

Dynamic Inverted Indexes for a Distributed Full-Text Retrieval System

Charles L. A. Clarke* Gordon V. Cormack†

MultiText Project
Dept. of Computer Science
University of Waterloo, Waterloo, Ontario N2L 3G1, Canada

MultiText Project Technical Report MT-95-01
February 24, 1995

Abstract

We describe data structures and an update strategy for the implementation of dynamic inverted indexes in the context of a dedicated index engine for a distributed full-text retrieval system. Except in rare cases, retrieval operations require a single disk access per query term. The on-line update strategy guarantees the consistency of on-disk data structures across node failures. Index compression integrates smoothly. We examine the performance of the system both experimentally and through an analytical comparison with a competing B-tree based approach.

1 Introduction

1.1 Environment

Our general concern is the construction of a distributed full-text retrieval system. The architecture consists of a group of LAN-connected processors, each managing its own separate disk and memory. Individual processors act as either *text servers*, storing documents and servicing requests for portions of these documents, or as *index engines*, identifying the portions of documents that match client-generated search criteria. To external clients, the group of machines appears to be a single large retrieval system. A

front-end processor, the *Marshaller/Dispatcher*, coordinates the activities of the group of processors, interacting with client applications, dispatching queries to the index engines and text servers, marshalling query results and returning the results to clients. Figure 1 provides a schematic overview of the architecture.

Design aspects of distributed full-text retrieval systems have been the subject of earlier research. Burkowski [2] studied the division of processors into text servers and index engines, with the conclusion that such a split can implement a system which provides better overall response time than a system in which each processor acts as both text server and index engine. Tomasic and Garcia-Molina [16] examine the allocation of index terms from a set of documents to index engines. They conclude that all indexing of terms for an individual document is best allocated to a single processor.

1.2 Inverted Indexes

Our specific concern is the data structure design and update strategy used by the index engines. The basic data abstraction implemented by an index engine is an inverted index [9]. File structures based on inverted indexes are standard for implementing retrieval systems [8, 14, 18, 20]. An inverted index is a function that maps index terms into positions in documents where the terms occur. Index terms are typically words, but may include document markup

*email: claclark@plg.uwaterloo.ca

†email: gvcormac@plg.uwaterloo.ca

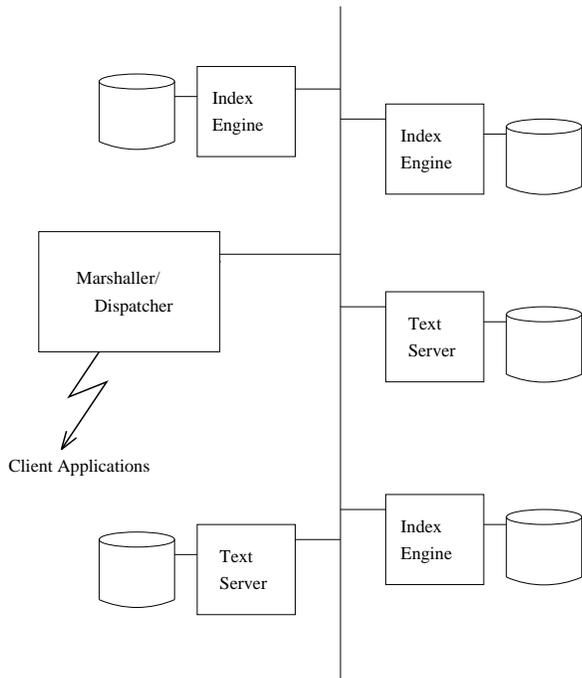


Figure 1: Architecture of the retrieval system.

tags and other structural information of importance to database clients.

Figure 2 presents a simple and widely used realization of the inverted index data abstraction. The *dictionary* maps terms into a pair of offsets into the *postings file*. Between these start and end offsets in the postings file is a sorted list of *postings*, positions within the database where the term occurs. Both the postings file and the dictionary are large enough to require disk storage. Using this realization, a mapping of a given index term into its postings list consists of a binary search of the dictionary (requiring $O(\log w)$ disk accesses, where w is the number of index terms in the dictionary) followed by a single access into the postings file.

1.3 Dynamic Inverted Indexes

Text retrieval systems are primarily viewed as archival. This view leads either to the assumption that the text collection is static or to an append-only (or append-dominated) model of update. If the retrieval system maintains an evolving collection — the current documentation for a large project, for exam-

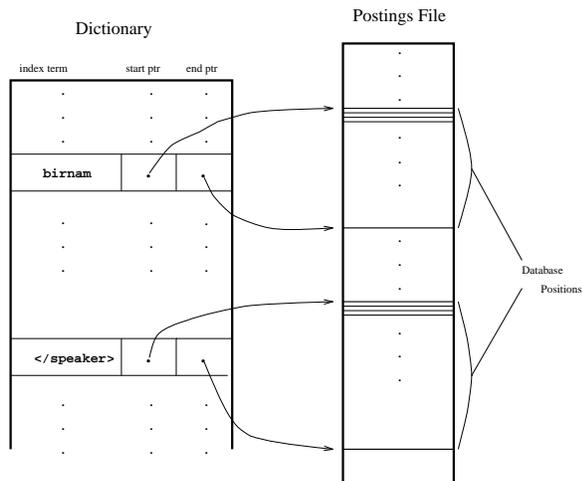


Figure 2: Simple realization of an inverted index.

ple — this view will not be valid. In some instances, we may wish to maintain the system in equilibrium. In a system that stores current netnews, deletions of articles will be as frequent as the addition of articles.

In the context of a distributed text database, an append-only model is not valid for individual index engines even if the system as a whole is archival. As hardware resources are added or removed from the system, or when portions of the collection become “hotspots” of interest, it will be necessary to reorganize the system for load-balancing purposes.

1.4 Design Issues

In an operational environment there are a number of practical issues to be considered when implementing inverted lists.

Retrieval Response Retrieval operations far outnumber update operations. Querying the retrieval system is the primary operation used by external clients and response time is of utmost importance. The mapping of an index term into its postings list must require as few disk accesses as possible. Ideally, a single disk access would be sufficient to translate any term, independent of the size of the dictionary and postings file, and independent of the size of the postings list for the particular word.

Update Throughput Since update is primarily a maintenance function rather than an external client service, update throughput, not response time, is of importance. The simple file structures of figure 2 require a complete rebuild to apply updates.

Index Compression Compression will increase the amount of dictionary and postings data that can be stored on available disk [1, 6, 20, 21]. Since compression and decompression techniques operate by linearly processing a range of data, this property creates a potential decrease in retrieval response time if random access into the data is thereby restricted.

Consistency The database must be maintained in a consistent state at all times. For example, if a failure occurs during update of the postings file, the dictionary must not be left pointing to an incorrect range of postings.

1.5 Existing approaches

Most discussions of inverted indexes for text retrieval assume that the file structures are static — created by an initial database load operation and not modified thereafter. These file structures generally require multiple disk accesses for term translation [9].

Discussions of updatable inverted indexes generally adopt an append-only model of update. Tries, hashing and the ubiquitous B-Tree can all be used to implement an updatable dictionary. An updatable postings file can be implemented using a variety of free space management techniques. The postings for a particular term may be maintained in chained buckets. A new bucket is added to the chain each time a document append causes a bucket to overflow. Alternately, postings may be stored in contiguous extents with free space left after each extent. If the extent overflows, the postings list is copied to a larger extent.

Cutting and Pedersen [7] examine in detail the use of B-Trees to implement an efficiently updatable inverted index. Using their optimizations to the basic B-Tree file structure, a retrieval operation can require as little as one disk access, but may require more, depending on the caching strategy used and the branching factor of the B-Tree. Updates are buffered in main memory and applied in large batches. Postings may be compressed, but the ability to do delete operations is then sacrificed.

Burkowski [3] discusses free space management for the postings file and proposes an organization that groups postings into subsets and leaves free space for appending new postings at the end of each subset. Postings for several different index terms may be combined in a single subset. Postings for different index terms are identified by unique markers assigned at build time. An append that overflows available free space triggers a complete rebuild of the postings file. Retrieval operations cannot be processed while a rebuild is taking place. During a rebuild, the free space usage since the last rebuild is used as a predictive model for free space allocation in the rebuilt postings file. Translating an index term requires at least two disk accesses: one (or more) into the dictionary, and one into the postings file to retrieve the appropriate subset. Update operations other than append can only be performed during a rebuild.

Tomasic, Garcia-Molina and Shoens [15, 17] examine in detail inverted index data structures and update policies under an append-only update model. Like Cutting and Pedersen, they argue that updates must be buffered and applied to the disk in batches. They propose a dual-structured index: one data structure for short postings lists and another for long postings list, recognizing that a typical term frequency distribution will result in a few lists with many postings and a large number of lists with few postings. Data structures for each type of list are then individually developed and optimized. Through a careful experimental study based on real data they show that this approach can exhibit good performance for both query processing and update.

In their discussion they raise interesting questions concerning on-line re-organization of inverted index data structures and point out that very little is known in general about implementation and performance of systems that update inverted list data structures in place. Our paper addresses some of these questions and complements their work by demonstrating similarly good performance for a more general insert/delete update model using a significantly different data structuring approach.

1.6 Our approach

In the following sections we present data structures that efficiently realize the inverted index data abstraction and permit the continuous on-line applica-

tion of updates without significantly disrupting retrieval performance. Accessing the inverted index for a term almost always requires a single disk operation. Update is an on-going background process, rewriting the database on an on-going basis. Updates are maintained in main-memory data structures until they can be applied to disk. The update process is occasionally checkpointed. If a processor failure occurs, the update process may be restarted at the last checkpoint. The integration of caching and index compression is straightforward. Experimental results indicate that the update strategy has very little impact on retrieval performance. An analytic model indicates that competing update strategies are unlikely to exhibit better performance.

The work described in this paper is part of the Waterloo Multi-User Multi-Server Very Large Text Database Project (the MultiText Project). The primary goal of this project is the creation of a prototype distributed full-text retrieval system. Where exposition is simplified and no generality is lost, we use the concrete data structures of the MultiText Project in our discourse.

2 Index Organization

2.1 Interfaces

The index engine ignores document boundaries, treating the text in the database as one continuous sequence. Document boundaries are treated the same as any other structural element. A position in the database is a single positive integer. The granularity of this position — whether it refers to a character, word or to an entire document — depends on the needs of the query language used by the retrieval system. The text structure model used by the MultiText project allows multiple terms to be indexed at a particular location. This property is crucial to schema-independent retrieval [5].

An index engine is responsible for implementing two classes of operations: *retrieval operations* and *update operations*.

A retrieval operation is a request to solve a query, expressed in some query language. Inverted indexes easily implement the boolean-algebra-based languages used by most commercial text retrieval systems [14, 19]. Inverted indexes are also appropriate for implementing the schema-independent heteroge-

neous structured text query language used by the MultiText project [5, 4].

Each query solution is a pair (*start*, *end*) indicating a range in the database that satisfies the query. Each retrieval operation may result in one or more solutions of this form.

An update operation is a request to add or remove indexing. An *add indexing* operation has the form:

- Add (*position*, *term*)

Which indicates that the specified term occurs at the specified position in the database. A *remove indexing* operation can take one of three forms:

- Remove (*start*, *end*, *term*)
Removes indexing of the specified occurrence of the term in the range specified by the start and end positions.
- Remove (*, *, *term*)
Removes all indexing of the specified term.
- Remove (*start*, *end*, *)
Removes all indexing in the specified range. This operation may be used to remove an entire document from the index engine.

The add operation and the first two remove operations change the postings list for a single index term and are referred to as *local update operations*. The third remove operation may affect many postings lists and is referred to as a *global update operation*. It is worth noting here that all the update operations are idempotent, they may be applied one or more times with the same effect. As an aside, the retrieval operation implemented by the text server is:

- lookup (*start*, *end*)

which returns the text associated with the range.

2.2 Resources

A primary concern is the management of the index engine's memory resources: disk storage and main memory RAM.

$$M_D = \text{Quantity of disk storage (words)}$$

$$M_R = \text{Quantity of main memory (words)}$$

We have $M_R \ll M_D$ where the size difference is typically a factor in the range 16 to 256.

The quantity M_R does not include the memory required for the operating system and its internal data structures, or for the index engine application software. To permit control over these resources, we disable all operating system paging, swapping and other memory management. The operating system is used only to provide address translation, physical I/O to the disk, network access and multiprocessing. Disk storage may consist of several physical disks. We assume striping if this is the case.

A word must be big enough to hold:

- any database position.
- any index term
- any integer in the range $[0, M_D]$

A 32-bit word size is sufficient to index approximately 20GB of (uncompressed) English-language text. A 64-bit word is sufficient for all currently conceivable applications. An index term is represented by a word-sized *term symbol*. This mapping of index terms to term symbols is a global function implemented by the Marshaller/Dispatcher. Ordering of the term symbols corresponds to the lexical ordering of the index terms.

2.3 Data Structures

Before examining the on-line update algorithm, we examine the static organization of the index engine data structures as they appear between cycles of the update process.

The vast majority of disk storage is allocated for *index blocks*. Together, the index blocks make up the *index*, which combines the functions of the dictionary and postings file of figure 2. The index contains both index terms and postings. An index block is of fixed size B . The size of the index is M_{DIdx} . The total disk storage allocated for the index is M_{DIdxS} . Both M_{DIdx} and M_{DIdxS} are multiples of B .

When in use, each index block contains one or more *index entries*. Each index entry consists of a term symbol followed by an occurrence count followed by an occurrence list, indicating positions in the database where the term occurs (see figure 3). The occurrence count is never bigger than the space remaining in the index block. Since an entry is always at least three words, up to two words at the end of an index block may be unused. By using the

occurrence counts as relative pointers, we may treat an index block as a linked list of entries. The length of this linked list is at most $\lfloor B/3 \rfloor$.

Overall, the index is ordered first by term symbol and then by database position, and divided into index blocks as appropriate. If an entry for a particular term symbol would not fit in a single index block, it is divided into multiple entries and stored in a range of index blocks.

Disk space allocated for index blocks is treated as a circular array. The index may start at any point in the circular array and does not completely fill it. The dynamic update algorithm operates by modifying the index according to an update list and writing the modified index into the unused portion of the array, releasing storage used by the unmodified index as it is modified and written back (figure 4).

Main memory is used for four purposes:

1. Working storage for retrieval operations (M_{RCache}). This storage is largely a cache of index blocks for solving queries, but also includes memory for queries and partial results.
2. Storage for the *index map* (M_{RMapS}). Each entry in the index map describes a block in the index storage.
3. Working storage for the update algorithm ($M_{RUpdate}$). This storage buffers updates until they are applied to disk.
4. I/O buffers for the update algorithm (M_{RUio}). This will be fairly large as it is crucial to performance that the update algorithm perform I/O in large segments.

The index map contains a two-word description of each block in the index, giving $M_{RMapS} = 2M_{DIdxS}/B$. The index map itself has size M_{RMap} , with the relation

$$\frac{M_{DIdx}}{M_{DIdxS}} = \frac{M_{RMap}}{M_{RMapS}}$$

The first word of each index map entry contains the first term symbol indexed in the block, and the second word contains the first posting for this term symbol indexed by the block. A binary search of the index map allows the determination of the range of index blocks containing the postings list for any particular

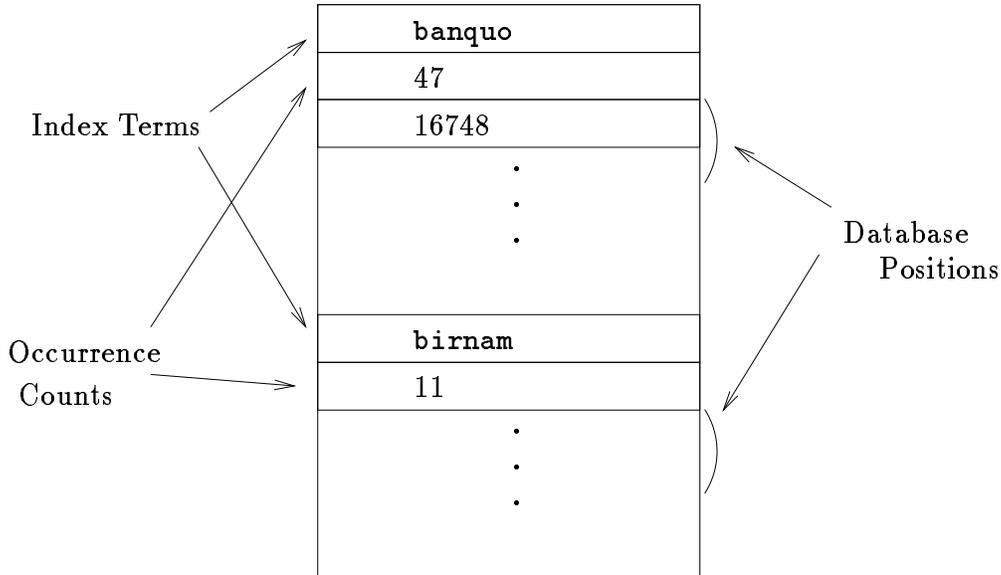


Figure 3: Index block organization.

term. This postings list may then be read with a single disk access. If we are interested only in postings within a limited range of values, either because of restrictions placed by the client or because of earlier partial query results, keeping the first posting in the index map assists in applying this restriction, particularly if the indexing for the term is divided across multiple index blocks.

Space must be allocated on the disk for the storage of a non-volatile copy of the index map. A small amount of additional space must be allocated for a *configuration block* containing parameters describing the disk layout: the start and end locations of the updated and yet-to-be-updated segments of the index. Consistency of the file structures requires that a write of the configuration block be atomic. This space is allocated to a single physical block on the assumption that a write of a physical block will either correctly overwrite the block or will not change the block.

3 Update Strategy

3.1 Update Application

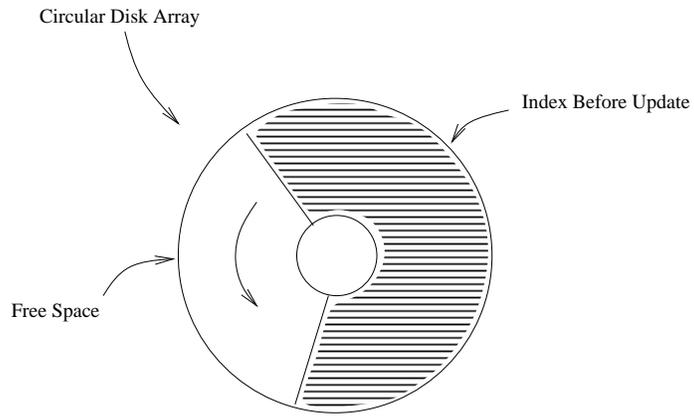
Updates are buffered in main memory until they can be applied to disk. A background process continuously cycles through index storage applying updates and re-writing the index. Update throughput is thus a function of the size of the main memory buffer and the period of an update cycle.

The organization of the index makes update application a simple process. Modified index blocks are written into the free portion of the index storage without affecting the consistency of the index itself. The free portion of the main memory index map is modified in parallel.

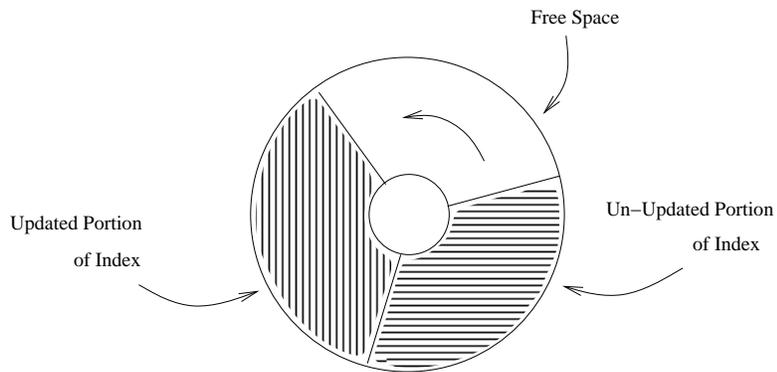
At any point in time, the index blocks are grouped into two segments: an updated segment and a yet-to-be-updated segment. The index map parallels this structure. During a retrieval operation the appropriate segment is selected before performing the binary search to determine the actual blocks to retrieve. In the rare case of an access to the term that stretches across the segments, two accesses to disk are required.

When all available free space is consumed with up-

Before Update Cycle



During Update Cycle



After Update Cycle

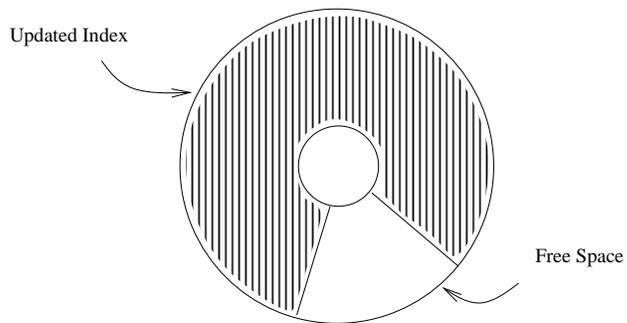


Figure 4: Overview of an update cycle.

dated index blocks, or at any other time deemed appropriate, the index is checkpointed and a range of index blocks is freed. Checkpointing consists of five steps:

1. Update the non-volatile copy of the index map. This involves changes only to portions of the map that are currently free.
2. Disable query access.
3. Invalidate cached indexed blocks that are about to be freed.
4. Re-write the configuration block with new file structure parameters.
5. Enable query access.

Steps 3 and 4, which are performed with query access disabled, proceed quickly as they involve only minor modifications to main memory data structures and a single write to disk. A failure before step 4 leaves the index in a consistent state corresponding to the previous checkpoint. A failure after step 4 leaves the index in a consistent state incorporating the updates.

3.2 Update Management

We assume that update operations are used only by an external maintenance agent responsible for maintaining the overall state of the database. The maintenance agent may run on the Marshaller/Dispatcher and may multiplex many sources of updates into a single source. For the purpose of query operations, updates should take effect as soon as they are received from the maintenance agent but latency in applying the updates to disk is acceptable. An acknowledgement is sent to the maintenance agent when outstanding updates are applied to disk.

Updates are buffered in main memory until they are applied to disk. During this time, we use the collection of unapplied updates as an auxiliary database, modifying the results of retrieval accesses into the main index. When buffered in memory, we assume that information relating to the update will occupy four words: three words for the update parameters and one word to encode the operator type and the information necessary for acknowledging the update.

Global operations and local operations are maintained in separate data structures, the *global update list* and the *local update list* respectively. We require identical properties for both data structures:

- Sequential access to updates in sorted order. The global list is sorted by start range and scanned once per index term during an update cycle. The local list is sorted by index term and then by database position and scanned once per update cycle.
- Insertion and deletion.
- Search. Results from querying the main index will be modified by searches in the local and global update lists.

It must be possible to perform these operations concurrently — sequential scan is an on-going process, updates are continuously arriving from the maintenance agent, retrieval requests are constantly being processed for clients.

The update process acknowledges and deletes updates at each checkpoint. It is possible that a failure after a checkpoint but before updates are acknowledged might result in the maintenance agent having an incorrect view of the index engine's file structures. Idempotency of the update operations ensures that these updates could be re-applied without harm.

Skip lists [12, 13] provide a simple and nearly ideal implementation for the update lists. For the purposes of sequential access, skip lists are effectively linked lists. Insertion, deletion and search are all $O(\log n)$ operations in the average case. Concurrent access to skip lists is simple and efficient. Overhead for pointers is in the range of one to two words per buffered update (depending on an implementation parameter). Total storage overhead for a buffered update is thus at most six words. Data structures other than skip lists may be used, but Pugh's discussion of concurrent maintenance of data structures [12] provides a strong argument for skip lists.

Global operations represent deletions of ranges of text; local operations represent changes to individual postings lists. Adding a document requires many local operations, but deleting the same document requires only a single global operation. For these reasons, global operations tend to be few in number and local operations tend to be relatively many in number. Global operations must be buffered in main memory, but if there is no requirement for updates to take immediate effect there is no need to buffer local updates in main memory. Instead, the local updates may be sorted by the maintenance agent and presented in multiple batches to the index engine over

the course of an update cycle. If this technique is used, disk access is no longer a factor limiting maximum update throughput.

3.3 Failure Recovery and Atomicity

After a checkpoint, updates that were committed to disk during the checkpoint process are acknowledged to the maintenance agent. During failure recovery the maintenance agent must assist the index engine by re-applying updates that were not acknowledged before the failure. As indicated earlier, all update operations are idempotent. It is therefore acceptable for the maintenance agent to re-apply updates that may have been applied to the file structures but had not been acknowledged when the failure occurred.

During update it may be the case that an index engine incorporates only partial indexing for a small number of documents. It is the responsibility of the Marshaller/Dispatcher to recognize this situation and filter from query results any references to partially indexed documents before transmitting the results to clients. The division of responsibility is clear: the Marshaller/Dispatcher maintains the mapping of database positions to document number, using this mapping in its role as coordinator of the system as a whole. Document boundaries are ignored by both text servers and index engines.

4 Performance Analysis

4.1 Implementation Measurements

The update strategy described in this paper has been implemented as part of the MultiText project. However, we do not at this time have experience with the strategy on a large scale under production loads.

The most unusual aspect of our approach is the continuous re-organizational cycle through the index. In order to convince ourselves that this strategy is feasible for use in a production environment, we have performed experiments to understand the possible performance impact of continuous update under a heavy retrieval load. We assumed the worst case: continuous random query accesses with no hits in the cache. Each query access reads a single index block. We varied the update cycle time and examined its effects on the sustainable query access rate. To provide baselines for the measurements we determined

<i>Word Size</i>	64 bits
M_D	2.1 GB
M_R	51 MB
M_{DIdx}	1 GB
M_{DIdxS}	2 GB
M_B	64 KB
M_{RIdx}	256 KB
M_{RIdxS}	512 KB
M_{RCache}	36 MB
$M_{RUpdate}$	12 MB
M_{RUIo}	2 MB

Figure 5: Retrieval system parameters.

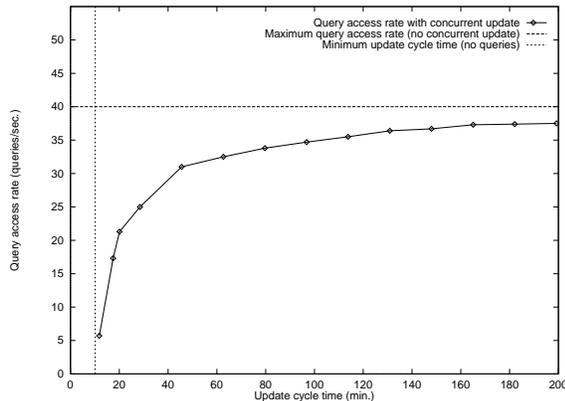


Figure 6: Performance impact of continuous update cycles.

the maximum sustainable access rate that could be achieved with no concurrent update and the minimum possible update cycle time with no concurrent query access. Parameters of the retrieval system for the experiments are given in figure 5. Not all of these parameters are strictly relevant to the experiment; the values for M_{RCache} and $M_{RUpdate}$ should be taken as nominal. The experiments were run on a dedicated DEC Alpha 2000-300 running OSF/1 V1.3 with a 2.9 GB Seagate ST43400N SCSI disk.

The results of the experiment are shown in figure 6. With an update cycle time of as little as 45 minutes, query performance degrades by only 23%. With an update cycle time of 3 hours, query performance degrades by only 7%.

We use the nominal value of 12 MB for $M_{RUpdate}$. Since each update may require six words (and each

word is 8 bytes), this 12 MB allows 256 K of updates to be buffered. Applying 256 K of updates over a 3 hour update cycle gives an update throughput of 24 updates/sec.; if applied over a 45 minute update cycle the update throughput is 97 updates/sec. When considering these numbers, the reader should keep in mind that any increase in the number of buffered updates (by increasing $M_R Update$ or by making more thrifty use of buffer storage) will have a proportional effect on the update throughput. Remember as well, that if updates are batched and fed to the index engine by the maintenance agent, the update throughput is not limited by any storage characteristic of the index engine.

4.2 Analytic Comparison

We analytically compare the performance of our update strategy to the performance of a B-tree for which updates are buffered and merged in batches as recommended by Cutting and Pedersen [7]. We first develop an analytic model for B-tree updates. In this model we make a number of simplifying assumptions, but the overall effect of these assumptions is to improve the performance of the B-trees.

We begin by assuming that all upper-level B-tree nodes are cached and that all postings for a particular term are contained in a single leaf node. When determining performance we count only accesses to leaf nodes.

We base our estimate of leaf nodes accessed on the number of distinct terms V (the *vocabulary size*) appearing in an update batch of size U . Since an update batch is most likely comprised of number of documents to be added or deleted, we may expect that the distribution of terms in the update batch follows Zipf's law [10]. Zipf's law predicts that the product of a term's rank r and its frequency f_r will be constant and consequently that this constant is equal to the vocabulary size:

$$f_r = V/r$$

Summing over rank:

$$\sum_{r=1}^V f_r = \sum_{r=1}^V V/r$$

gives

$$U = V \sum_{r=1}^V 1/r = V H_V \approx V(\ln V + \gamma)$$

where H_V is the V th harmonic number and γ (≈ 0.57722) is Euler's constant.

Let L be the number of leaves in the B-tree. If the distribution of terms over the B-tree leaves is uniform, then the probability that a particular B-tree leaf will be not be updated by a particular batch of updates is

$$\left(\frac{L-1}{L}\right)^V = \left(\frac{L-1}{L}\right)^{\rho L} \approx \frac{1}{e^\rho}$$

where $\rho = V/L$. It follows that the number of leaves that will be modified may be estimated by

$$L(1 - 1/e^\rho)$$

Figure 7 shows the predicted number of B-tree leaves modified while merging update batches of various sizes with B-tree block sizes of 1KB, 8KB and 64KB. Database parameters are equivalent to those of the experimental study: The database is 2GB in size and each B-tree leaf node is 50% full. Each posting is one word in size, for a total of 128M of postings in all. Updating a leaf node requires a disk seek, followed by a read and a write of a single B-tree block. The procedure suggests that large block sizes are to be favoured over small block sizes to reduce the number of seeks, which are costly in comparison to small data transfers. The flattening of the curves as the number of buffered updates becomes larger indicates that all leaves are being modified. This flattening becomes significant with a batch as small as 128K updates when the B-Tree block size is 64KB. A batch of this size represents a change to only 0.1% of the database.

To illustrate the point more clearly, figure 8 shows the percent of B-tree leaves accessed during an update. Once at the 100% level, the entire disk is being read and written while applying the update batch. It is only by using small block sizes that we can hope to avoid re-writing the a large portion of the disk to merge an update batch of any reasonable size. Our update strategy gains its performance benefits by assuming *a priori* that the entire disk will be re-written and by performing I/O in large segments rather than in individual blocks. This approach is not possible

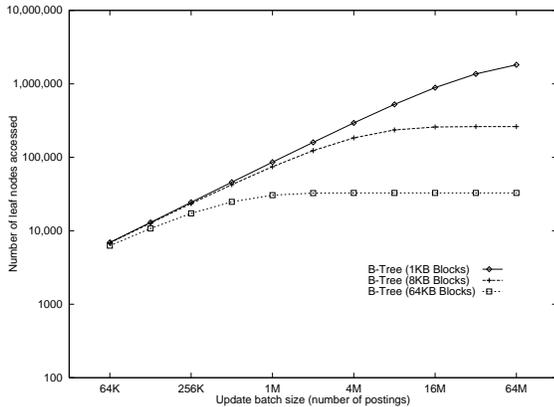


Figure 7: Number of B-Tree leaves accessed while merging buffered updates (database size: 128M postings).

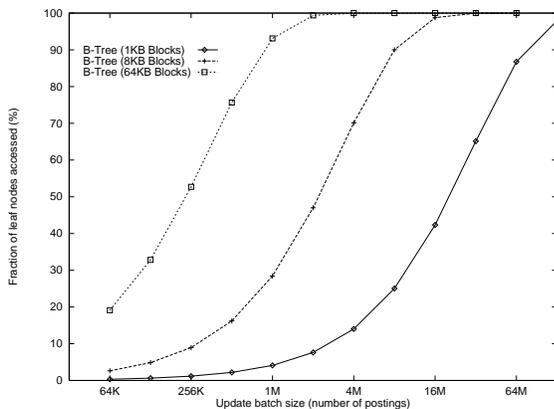


Figure 8: Percent of B-Tree leaves accessed while merging buffered updates (database size: 128M postings).

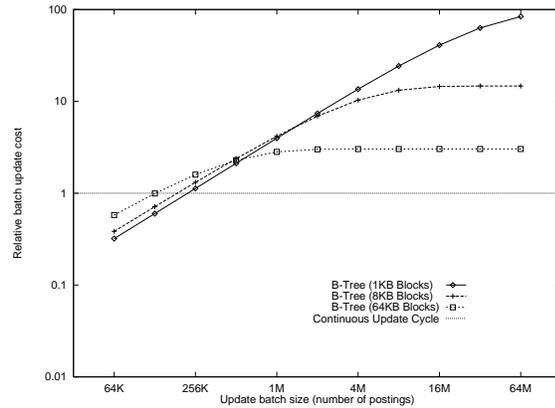


Figure 9: Update performance comparison (database size: 128M postings).

for a B-tree file structure as sequential blocks will be scattered across the disk.

Figure 9 combines the results of our experimental study with our analytic B-tree model and directly compares our update strategy with the use of B-trees. We measured disk access rate for various block sizes and from this predict the time required to apply various sizes of update batches for various B-tree block sizes. These predicted times are scaled by the measured time to apply an update batch using a continuous update cycle. Each plotted line indicates the factor of improvement seen when our update strategy is used. For any reasonable batch size this factor is greater than one. For B-trees with small block size, this factor may grow as large as 100.

Although the results of this section are based on an analytic model and on specific database parameters, it seems reasonable to generalize from the apparent trends. If B-trees are used to implement a dynamic inverted index a design trade-off exists: Merging even a small batch of updates into a B-tree with large block size will require a re-write of the entire database. If a small block size is used, merging a batch of updates will require a large number of expensive disk operations. Our approach avoids dependency on block size and provides consistent update performance regardless of database or disk parameters.

5 Further Issues

5.1 Index Compression

Index compression integrates smoothly into the scheme. Each index block is individually compressed to a variable-length segment. The index map references compressed blocks rather than fixed-size blocks. We add a word to each entry in the index map that indicates the offset of the compressed block in the index. Blocks are decompressed as they are brought into the cache or read by the update process. Since all blocks are cyclically re-written, compression does not hamper update.

5.2 System Considerations

This paper has concentrated on the design of the index engine. We look briefly at a few relevant aspects of the remainder of the system.

The text server translates a range in the database into the associated text. While the details differ, we organize the text server using data structuring principles similar to those used in the index engine.

While our data structures efficiently implement the inverted index data abstraction, they do not efficiently implement queries that are based only on the dictionary. Generally, these queries consist of identifying terms that match specified patterns. In systems that separate the dictionary from the index, this type of query can be satisfied by a dictionary search. Besides its other duties, the Marshaller/Dispatcher is responsible for handling these dictionary-based queries by maintaining a separate dictionary database of all words in the system. In addition, the Marshaller/Dispatcher is responsible for implementing a term thesaurus.

5.3 Wider Application

While our exposition has been in the context of a distributed text retrieval system the data structures and update strategy have wider applicability. Inverted indexes are used in applications other than text retrieval. Even if file structures are static and update is not a requirement (in the case of a CD ROM, for example) our data structures provide an efficient realization of inverted lists. The update strategy has applicability to other databases with similar update

characteristics, with the text server being a ready example.

6 Conclusions

The data structures presented in this paper efficiently realize the inverted index data abstraction. A retrieval operation requires a single disk access in all but rare cases. The update strategy provides high throughput with little impact on retrieval performance. The file structures may be compressed to increase the size of index that can be stored on available disk. Although discussed in the context of a distributed full-text retrieval system, the results of this paper have applicability to any use of inverted indexes and any database with similar update characteristics.

Acknowledgements

The MultiText project is funded by the Government of the Province of Ontario through its Information Technology Research Centre. The Natural Sciences and Engineering Research Council of Canada provided additional financial support.

References

- [1] Timothy C. Bell, Alistair Moffat, Craig G. Nevill-Manning, Ian H. Whitten, and Justin Zobel. Data compression in full-text retrieval systems. *Journal of the American Society for Information Science*, 44(9):508–531, October 1993.
- [2] Forbes J. Burkowski. Retrieval performance of a distributed text database utilizing a parallel processor document server. In *Proceedings 2nd International Symposium on Databases in Parallel and Distributed Systems*, pages 71–79, 1990.
- [3] Forbes J. Burkowski. Surrogate subsets: A free space management strategy for the index of a text retrieval system. In *Proceedings of the 13th Annual International ACM SIGIR Conference*, pages 211–226, Brussels, 1990.
- [4] Charles L. A. Clarke, Gordon V. Cormack, and Forbes J. Burkowski. An algebra for structured

- text search and a framework for its implementation. *The Computer Journal*, 1995. To appear. An early version of this paper was distributed as University of Waterloo Computer Science Department Technical Report number CS-94-30 [11].
- [5] Charles L. A. Clarke, Gordon V. Cormack, and Forbes J. Burkowski. Schema-independent retrieval from heterogeneous structured text. In *Fourth Annual Symposium on Document Analysis and Information Retrieval*, Las Vegas, Nevada, April 1995. An early version of this paper was distributed as University of Waterloo Computer Science Department Technical Report number CS-94-39 [11].
- [6] Gordon V. Cormack. Data compression on a database system. *Communications of the ACM*, 28(1):1336–1342, December 1985.
- [7] Doug Cutting and Jan Pedersen. Optimizations for dynamic inverted index maintenance. In *Proceedings of the 13th Annual International ACM SIGIR Conference*, Brussels, 1990.
- [8] Christos Faloutsos. Access methods for text. *Computing Surveys*, 17(1):49–74, March 1985.
- [9] Donna Harmon, Edward Fox, R. Baeza-Yates, and W. Lee. Inverted files. In William B. Frakes and Ricardo Baeza-Yates, editors, *Information Retrieval: Data Structures and Algorithms*, chapter 3, pages 28–43. Prentice Hall, Englewood Cliffs, NJ, 1992.
- [10] Wentian Li. Random texts exhibit Zipf’s-law-like word frequency distribution. *IEEE Transactions on Information Theory*, 38(6):1842–1845, November 1992.
- [11] The MultiText Project. Project repository: <ftp://plg.uwaterloo.ca/pub/mt>.
- [12] William Pugh. Concurrent maintenance of skip lists. Technical Report TR-CS-2222, Department of Computer Science, University of Maryland, College Park, Maryland, April 1989.
- [13] William Pugh. Skip lists: A probabilistic alternative to balanced trees. *Communications of the ACM*, 33(6):668–676, June 1990.
- [14] Gerard Salton and Michael J. McGill. *Introduction to Modern Information Retrieval*, chapter 2, pages 24–51. McGraw-Hill Computer Science Series. McGraw-Hill, New York, 1983.
- [15] Kurt Shoens, Anthony Tomasic, and Hector Garcia-Molina. Synthetic workload performance analysis of incremental updates. In *Proceedings of the 17th Annual International ACM SIGIR Conference*, pages 229–338, Dublin, Ireland, July 1994.
- [16] Anthony Tomasic and Hector Garcia-Molina. Performance of inverted indices in shared-nothing distributed text document information retrieval systems. In *Proceedings 2nd International Conference on Parallel and Distributed Information Systems*, pages 8–17, San Diego, January 1993.
- [17] Anthony Tomasic, Hector Garcia-Molina, and Kurt Shoens. Incremental updates of inverted lists for text document retrieval. In *Proceedings of the 1994 ACM SIGMOD International Conference*, pages 289–300, Minneapolis, Minnesota, May 1994.
- [18] C. J. van Rijsbergen. *Information Retrieval*. Butterworths, London, second edition, 1979.
- [19] Steven Wartik. Boolean operations. In William B. Frakes and Ricardo Baeza-Yates, editors, *Information Retrieval: Data Structures and Algorithms*, chapter 12, pages 264–292. Prentice Hall, Englewood Cliffs, NJ, 1992.
- [20] Ian H. Whitten, Alistair Moffat, and Timothy C. Bell. *Managing Gigabytes*. Van Nostrand Reinhold, New York, 1994.
- [21] Justin Zobel, Alistair Moffat, and Ron Sacks-Davis. An efficient indexing technique for full-text database systems. In *Proceedings of the 18th International Conference on Very Large Databases*, pages 352–362, Vancouver, August 1992.