this purpose we slightly modified C³M: Let c_{ij} be the CC value used to cluster d_i , i.e., let us assume that d_i joins to the cluster initiated by d_j since d_j provides the highest coverage of d_i among all cluster seeds. In the

overlapping version we define a *tolerance threshold* h ($1 > h > 0$) and assign d_i also to the cluster initiated by seed document d_k if $c_{ik} > h \times c_{ij}$. Furthermore, non-seed documents can be assigned to (at most) a preset maximum number of clusters. We intuitively set the tolerance threshold to 0.9 and the maximum number of clusters to be assigned, *cluster ceiling* (k), to five clusters. The rationale for our choice of tolerance threshold value is to limit the number of candidate clusters to only those clusters that have reasonably close CC values to the original cluster's CC value (i.e., the CC value observed between the cluster seed that provides the highest coverage for the document to be clustered). For the cluster ceiling, k , we aim to choose an appropriate value so that the index structures would have a reasonable storage size; especially for very large document collections index size can be a concern. Our experimental results (reported in Section 4.1.1) show that as long as the tolerance threshold is high (e.g., 0.9 as in our case), for a document to be clustered only a few cluster seeds can reach the cover coefficient value set by this threshold. In general, h and k values may be determined by some preliminary experiments: cluster overlap performance can be predicted by observing the behavior for a small sample of non-seed documents.

2.2. File structures for the implementation of CBR

2.2.1. Previous file structures for the implementation of CBR

The (best-match) CBR search strategy has two components (a) selection of n_s number of best-matching clusters using centroids, (b) selection

of d_s number of best-matching documents of the selected best-matching clusters. For item (a) we have two file structure possibilities: centroid vectors (CV) and inverted index of centroids (IC). For item (b) we again have two possibilities: document vectors (DV), and inverted index of *all* documents (IIS). One remaining possibility for (b), a separate inverted index for the members of each cluster, is ignored due to its excessive cost in terms of disk accesses (for a query with k number of terms it would involve k disk accesses for each selected cluster). Hence, possible combinations of (a) and (b) define the following CBR implementation policies: CVDV, ICDV, CVIIS, ICIIS.

As summarized in Fig. 1, CVDV means that for cluster match use centroid vectors as they are and for document selection from the selected clusters use the document vectors of the member documents. In ICIIS the documents of the best-matching clusters are selected using the results of FS, which is implemented by IIS. Notice that ICIIS is somewhat counter intuitive to the concept of CBR, since CBR considers only a subset of the database for retrieval purposes, but the IIS component of ICIIS will be performed on the complete database. However, ICIIS still has the potential of being efficient, since query vectors may contain a limited number of terms.

In [7], the efficiency of these methods is measured in terms of CPU time, disk accesses, and storage requirements in a simulated environment defined in [9]. The implementations from best to worst efficiency performance are ordered in the following way: ICIIS, ICDV, CVIIS, CVDV. It is observed that the ICIIS strategy is significantly better than the others. It is also shown that ICIIS is significantly better (5.42 times faster) than a hierarchical cluster search technique, which is

based on a complete link hierarchy [7]. However, this earlier study has further revealed that ICIIS is inferior to FS (1.5 times slower) in terms of efficiency. In this study our aim is to introduce a CBR implementation strategy that would outperform ICIIS and achieve comparable efficiency and effectiveness with FS, and measure its performance in a large document collection.

2.2.2. The new CBR implementation using skips

If we could generate a separate inverted index for the members of individual clusters, this would provide the most efficient computational environment for CBR. However, for a query with k terms if we select n_s best clusters this file structure implies ($k \times n_s$) number of disk accesses, which is large since n_s would be large. To keep both the number of computations and number of disk accesses at its minimum, we have introduced a new CBR implementation structure that we call ICsIIS (IC skip IIS). In this structure IC has its usual structure; however, the IIS component stores not only the traditional posting list information but also the cluster membership information. In this organization, posting list information associated with the members of a cluster are stored next to each other, and this is followed by those of the next cluster's. At the same time we keep a pointer from the beginning of one cluster sub-posting list to the next one. During query processing we use these pointers to skip the clusters, which are not selected as a best-matching cluster. In the literature, another skip idea introduced by Moffat and Zobel is used for efficient decompression of inverted indexes [17]. They compress posting lists by using some fixed length skips, which serve as synchronization points, and are able to decompress posting lists from any point of skips without decompressing the unwanted parts. For example, the (posting) list (1, 5, 10, 13, 18, 23, 50, 57, 58, 60) could be given four synchronization points: 1, 13, 50, and 60. For simplicity let us assume that we are in a conjunctive Boolean query environment, and also assume that another list has already been processed and it is known that the query has no answers between documents 13 and 50, then the original list only needs to be accessed (and decompressed) up to document 13 and after

CVDV	Use Centroid Vectors for cluster selection and Document Vectors for document selection.
ICDV	Use Inverted index of Centroids for cluster selection and Document Vectors for document selection.
CVIIS	Use Centroid Vectors for cluster selection and Inverted Index Search for document selection.
ICIIS	Use Inverted index of Centroids for cluster selection and Inverted Index Search for document selection.

Fig. 1. Summary of possible file structure strategies for CBR implementation (adapted from [7]).

document 50—this means seven posting list positions instead of 10.

An example file structure for our approach is provided in Fig. 2 for a D matrix, which is clustered using C^3M . In this figure each posting list header contains the associated term, the number of posting list elements associated with that term, and the posting list pointer (disk address). The posting list elements are of two types, “cluster number—position of the next cluster,” and “document number—term frequency” for the documents of the corresponding clusters.

Our skip structure is simple yet novel. In the previous CBR research a similar approach has not been used. For example, Salton and McGill’s classical textbook [5, pp. 223–224] defines three cluster search strategies. Two of them are related to hierarchical cluster search and their concern is

the storage organization of the cluster centroids. In the third CBR strategy, documents (not their inverted lists) are stored in cluster order, that is, one access to the “document file” retrieves a cluster of related documents. Our skip idea provides a completely new way of implementing CBR by clustering the individual posting lists elements. This is certainly different than accessing the “documents” in cluster order.

Salton wrote [4, p. 344]:

“In general, the efficiencies of inverted-file search techniques are difficult to match with any other file-search system because the only documents directly handled in the inverted-list approach are those included in certain inverted lists that are known in advance to have at least one term in common with the queries. In a

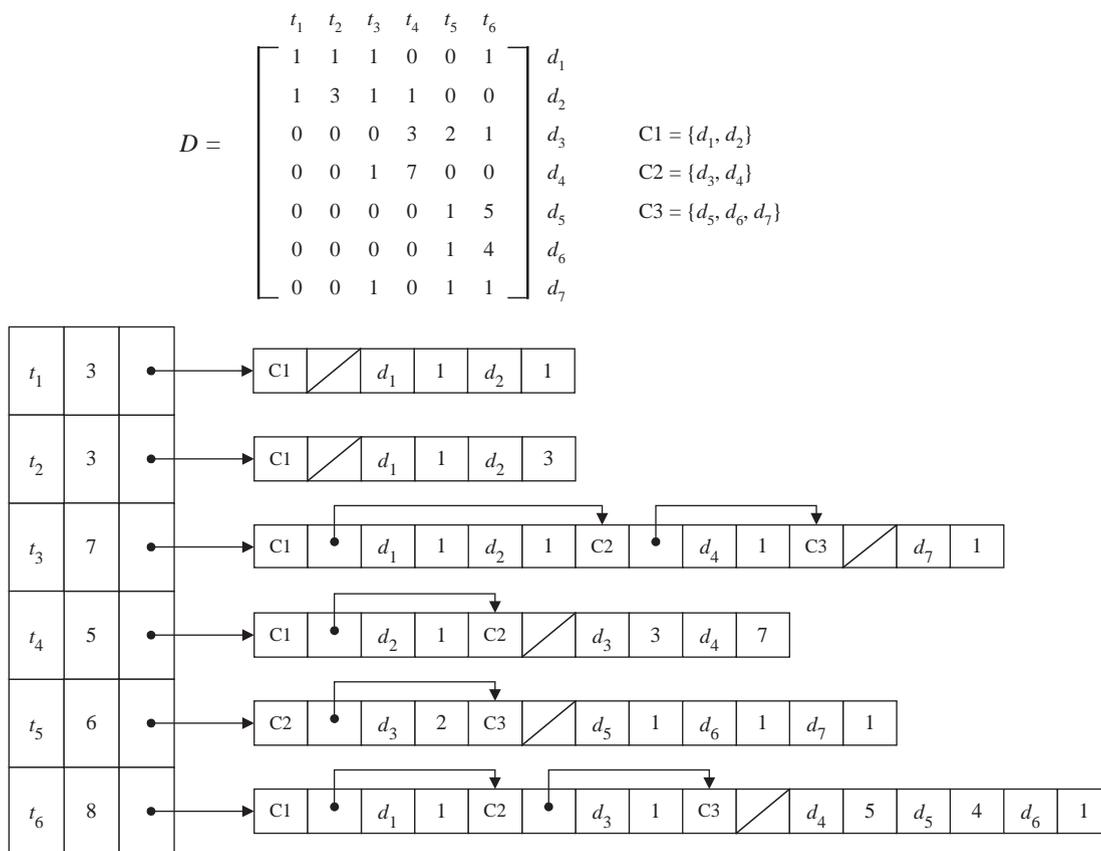


Fig. 2. Example inverted file structure with skips.

clustered organization, on the other hand, many cluster centroids, and ultimately many documents, must be compared with query formulations that may have little in common with the queries.”

The CBR using the skip-based inverted index search technique overcomes the problem stated by Salton, i.e., it prevents matching many unnecessary documents with the queries. For example, in the clustering environment of Fig. 2, if we assume that the user query contains the terms $\{t_3, t_5\}$ and the best-matching clusters for this query are $\{C1, C3\}$, using the ICsIIS approach during query processing after selecting the best-matching clusters we only consider the posting lists associated with t_3 and t_5 . While processing the posting list of t_3 we skip the portion corresponding to $C2$ (since it is not a best-matching cluster). Similarly, while processing the posting list of t_5 , we again skip the unnecessary $C2$ portion of the posting list and only consider the part corresponding to $C3$. In other words, by using the skip approach we only handle the documents that we really need to match with the query.

In the implementation of the skip idea another alternative is to store the cluster number and skip information at the start of the posting lists. Here we adopt the former approach illustrated in Fig. 2. Practically, these two alternatives have no major difference in terms of posting list I/O time, since in almost all cases query term posting lists are read in their entirety because a term usually appears in enough number of different clusters that would require inputting its whole posting list. In query processing, a significant portion of the time cost comes from similarity calculations for ranking, and skipping information helps us in considerably decreasing the cost of these calculations.

3. Experimental environment

3.1. Document database

In the experiments, Financial Times Limited (1991–1994) document collection (referred to as FT database) of TREC Disk 4 is used. The

document database includes 210 158 newspaper articles published between 1991 and 1994. During the indexing stage, we eliminated English stopwords and numbers, indexed the remaining words, and no stemming is used. The resulting lexicon contains 229 748 terms. The D matrix contains 29 545 234 non-zero elements. The average number of distinct terms per document, or depth of indexing x_d , is 140.6, and the longest and shortest documents contain 3220 and *four* distinct terms, respectively. On the average each term appears in 128.6 different documents. This is the average number of distinct documents per term (term generality, t_g).

For easy reference statistical characteristics of the FT collection are provided in Table 2 along with some other databases to give some sense of sizes of the important variables in traditional (INSPEC, NPL), and OPAC (BLISS, MARIAN) [20,21] collections. In this table the number of clusters, n_c , is obtained by using C³M. The numbers show that databases, more specifically their vector spaces, show various degrees of sparsity as indicated by the number of clusters. For example, FT database is quite cohesive and the number of clusters is not that high. On the other hand, OPAC (library), BLISS-1 and MARIAN, vector spaces are sparse and contain relatively large number of clusters, since they cover documents in many different subject areas. The content cohesiveness of a database may be uniformly distributed and clusters may contain approximately the same number of documents or it can be skewed and it may contain a few number of large clusters containing relatively high number of related documents. We will revisit this issue later in Section 4.1 from our database’s point of view.

3.2. Queries and query matching

We used the TREC-7 query topics corresponding to the FT database of TREC Disk 4 collection (queries 351–400) along with their relevance judgments; on the average, there are 38.1 relevant documents per query. In the experiments we used four different types of query sets first two of which are created from the TREC queries.

1. *Qshort* (short queries) created from the title field of the TREC queries, i.e., these are title-only queries.
2. *Qmedium* (medium length queries) created from the title and description fields (combined) of the TREC queries.
3. *Qlong* created from the top retrieved document of each *Qmedium* query. We assume that the relevance judgments of the original query also apply to them.
4. *Qgiant* created by combining a number of random documents from the original data set, and is used for the purpose of evaluating efficiency in its theoretical limits. For this single query we do not measure effectiveness since we have no relevance information for it.

There are 50 queries in each of the query sets *Qshort*, *Qmedium* and *Qlong*, but only one query in the *Qgiant* set. Table 3 provides query sets summary information.

There are several query matching functions that depend on the term weighting used for document and query terms [22]. In this study, the document term weights are assigned using the *term frequency \times inverse document frequency* (IDF) formulation. While computing the weight of term t_j in document d_i , term frequency is computed as the number of occurrences of t_j in d_i , and IDF is $\ln(\text{number of all documents}/\text{number of documents containing } t_j) + 1$. Once the term weights are obtained, document vector is normalized using cosine normalization [22].

The term weights for query terms are calculated in a similar fashion to document term weights. In this case, for computing term frequency component, we use augmented normalized frequency formula defined as $(0.5 + 0.5tf/\max tf)$. Here

$\max tf$ denotes the maximum number of times any term appears in the query vector. IDF component is obtained in exactly the same manner with the document terms. No normalization is done for query terms since it does not affect document ranking.

After obtaining weighted document (d) and query (q) vectors in an n dimensional vector space the query-document matching is performed using the following formula.

$$\text{similarity}(q, d) = \sum_{j=1}^n w_{qj}w_{dj}.$$

The members of the best-matching clusters (note that in CBR a subset of the entire collection is under consideration) are ranked according to their similarity to the query, and for the top 10 (20, 100) documents the effectiveness measures precision and recall are calculated. *Precision* is defined as the ratio of retrieved relevant documents to the number of retrieved documents, and *recall* is defined as the ratio of retrieved relevant documents to total number relevant documents in the collection.

4. Experimental results

In this section, we present various experiments to compare the efficiency and effectiveness of three retrieval strategies: FS, CBR combined with a full inverted index (ICISS), and CBR incorporating the skipping concept (ICsIIS). As stated before, it has been shown that ICISS is more efficient than some other CBR techniques in terms of paging and CPU time, but inferior to FS [7]. In the following set of experiments, we first investigate the validity of C³M clustering for the FT database, and then

Table 3
Query sets summary information (the last three columns indicate no. of terms)

Query set	Source	Average	Min	Max
<i>Qshort</i>	TREC query titles	2.38	1	3
<i>Qmedium</i>	TREC query titles and descriptions	8.16	2	19
<i>Qlong</i>	Top relevant document	190.04	13	612
<i>Qgiant</i>	Random documents	2175.00	2175	2175

examine the effectiveness and efficiency of the three retrieval strategies (namely, FS, ICiIS, and ICsIIS) as we vary several environment parameters. Actually, ICiIS and ICsIIS are the same in terms of their effectiveness since they are two different implementations of the same CBR operation; therefore, for these two we can only compare their efficiency. We also study the scalability of our results. In the rest of the paper, we use CBR interchangeably with ICiIS and ICsIIS when it is appropriate.

The experiments are performed on dual processor Pentium III 866 PC with 1 GB main memory and 20 GB SCSI HDD. The operating system installed on this PC is Windows NT 4.0™. All three retrieval strategies are implemented by using the C programming language. The common data structures and operations are implemented in exactly the same way for all of these strategies, to provide a fair basis of comparison. The source code for the prototype implementation is available at <http://www.cs.bilkent.edu.tr/~ismaila/jrcode.htm>.

4.1. Clustering structure: generation, characteristics and validation

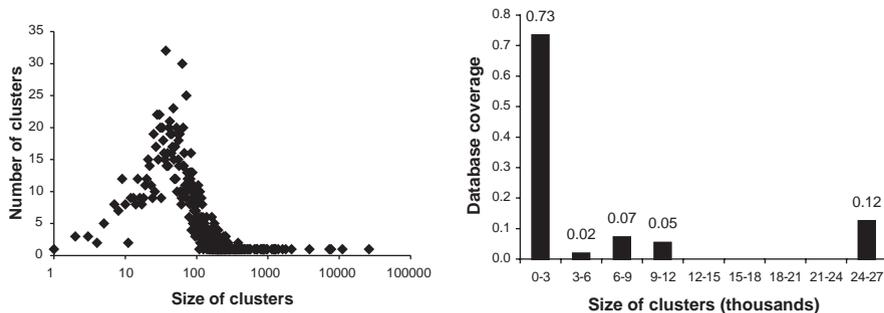
4.1.1. Cluster generation and characteristics of the generated clustering structure

Our experiments yield 1640 clusters (in both non-overlapping and overlapping cases) for the FT collection. In the non-overlapping case the average cluster size is 128 (vs. 176 in overlapping), and the average number of distinct terms in a cluster is

4700 (vs. 5560). Note that in the overlapping case the total number of documents in the clusters is 288 685 (vs. 210 158), which means 37% document duplication.

The generated clustering structure of the non-overlapping case follows the indexing–clustering relationships implied by the CC concept. For example, the indexing–clustering relationships $n_c = (m \times n)/t = m/t_g = n/x_d$, and $d_c = t_g$ are all observed in the experiments (for easy reference the values of these variables are repeated here, $m = 210\ 158$, $n = 229\ 748$, $t = 29\ 545\ 234$, $x_d = 140.6$, $t_g = 128.6$ and the values obtained for n_c and d_c after clustering are 1640 and 128). For example, by substituting the corresponding values (m , n , and t) to the above formula, n_c was implied as 1634 by the relationships, which shows only a 0.4% deviation from the real value obtained by actual clustering. Similarly, the d_c (128) value is almost identical with t_g . As shown in our related previous work [7,10,11] for a given D matrix the clustering structure to be generated by C³M is predictable from the indexing characteristics of a database.

The size distribution of the clusters for the non-overlapping case is presented in Fig. 3. In Fig. 3a the x -axis (in logarithmic scale) shows the cluster size in terms of documents and y -axis shows the number of clusters for the corresponding size. The figure shows that cluster sizes show variety, there are a few large clusters (largest one containing 26 076 documents) and some small clusters, and there are many clusters close to the average cluster size. Fig. 3b shows that majority of the



(a) Cluster distributions in terms of no. of cluster per cluster size (logarithmic scale).

(b) Ratio of total no. of documents observed in various cluster size windows.

Fig. 3. Cluster size distribution information: (a) Cluster distributions in terms of no. of clusters per cluster size (logarithmic scale). (b) Ratio of total no. of documents observed in various cluster size windows.

documents (about 73% of them) are stored in clusters with a size 1–3000. Please note that for only 10% of the queries top 10 results include documents from the largest cluster, which means that our results are not significantly biased by the existence of a large cluster.

4.1.2. Validation of the generated clustering structure

Before using a clustering structure for IR we must show that it is significantly *different* from, or better than, random clustering in terms of reflecting the intrinsic nature of the data. Such a clustering structure is called valid. Two other cluster validity issues, clustering tendency and validity of individual clusters, are beyond the scope of this study [18].

Our cluster validation approach is based on the users' judgment on the relevance of documents to queries and follows the methodology defined in [10]. Given a query, a cluster is said to be a *target cluster* if it contains at least one relevant document to the query. Let n_t denote the *average number of target clusters for a set of queries*. Next, let us preserve the clustering structure and distribute all documents randomly to these clusters. The average number of target clusters for this case is shown by n_{tr} and its value can be calculated without creating random clusters by the modified form [10] of Yao's formula [23]; however, for the validity decision we need the distribution of the n_{tr} values. The case $n_t > n_{tr}$ suggests that the tested clustering structure is invalid, since it is unsuccessful in placing the documents relevant to the same query into a fewer number of clusters than that of the average random case. The case, $n_t < n_{tr}$, is an indication of the validity of the clustering structure; however, to decide validity one must show that n_t is significantly less than n_{tr} .

According to our validity criterion, we must know the probability density function of n_{tr} . For this purpose, we perform a Monte Carlo experiment and randomly distribute the documents to the cluster structure for 1000 times and for each experiment compute the average number of target clusters. The minimum, maximum, and average n_{tr} values are observed as 27.78, 29.02 and 28.41 (see Fig. 4 for the probability density function of the n_{tr}

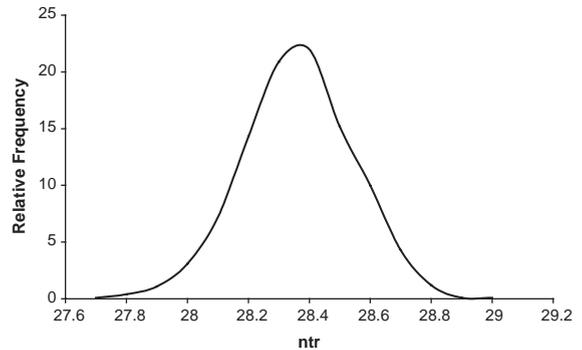


Fig. 4. Histogram of n_{tr} values for the FT database ($n_t = 20.1$).

values). Then, we compute the n_t value, and it is 20.1. Clearly, n_t is significantly different than the random distributions n_{tr} , since it is less than all of the observed random n_{tr} values. These observations show that the clustering structure used in the retrieval experiments is not an artifact of the C³M algorithm, on the contrary, significantly better than random and valid.

4.2. Determining number of best-matching clusters for CBR

The experiments show that selecting more clusters increases effectiveness since as we increase n_s (i.e., the number of selected clusters) more relevant documents would be covered [3, p. 376]. In our previous research, it is observed that effectiveness increases up to a certain n_s value, after this (saturation) point, the retrieval effectiveness remains the same or improves very slowly [10, Fig. 6]. For the INSPEC database, this saturation point is observed when n_s is about 10% of the clusters and during the related experiments about the same percentage of the documents is considered for retrieval. This percentage is typical for (best-match) CBR [3, p. 376].

In our experiments, for a range of n_s values, we retrieved top 10 documents for the query set Q_{medium} and measured the effectiveness in terms of mean average precision (i.e., average of the precision values observed when a relevant document is retrieved) [24, p. 80]. The results depicted in Fig. 5 also confirm the above observation

regarding INSPEC, where the effectiveness increases up to 164 clusters (10% of the cluster number n_c) and then no major change occurs. Therefore, we use 10% of n_c ($n_s = 164$ clusters) as the number of clusters to be used in the retrieval experiments.

In Fig. 6, we report the total number of documents in the clusters for each value of n_s . Both figures show that, for example, if we select the first best-matching 164 clusters (10% of the existing clusters) we need to match 9.09% of the documents with the queries, since this many documents exists in the selected clusters (the numbers are averages for all queries). The

observations show that there is a linear relationship between the percentage of clusters selected and the percentage of the database covered by them.

Determining the centroid terms is also an issue, since they may influence the effectiveness and efficiency of CBR. In this paper, the most frequent terms in clusters are chosen as centroid terms. The weight of a centroid term t_j is computed by $term\ frequency \times IDF$ formula, where term frequency is set to 1 and IDF is $\ln(\text{number of centroids}/\text{number of centroids including the term } j)+1$. In Sections 4.3 and 4.4, we use the ad hoc centroid length value of 250 terms for both overlapping and non-overlapping cases. In Section 4.5, we further investigate the impact of various centroid length and term weighting strategies on the efficiency and effectiveness of query processing.

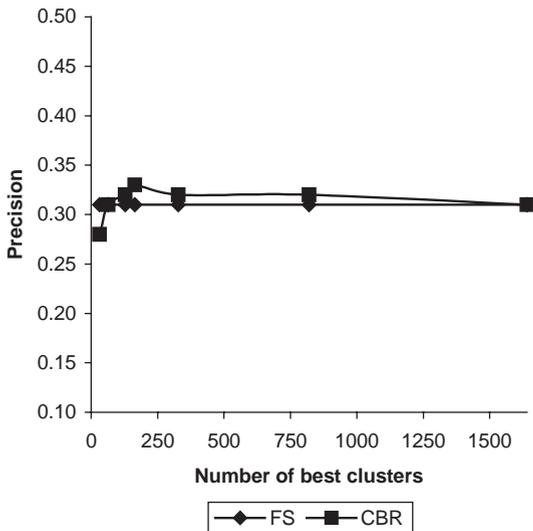


Fig. 5. For $d_s = 10$ and query set *Qmedium*, mean average precision versus n_s .

4.3. Effectiveness experiments

To evaluate the effectiveness of three IR strategies, we retrieved the top 10, 20, and 100 documents for each of the query sets, namely, *Qshort*, *Qmedium* and *Qlong*. The experiments are conducted over both overlapping and non-overlapping clustering structures. The effectiveness results are presented by using both a TREC-like interpolated 11-point precision-recall graph [24, pp. 76–77] and a single mean average precision value (defined in the previous section) for each of the experiments. For the sake of saving space, we provide only top 10 effectiveness results for the

No. (%) of Selected Clusters	Avg. No. (%) of Selected Documents
32 (1.95)	3857.08 (1.84)
64 (3.90)	8608.14 (4.10)
96 (5.85)	12 041.18 (5.73)
128 (7.81)	15 701.94 (7.47)
164 (10.0)	19 107.12 (9.09)
820 (50.0)	102 016.14 (48.54)
1640 (100.00)	210 158.00 (100.00)

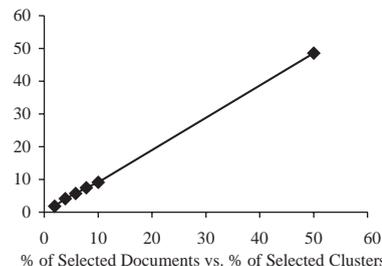


Fig. 6. Relationship between number of selected clusters and number of documents in the selected clusters.

experiments of the non-overlapping and overlapping clustering. For top 20 and 100 documents we have similar results.

Table 4 provides the mean average precision values for the retrieval strategies. For short queries, FS gives the best performance and it is followed by non-overlapping cases. In the case of medium size queries, CBR outperforms FS. For long queries, the reverse is true. For a more detailed comparison consider Fig. 7. They illustrate that the effectiveness of FS and CBR are quite close to each other for different sets of queries with varying lengths. The effectiveness achieved over the overlapping cluster structure can be comparable or sometimes better than non-overlapping CBR and FS. For instance, Table 4 shows that for Q_{medium} , non-overlapping CBR is better than FS, and overlapping CBR is even better than the non-overlapping case.

In Table 5, for the same query sets and top 10 documents, we provide the effectiveness comparisons of individual queries during FS and CBR in non-overlapping case. For instance, CBR achieves better than FS in 6% of the Q_{short} queries. These results further indicate that there is no single best approach for IR, and either one of CBR

Table 4
Mean average precision values for retrieval strategies ($n_s = 164$, $d_s = 10$)

Query set	FS	CBR (non-overlap.)	CBR (overlap.)
Q_{short}	0.307	0.296	0.268
Q_{medium}	0.314	0.326	0.348
Q_{long}	0.383	0.354	0.350

Table 5
Effectiveness comparison of FS and CBR (ICIS and ICsIIS), for non-overlapping clusters

Query Set	CBR = FS (%)	CBR > FS (%)	CBR < FS (%)
Q_{short}	76	6	18
Q_{medium}	70	10	20
Q_{long}	88	4	8

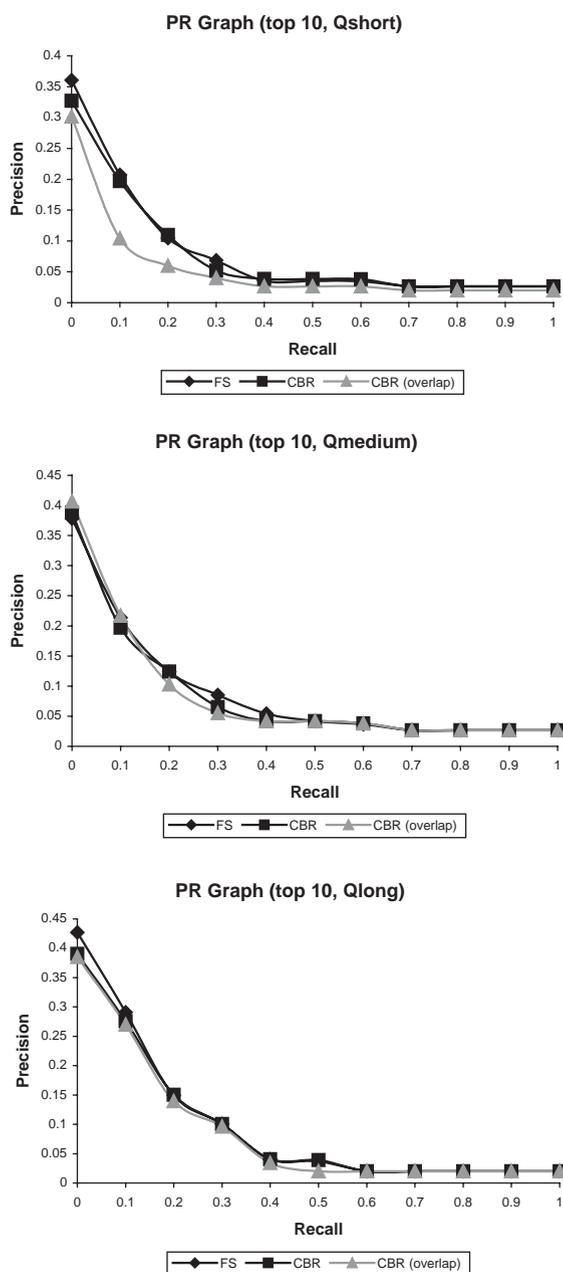


Fig. 7. Interpolated PR graph for all query sets using top 10 ($d_s = 10$) documents.

or FS can perform better for different queries. Note that our CBR approaches that blend inverted indexes with cluster-based retrieval lead to new opportunities for combining the best results of both strategies, in a way that has not been done

Table 6

Average number of document posting list elements processed by each retrieval strategy for each query set and percentage savings provided by ICsIIS

Query set	FS and ICiIS	ICsIIS	% ICsISS savings w.r.t. FS and ICiIS
<i>Qshort</i>	9791	4238	57
<i>Qmedium</i>	49415	16342	67
<i>Qlong</i>	1 813 734	784 005	57
<i>Qgiant</i>	12 398 355	6 637 800	47

before. For example, during query processing we can handle query terms as in FS or CBR like a mixture depending on the query term properties.

4.4. Efficiency experiments

4.4.1. Results in terms of processing requirements

We measure the efficiency of each retrieval strategy for top 10, 20 and 100 documents for all query sets. We evaluate the efficiency by using two different measures: (i) number of processed (accessed) posting list elements, and (ii) actual query processing time. In the following we provide the results for only non-overlapping case. For the overlapping case although the processing requirements are higher (due to the longer posting lists), the relative efficiency performance of compared algorithms does not exhibit a significant difference.

In Table 6, we present the average number of document posting list elements processed for each query set while ranking documents using the query matching function. The posting lists brought to memory for ICsISS are longer than those for FS and ICiIS, as the skipping inverted index elements of ICsIIS include cluster information. On the other hand, in all cases, during the similarity calculations the ICsIIS strategy visits much less posting list elements than IIS and ICiIS, since most of the posting list elements of ICsIIS are skipped due to our (skipping) storage structure. The last column shows the percentage savings provided by ICsIIS with respect to FS and ICiIS in posting list processing (the entries of this column are obtained from those of the second and third columns). The savings in terms of realistic query cases (*Qshort* to *Qlong*) savings range between 57% and 67%.

Table 7

Average in-memory processing time (sec) per query for each retrieval strategy and relative performance of ICsIIS with respect to FS

Query set	FS	ICiIS	ICsIIS	ICsIIS/FS
<i>Qshort</i>	0.051	0.052	0.038	0.75
<i>Qmedium</i>	0.141	0.143	0.055	0.39
<i>Qlong</i>	1.090	1.107	0.442	0.41
<i>Qgiant</i>	4.319	4.385	2.418	0.56

The average in-memory processing time per query is reported in Table 7. The results reveal that the savings indicated in Table 6 are proportionally reflected to the actual execution times. In all cases, ICsIIS performs faster than its competitors as the candidate result set to be considered is significantly reduced by the skipping technique proposed in this paper. Also note that in our database the average posting list length (or term generality t_g) is short, 129 elements. Our heuristics save more time for longer posting lists; therefore, we anticipate that efficiency results would be even better in databases with longer lists. This also explains why the savings for *Qshort* is relatively less than it may be expected.

Please note that our skipping optimization is in-memory, whereas both ICiIS and ICsIIS have an extra cost of disk access for inverted centroid index entries. So, from a theoretical point of view CBR approaches discussed here suffer from this extra I/O cost. However, in practice, we observed that the extra I/O operations associated with accessing inverted centroid index entries are mostly compensated with today's file caching capabilities. In particular, the size of the inverted centroid index is only 1.5% of IIS (see Table 9 of Section 4.4.2) and

Table 8
Storage requirements (size in MB) and posting list (PL) information for inverted index files

Inverted index file	Size	Avg. posting list length (docs/term)	Max. no. of docs./PL	Min. no. of docs./PL
IIS	338	129	93 693	1
Skip ISS (non-overlapping)	426 (26% > FS)	162	95 329	2
Skip IIS (overlapping)	560 (65% > FS)	215	132 328	2

it can be effectively buffered or even totally stored in main memory. For instance, the average overall query processing times (in-memory computations + I/O overhead) for Q_{medium} is measured as 0.265, 0.301 and 0.240 sec per query for FS, ICIIS and ICsIIS, respectively. In this case, all index structures are kept in the disk medium and extra I/O cost is reduced by OS buffering mechanism. Thus, we claim that ICsIIS is a worthwhile retrieval strategy also in terms of efficiency considerations.

4.4.2. Results in terms of storage requirements

As it is mentioned before, the ICsIIS strategy proposed in this paper incorporates cluster membership information into the inverted index posting lists. In Table 8, we present statistics about the inverted index files stored on the disk for FS and the non-overlapping and overlapping clustering cases. It may be seen that the storage requirement for cluster-skipping inverted index is modestly higher than the ordinary inverted index file. The index creation time is 182 min for all structures (i.e., IIS, skip IIS non-overlapping, and skip IIS overlapping, the effect of skips on indexing time is negligible). (Centroid generation time for both non-overlapping and overlapping structures is about 20 min.)

The storage requirements of FS (using IIS) is simply equal to “the total number of elements in the posting lists” times “the size of an elements in the posting list.” A posting list element consists of a 4 byte (integer) document number and 8 byte (double) term weight. There are 29 545 234 inverted index elements where each costs 12 bytes, leading to a total of 338 MB.

In ICsIIS, posting lists include extra elements consisting of cluster number and the skip pointer.

These additional elements also take 12 bytes to conform to the ordinary posting list elements (including document number and term weight). The average number of terms per cluster (*avg. terms/cluster*) is equal to 4700 and 5560, respectively, for the non-overlapping and overlapping clustering cases. This means that in the non-overlapping case the cluster number and the skip pointer address are included in 4700 different posting lists. This makes an additional cost of 88 MB (4700×1640 (total number of clusters) $\times 12$ bytes), i.e., total of 426 MB.

In the ICsIIS overlapping case the extra cost with respect to FS increases, *avg. terms/cluster* is equal to 5560; therefore, the cost due skip information is 104 MB ($5560 \times 1640 \times 12$). Furthermore, in the overlapping case we have 78 527 additional (overlapping) documents (288 625–210 158). For projection purposes if we assume that each overlapping document is an average document containing 141 terms, then the cost these additional documents will introduce to IIS is 127 MB ($78 527 \times 141$ (average number of terms per document) $\times 12$), thus both (104 + 127 MB) make a total of 231 MB. In Table 8, the difference between (actual) overlapping ICsIIS and FS is slightly less than this number and is equal to 222 MB (about 4% less than our projection).

It is possible to decrease the storage cost of inverted file structures by almost 50% by replacing the term weight information (8 bytes) by term document frequency (2 bytes) [17]. If we do that, relative retrieval performances are expected to remain the same and the cost of individual posting list elements drops from 12 (4 + 8) bytes to 6 (4 + 2) bytes.

The detailed disk storage requirements for the important file structures of each strategy are

Table 9
Storage requirements (in MB) for individual components

Storage component	Size	FS	ICISS	ICsISS
IIS (inverted index for docs.)	338	a	a	
IIS with skip information	426			a
IC (centroid length 250)	5		a	a
CM (cluster membership)	3		a	
Wordlist	5	a	a	a
Centroid word list (cent. size 250)	5		a	a

^a Means that corresponding component is required.

shown in Table 9. The last two rows of the table show the storage overhead of storing the indexing terms to find the posting lists of the query terms.

The in-memory requirements of ICsIIS are similar to that of ICiIS. However, ICsIIS does not require cluster membership information to be kept in the memory, since it is blended into the posting lists, whereas ICiSS does. Accordingly, the most demanding internal storage requirement for ICsIIS is for the so-called *accumulator array*, which is used to store the similarity of documents to the processed query. This requirement is clearly the same for all three strategies described in this paper. From these discussions, we can conclude that ICsIIS is feasible in terms of memory and disk storage costs.

4.5. Effects of centroid generation strategies

In the experiments, we investigated the impact of centroid length and centroid term weighting schemes on the effectiveness and efficiency of cluster-based retrieval by using the *Qmedium* case as a representative. All the experiments are performed over clusters generated by both non-overlapping and overlapping versions of the C³M algorithm. We generated four sets of centroids with fixed lengths 250, 500, 1000 and 2500. In an additional experiment, we set the centroid length of each cluster to the 10% of its unique terms. For each of these lengths, we applied three different centroid term weighting schemes: CW1, CW2, and CW3, where the weight of a centroid term is computed by the formula $term\ frequency \times IDF$. In CW1, *term frequency* is taken as 1, in CW2 and

CW3 it is taken as the number of occurrence of the term in the cluster. In CW1 and CW2, *IDF* is taken as $\ln(\text{number of clusters}/\text{number of centroids including the term})+1$, in CW3, it is taken as $\ln(\text{sum of occurrence numbers in the centroids}/\text{number of occurrence in the cluster})+1$. All weights are normalized once they are assigned.

For all of these experiments, the effectiveness of CBR remains almost the same; whereas the efficiency slightly degrades as accessing inverted index elements for centroids requires more time with increasing centroid length. However, in all experiments, ICsIIS still outperforms ICiSS in terms of query processing time. Also, in most of the experiments, ICsIIS achieved comparably well as FS, which is not influenced from the change of centroids.

4.6. Scalability experiments

The scalability of C³M, especially from an incremental clustering point of view, has been thoroughly studied in our previous work [11,20]. In this section we consider the scalability of our skip-based CBR strategy in terms of its efficiency, effectiveness, and storage structures. For obtaining the clusters, we use a naïve implementation of C³M based on ASCII files. In a PC environment, this unrefined implementation clusters the FT database in approximately 114 min.

For the scalability experiments we obtained two smaller versions of the FT database containing approximately one-third and two-thirds of the original collection. We refer to them as FT small (FTs) and FT medium (FTm). The characteristics

Table 10
Characteristics of the FT databases

Database	m , no. of documents	n , no. of terms	x_d , avg. no. of distinct terms/doc.	n_c , no. of clusters	d_c avg. no. of docs./clust.
FTs	69 507	144 080	145.7	955	73
FTm	138 669	191 112	142.1	1319	105
FT	210 158	229 748	140.6	1640	128

of all FT databases are given in Table 10 (for easy reference the original FT database is also repeated in the same table). FTs and FTm, respectively, contain the first 69 507 and 138 669 documents of the original FT database. It may be noted in passing that the indexing-clustering relationships are again observed. For example, the indexing-clustering relationship $n_c = n/x_d$ implies 989 and 1345 clusters, respectively, for the FTs and FTm databases. The difference between actual numbers and projected numbers is less than 4% as in the case of FT (see Section 4.1.1).

In the scalability experiments, as a representative case, we only consider the non-overlapping clustering structure and use the *Qmedium* query set, which is the mid-way in terms of the query sizes we used. In the experiments we retrieve 10% of the clusters ($n_s = 0.1n_c$), examine the top 10 documents ($d_s = 10$) for performance measurement, and use centroids with 250 terms as in the previous experiments.

4.6.1. Scalability of effectiveness

The experimental results in terms of single mean average precision value are reported here. Table 11 shows that when we use the small database, FTs, the CBR effectiveness is about 16% lower than that of FS. In the case of FTm the performance of CBR in terms of effectiveness improves and it lags behind FS by only 1%. Finally, with the full and largest database, CBR outperforms FS by 4%. These observations confirm that our CBR methodology scales well with the database size and has the tendency of showing slightly better performance than that of FS with larger databases. This improvement of CBR effectiveness can be attributed to the refinement of cluster structures with increasing database size.

Table 11

Mean average precision values for retrieval strategies FS and CBR with different databases ($n_s = 10\%$ of n_c , $d_s = 10$) and performance of CBR w.r.t. FS

Database	FS	CBR	CBR/FS
FTs	0.285	0.238	0.84
FTm	0.335	0.332	0.99
FT	0.314	0.326	1.04

4.6.2. Scalability of efficiency

Table 12 provides the average number of posting list elements processed for a query with each database. The values of the ICsIIS column show that this approach is much more efficient than the other two approaches. The last column of the table shows the savings provided by ICsIIS with respect to FS and ICiIS in terms of the posting list elements processed. For FTs, it provides 46% savings and savings increase with the increase of the database size and finally for FT the savings provided by ICsIIS are a substantial 67%. As shown in Table 13, these savings translate themselves to in-memory processing time savings. The last column of the table shows that the efficiency of ICsIIS becomes more prevalent with increasing database size. This again shows that ICsIIS scales well with increasing database sizes. Note that comparable efficiency results for overall query processing times (with I/O) are also observed in our experimental environment.

4.6.3. Scalability of storage and indexing structures

For the experimental databases FTs, FTm, and FT the requirements of the individual storage components are shown in Table 14. As the

Table 12

Average number of “document” posting list elements processed by each retrieval strategy for each database and percentage savings provided by ICsIIS

Database	FS and ICIS	ICsIIS	% ICsISS savings w.r.t. FS and ICIS
FTs	16 875	9214	46
FTm	32 916	14 161	57
FT	49 415	16 342	67

Table 13

Average in-memory processing time (sec) per query for each retrieval strategy with each database and relative performance of ICsIIS with respect to FS

Database	FS	ICISS	ICsIIS	ICsIIS/FS
FTs	0.043	0.044	0.022	0.51
FTm	0.092	0.091	0.045	0.49
FT	0.141	0.143	0.055	0.39

numbers show, the overhead of the secondary storage structures (namely inverted index for centroids (IC), cluster membership information (CM), and the word lists used to find the posting lists associated with the query terms, the last two rows) is negligible. For example, the storage cost of IC with respect to IIS is 2.6%, 1.8%, and 1.5% for the databases FTs, FTm, and FT, respectively. As the size of the database increases, the relative cost of IC decreases, since the rate of increase in number of clusters is lower than that of documents. The size of the IC storage structure also indicates that in query processing the cost of selecting the best-matching clusters is a small fraction of the query processing time. In terms of storage requirements, numbers are basically proportional to the sizes of the document vectors used for the creation of the index structures. As we

increase the size of the database, the cost of skip-based IIS slightly decreases (from 0.30 to 0.26) with respect to IIS. This is again due to the fact that the rate of increase in number of clusters is smaller than that of documents.

In terms of centroid (and IIS) generation, we have the following time observations respectively for FTs, FTm, and FT: 5(46), 11(109), and 20(182) min. The time requirements of generating IIS and skip-based IIS are almost the same. Since these are one-time costs and our concern was the scalability of efficiency and effectiveness, we did not try to optimize our implementations for the generation of these storage structures.

4.7. Discussion of results

From the experiments, we draw the following conclusions:

1. In the non-overlapping clustering experiments, all three retrieval strategies of FS, ICISS and ICsIIS achieve similar effectiveness values. In the efficiency experiments, the ICsIIS strategy incorporated with a skip-based inverted index outperforms the other strategies in terms of in-memory operations and performs comparably

Table 14

Storage requirements (in MB) for individual components

Storage component	FTs	FTm	FT
IIS (inverted index for docs.)	116	226	338
IIS with skip information (extra overhead w.r.t. IIS)	151 (0.30)	286 (0.27)	426 (0.26)
IC (centroid length 250)	3	4	5
CM (cluster membership)	1	2	3
Wordlist	3	4	5
Centroid word list (cent. size 250)	3	4	5

- well in terms of overall query processing times, i.e., including I/O, with effective OS file caching for centroid index.
2. In the overlapping clustering experiments; the effectiveness values of ICISS and ICsIIS are slightly improved in particular experiments, but the efficiency results are not as good as the non-overlapping case due to the increasing access costs for both CBR strategies.
 3. The results are independent of the centroid lengths and weighting schemes, as the variations over these parameters do not significantly affect the presented results.
 4. Storage requirements in the disk and memory for ICISS and ICsIIS are moderately higher than FS, and current compression techniques may further reduce these requirements. In ICsIIS, such a reduction has the potential of further improving the processing time, since by using our skipping approach the decompression time can be reduced significantly.
 5. The experiments show that our results are scalable: Effectiveness of CBR slightly increases and efficiency of ICsIIS can improve significantly with the increasing database sizes.

5. Previous and related work

A good survey of clustering in information retrieval is provided in [25]. This work comes with an impressive reference list. The books by Salton [3,4], Salton and McGill [5] and van Rijsbergen [1] also cover previous work on clustering in information retrieval. A new survey of clustering in various application areas can be found in [26]. A good discussion of algorithms for clustering data and cluster validation approaches is available in a beautiful concise book by Jain and Dubes [18].

Our previous work on C³M includes its concepts and effectiveness analysis [10], and how it works in dynamic databases [11,20]. The CBR effectiveness in terms of precision for the INSPEC database is reported in [10]. The study shows that C³M is 15.1–63.5 (with an average of 47.5)% better than four other clustering algorithms [27] in CBR. The same study also shows that the IR effectiveness of the algorithm is comparable with a demanding (in

terms of CPU time and main memory) complete link clustering method that is known to have good retrieval performance [8,9]. The experiments also show that the CBR using C³M is slightly less effective (1.0–6.9%) than FS. The experimental observations reported in [20] show that the incremental version of C³M is cost effective and can be used for many increments of various sizes.

C³M and its concepts have also attracted the attention of other researchers in various application areas, such as chemical information systems [28,29], clustering tendency testing [30], automatic hypertext structure generation [31], and search output clustering [32].

Most clustering research in IR is related to cluster search effectiveness [6,25,33–35]. The research on efficiency aspects of cluster searches is limited. For example, the works presented in [7,9] considers storage, CPU, and I/O efficiency in the same simulated environment. The Salton–McGill book [5] approaches to the efficiency problem in terms of page faults during information retrieval.

The studies reported in [12,13] provide experiments using databases larger than our collection; however, in their evaluations they assume that the user picks the optimal cluster or try to generate refined cluster via scatter/gather browsing paradigm based on an existing global clustering structure. In contrast to these approaches, in our work all decisions are made automatically using similarity measures based on a pre-existing clustering structure with no user interaction or any other assumption.

In FS only query term posting lists are accessed from the disk medium. As a result, the efficiency of FS decreases with increasing query length, since for each query term another posting list must be processed. It is possible to employ a partial evaluation (or pruning) strategy that skips some of the query terms to improve search efficiency with similar search effectiveness [14,16,17,36]. However, as it is stated before such an approach is beyond the scope of this study and inverted index search optimization in CBR (i.e., in ICISS and ICsIIS) is an interesting research possibility by itself, which can be further incorporated to our work.

6. Conclusions and future work

Our CBR implementation method employs a storage structure that blends the cluster membership information with the inverted file posting lists using the concept of *skips*. In the skip approach, posting lists contain the cluster membership information in addition to traditional term weighting information. During CBR, skip pointers embedded in posting lists provide the information to skip unnecessary (non-best matching) cluster members. The indexing structure of the skip approach can be used both for FS and CBR. Our skip-based CBR significantly improves the efficiency of query processing and this improvement is especially due to in memory similarity calculations. As Web search engines often need to traverse very long posting lists in memory, our skip-based CBR would improve the efficiency of Web search engines that may employ clustering. Our results are significant in the sense that the efficiency and effectiveness of CBR have been analyzed at this level for the first time for an existing global clustering structure (note that the clustering structure is static at the time of query processing; however it can be updated in an incremental manner at other times [11,20]).

We show that for large databases, CBR can achieve a time efficiency and effectiveness comparable with FS. The storage requirement for CBR is modestly higher than the ordinary inverted index file. The experiments show that our results are scalable: Effectiveness of CBR slightly increases and efficiency of ICsIIS can improve significantly with increasing database sizes.

There are several promising future research directions:

1. In the experiments it is observed that CBR and FS do not always return exactly the same set of relevant documents even when they achieve the same precision levels; therefore, our results are also important in terms of data fusion or mixing the results of FS and CBR [34,36]. Our skip-based storage structure is especially suitable for an unusual fusion method, which is a hybrid of FS and CBR. For instance, for *important* query terms with relatively high weights, we may *turn-off* skipping, to retrieve some of the documents that are not in the best-matching clusters but still qualify to be in the top 10 (20) documents. We are currently studying other possible heuristics that may allow combining the best possible results from FS and CBR, with the least additional overhead.
2. It would be interesting to study the update of the skip-based inverted index structures in a dynamic retrieval environment with new and deleted old documents.
3. Compression of the posting lists and its effect on the system efficiency both in terms of retrieval time and disk space is another promising research direction. There is every reason to expect that compression will have positive effects on performance, since with compression a similar skip approach gives good results [17].
4. Another research direction is definition of document vectors with different levels of indexing exhaustivity [33] or by latent semantic indexing (LSI) and measuring the system performance [35,36].
5. Indexing of documents at a lower level, such as paragraphs or sentences, looks promising from CBR's point of view. Since in such an environment FS inverted indexes could be extremely long, our optimization with skip concept combined with CBR may provide an important efficiency leap during query processing.
6. For the calculation of the similarity values instead of an *accumulator array* dynamic data structures can be used for memory efficiency [37]. A partial query evaluation or pruning strategy and its effectiveness and efficiency should also be investigated [15–17]. Our scalability experiments (Table 13) and the results reported in [17] imply that when our ICsIIS approach is combined with the *restricted accumulators* (*quit* and *continue* methods of Moffat and Zobel), it can further improve the efficiency performance.
7. For efficiency one may also consider a *lazy evaluation* approach that delays computations until they are requested. For example, the method introduced by Buckley and Lewit [14] calculates the approximate similarity values for a certain number of top matching documents.

It is conceptually a “lazy evaluation” method, since it does the computations if they are needed to satisfy the user requests (in this case certain number of top matching documents). Such an approach could be adapted to our framework to improve the system efficiency without lowering the system effectiveness.

Acknowledgements

We appreciate the comments made by a referee; they help us improve the presentation. We are grateful to Berkant Barla Cambazoglu of Bilkent University for always making himself available for numerous valuable discussions; his pointers have greatly improved the paper. We are thankful to our colleagues Ediz Saykol, Nazlı İkizler, Eray Özkural, M. Mustafa Özdal, and Atacan Conduroglu of Bilkent University for their support of the project. We greatly appreciate the TREC 4 CD and other information made available by NIST (<http://trec.nist.gov/>).

References

- [1] C.J. van Rijsbergen, *Information Retrieval*, 2nd Edition, Butterworths, London, 1979.
- [2] N. Jardine, C.J. van Rijsbergen, The use of hierarchical clustering in information retrieval, *Inform. Storage Retrieval* 7 (1971) 217–240.
- [3] G. Salton, *Dynamic Information and Library Processing*, Prentice-Hall, Englewood Cliffs, NJ, 1975.
- [4] G. Salton, *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer*, Addison Wesley, Reading, MA, 1989.
- [5] G. Salton, M.J. McGill, *Introduction to Modern Information Retrieval*, McGraw-Hill, New York, NY, 1983.
- [6] W.M. Shaw Jr, R. Burgin, P. Howell, Performance standards and evaluations in IR test collections: cluster-based retrieval models, *Inform. Process. Manage.* 33 (1997) 1–14.
- [7] F. Can, On the efficiency of best-match cluster searches, *Inform. Process. Manage.* 30 (1994) 343–361.
- [8] E.M. Voorhees, The effectiveness and efficiency of agglomerative hierarchical clustering in document retrieval, Ph.D. Thesis, Cornell University, Ithaca, NY, 1986.
- [9] E.M. Voorhees, The efficiency of inverted index and cluster searches, in: *Proceedings of the Ninth Annual International ACM-SIGIR Conference*, ACM, New York, 1986, pp. 164–174.
- [10] F. Can, E.A. Ozkarahan, Concepts and effectiveness of the cover-coefficient-based clustering methodology for text databases, *ACM Trans. Database Systems* 15 (1990) 483–517.
- [11] F. Can, Incremental clustering for dynamic information processing, *ACM Trans. Inform. Systems* 11 (1993) 143–164.
- [12] M.A. Hearst, J.O. Pedersen, Reexamining the cluster hypothesis: scatter/gather on retrieval results, in: *Proceedings of the 19th ACM-SIGIR Conference*, ACM Press, New York, 1996, pp. 74–81.
- [13] C. Silverstein, J.O. Pedersen, Almost-constant-time clustering of arbitrary corpus subsets, in: *Proceedings of the 20th ACM-SIGIR Conference*, ACM Press, New York, 1997, pp. 60–66.
- [14] C. Buckley, A.F. Lewit, Optimization of inverted vector searches, in: *Proceedings of the Eighth Annual International ACM-SIGIR Conference*, ACM, New York, 1985, pp. 97–110.
- [15] M. Persin, Document filtering for fast ranking, in: *Proceedings of the 17th ACM-SIGIR Conference*, ACM Press, New York, 1994, pp. 339–348.
- [16] E.W. Brown, Fast evaluation of structured queries for information retrieval, in: *Proceedings of the 18th Annual International ACM-SIGIR Conference*, ACM, New York, 1995, pp. 30–38.
- [17] A. Moffat, J. Zobel, Self-indexing inverted files for fast text retrieval, *ACM Trans. Inform. Systems* 14 (1996) 349–379.
- [18] A.K. Jain, R.C. Dubes, *Algorithms for Clustering Data*, Prentice-Hall, Upper Saddle River, NJ, 1988.
- [19] C.T. Yu, W. Meng, *Principles of Database Query Processing for Advanced Applications*, Morgan Kaufmann, San Francisco, CA, 1998.
- [20] F. Can, E.A. Fox, C.D. Snavelly, R.K. France, Incremental clustering for very large document databases: initial Marian experience, *Inform. Sci.* 84 (1995) 101–114.
- [21] S. Kocberber, F. Can, J.M. Patton, Optimization of signature file parameters for databases with varying record lengths, *Comput. J.* 42 (1999) 11–23.
- [22] G. Salton, C. Buckley, Term-weighting approaches in automatic text retrieval, *Inform. Process. Manage.* 24 (1988) 513–523.
- [23] S.B. Yao, Approximating block accesses in database organizations, *Comm. ACM* 20 (1977) 260–261.
- [24] R. Baeza-Yates, R. Ribeiro-Neto, *Modern Information Retrieval*, Addison Wesley, Reading, MA, 1999.
- [25] P. Willett, Recent trends in hierarchical document clustering: a critical review, *Inform. Process. Manage.* 24 (1988) 577–597.
- [26] A.K. Jain, M.N. Murty, P.J. Flynn, Clustering algorithms: a review, *ACM Comput. Survey* 31 (1999) 264–323.
- [27] P. Willett, A. El-Hamdouchi, Comparison of hierarchic agglomerative clustering methods for document retrieval, *Comput. J.* 32 (3) (1989) 220–227.
- [28] F. Can, E.A. Ozkarahan, Two partitioning type clustering algorithms, *J. Am. Soc. Inform. Sci.* 35 (1984) 268–276.

- [29] P. Willett, V. Winterman, D. Bawden, Implementation of non-hierarchical cluster-analysis methods in chemical information-systems—selection of compounds for biological testing and clustering of substructure search output, *J. Chem. Inform. Comput. Sci.* 26 (1986) 109–118.
- [30] A. El-Hamdouchi, P. Willett, Techniques for the measurement of clustering tendency in document retrieval systems, *J. Inform. Sci.* 13 (1987) 361–365.
- [31] R.B. Kellogg, M. Subhas, Text to hypertext: can clustering solve the problem in digital libraries, in: *Proceedings of the First ACM International conference on Digital Libraries*, ACM, New York, 1996, pp. 144–150.
- [32] Y. Kural, S. Robertson, S. Jones, Deciphering cluster representations, *Inform. Process. Manage.* 37 (2001) 593–601.
- [33] R. Burgin, The retrieval effectiveness of five clustering algorithm as a function of indexing exhaustivity, *J. Am. Soc. Inform. Sci.* 40 (1995) 562–572.
- [34] A. Griffiths, H.C. Luckhurst, P. Willet, Using interdocument similarity information in document retrieval systems, *J. Am. Soc. Inform. Sci.* 37 (1986) 3–11.
- [35] H. Schütze, C. Silverstein, Projections for efficient document clustering, in: *Proceedings of the 20th ACM-SIGIR Conference*, ACM Press, New York, 1997, pp. 74–81.
- [36] J.H. Lee, Analyses of multiple evidence combination, in: *Proceedings of the 20th ACM-SIGIR Conference*, ACM Press, New York, 1997, pp. 267–276.
- [37] I.H. Witten, A. Moffat, T.C. Bell, *Managing Gigabytes Compressing and Indexing Documents and Images*, 2nd Edition, Morgan Kaufmann, San Francisco, CA, 1999.