

Lawrence Berkeley National Laboratory

Recent Work

Title

BIT TRANSPOSED FILES

Permalink

<https://escholarship.org/uc/item/6pg9b0br>

Author

Wong, H.K.T.

Publication Date

1985-02-01



Lawrence Berkeley Laboratory

UNIVERSITY OF CALIFORNIA, BERKELEY

Information and Computing Sciences Division

RECEIVED
LAWRENCE
BERKELEY LABORATORY
JUN 26 1987
LIBRARY AND
DOCUMENTS SECTION

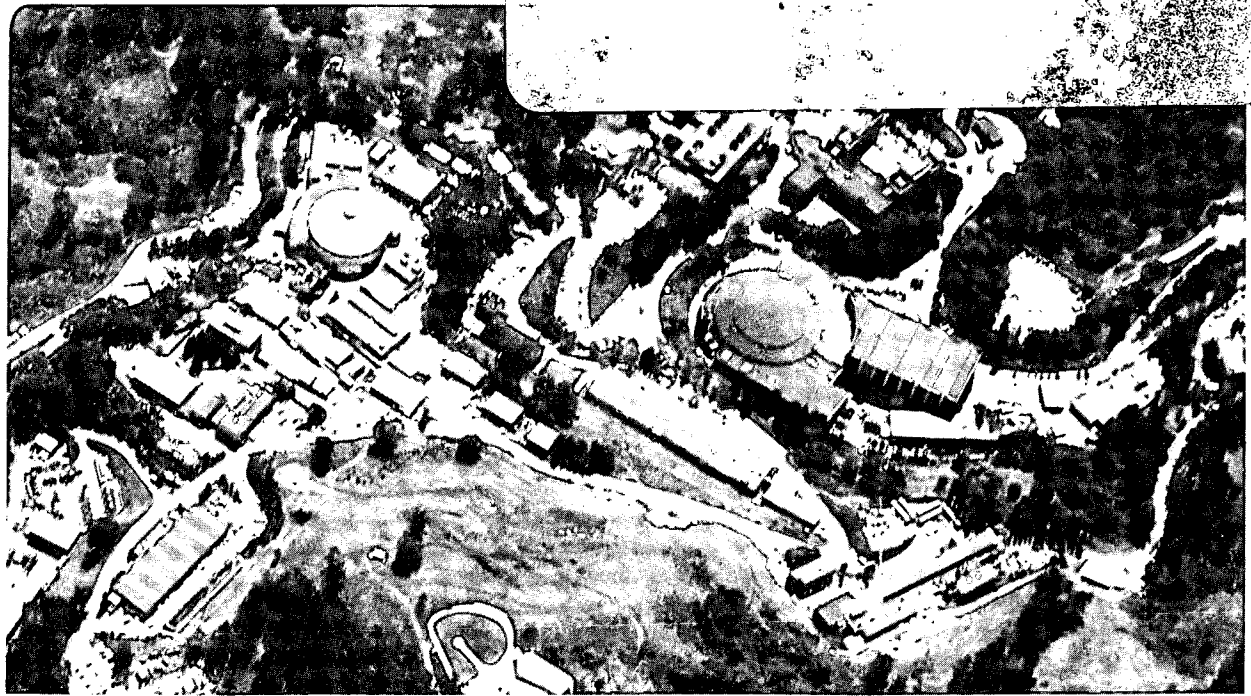
Presented at the Very Large Data Bases Conference,
Stockholm, Sweden, September 26-29, 1985

BIT TRANSPOSED FILES

H.K.T. Wong, F. Liu, F. Olken,
D. Rotem, and L. Wong

February 1985

TWO-WEEK LOAN COPY
*This is a Library Circulating Copy
which may be borrowed for two weeks.*



LIB-19149
c.2

DISCLAIMER

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor the Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or the Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or the Regents of the University of California.

Bit Transposed Files

H.K.T.Wong, F. Liu, F.Olken, D.Rotem, and L.Wong

**Computer Science Research Department
Lawrence Berkeley Laboratory
University of California
Berkeley, California 94720**

February, 1985

This research was supported by the Applied Mathematics Sciences Research Program of the Office of Energy Research, U.S. Department of Energy under contract DE-AC03-76SF00098.

B i t T r a n s p o s e d F i l e s

Harry K.T. Wong, Fanny Liu, Frank Olken, Doron Rotem*, Linda Wong

Lawrence Berkeley Laboratory,

University of California

Abstract

This paper first examines the reasons why sophisticated access methods are often not used in large Scientific/Statistical Database (SSDB) applications. A file structure (called bit transposed file) is proposed which offers several attractive features that are better suited for the special characteristics that SSDBs exhibit. This file structure is an extreme version of the transposed file where the data is stored by vertical bitwise partitions (rather than by attributewise). The bit patterns of attributes are assigned using one of several index encoding methods. Each of these encoding methods is appropriate for different query types and access requirements. The bit partitions can also be compressed using a version of the run length encoding scheme. Efficient operators on compressed bit vectors are available to form the backbone of a query language. In addition to selective power with low overhead for SSDBs, the bit transposed file also is amenable to special parallel hardware. Results from experiments with the file structure suggest that this may be a reasonable alternative file structure for large SSDBs.

Supported by the Office of Energy Research, U.S. DOE under Contract No. DE-AC03-76SF00098.

* Dept. of Computer Science, Univ. of Waterloo, Canada

1. Motivation and Overview

Scientific/Statistical Databases (SSDBs) exhibit many distinctive types of characteristics and usage [Shoshani,Olken,Wong84], [Wong84]. With the advent of many advanced access methods, the dominant file structure for very large SSDBs is still the simple sequential file. The major reason is that there is a "mismatch" between conventional access methods such as inverted files, B-trees, hashing, etc. and the characteristics of SSDBs. First, because the cardinality of SSDBs attributes is typically small, most access methods simply partition the database into a small number of still very large files, with prohibitively expensive overhead for the pointers, structures, tables, etc., with only limited selective power added. Second, since SSDBs are largely static, the expensive overhead associated with the dynamic facilities of most access methods is not justified. Third, the values of SSDBs attributes tend to cluster, and current access methods often do not take advantage of this opportunity for compression. Fourth, the access to SSDBs is typically long "sweep" in that a long sequence of individual records is fetched and a small number of attributes extracted, this kind of range access is not supported well by most access methods.

The search for an appropriate file structure begins with the fourth point mentioned above, which is the motivation for transposed files ([Wiederhold 83], [Batory 79]). The file structure we propose is an extreme form of the transposed file and we call it the Bit Transposed File (BTF).

The BTF has three major components: an index encoder, transposed bit vector loader, and a query processor on bit vectors.

The index encoder translates each field in each record in the database into a series of bits based on several encoding schemes. The result is that each record of the database is translated into a bit pattern.

The second component, called the transposer, stores the bit patterns in a transposed manner so that for each bit position of the bit pattern, a file is produced which contains the bit value of that bit position from all the records in the database. The result is n BTFs where n is equal to the number of bit columns that result after encoding. Because values in large statistical databases tend to cluster, we have developed a compression method to compress the BTFs so that long runs of 0's and 1's can be stored more efficiently.

The third component of this file structure is the query processor on BTFs. The processor translates the retrieval requests on the database into a boolean expression on the BTFs. The translation algorithm takes as input the encoding schemes for the attributes in the query and the query type in order to generate the shortest boolean expression. The boolean expression is then evaluated by using the primitive boolean operators AND, OR, and NOT. These operators are very efficient that can also take advantage of the compressed BTFs.

In section 2 the various index encoding schemes are described with examples. Section 3 gives details and examples to the transposition of records by bits. In Section 4 the query processing aspect is examined. Section 5 formalizes the problem of optimal index encoding assignment and experiment results with the algorithm are included in an appendix. Section 6 describes the implementation and experimentation of the file structure and results are listed in another appendix. Some interesting current work is mentioned in Section 7. Section 8 contains the summary and conclusion of the paper.

2. Index Encoding Schemes

In this section we will describe the available index encoding schemes in our current BTF transposed file structure. Index encoding schemes are crucial to BTFs because they ultimately decide how many boolean operations have to be performed on the bit vectors. There are four basic schemes: binary, k-of-n, unary, and

superimposed. Each one of these schemes can have a composite version for attributes with large number of values. Below we will describe each of them with examples and discuss the usage of the scheme for different kind of queries.

2.1. Binary Encoding

Given an attribute A with n possible values, the binary encoding of A is to use $\log_2(n)$ bits for each value v and the bit pattern for v is the binary number in the range of 0 and n , corresponding to the ordinal integer of v among the n values of A . As a convention, the bit positions are labeled b_0, b_1, \dots, b_n , from the rightmost bit to the leftmost. This scheme requires the minimum of storage but all bits have to be examined for retrieval.

As an example throughout this paper, we will use an application of radiation experiment on dogs. This experiment database contains information such as dog type, weight, age, dosage, location, etc. Assume that there are 10 dog types. To encode dog type using the binary encoding requires 4 bits and the bit patterns of these 10 values range from 0000 to 1010.

2.2. K-of-N Encoding

This encoding scheme assigns bit patterns to attribute values by turning on a distinct set of K bits from N bits. Hence it can encode up to $\binom{N}{K}$ values. For example, the 1-of-10 encoding for dog type mentioned above would involve the following bit patterns:

```
0000000001
0000000010
0000000100
...
1000000000
```

An 2-of-5 encoding for dog type has the following bit patterns:

```
00011
00101
00110
```



```
01001
01010
01100
10001
10010
10100
11000
```

Unlike binary encoding, this scheme requires examining only K bits for any value. It also allows a time-space tradeoff in the sense that more storage space (larger N) would mean less bits to examine (smaller K).

2.3. Unary Encoding

This scheme requires N bits to encode N values and it is useful for attributes that are involved mostly in range or inequality queries. For example, the following is the result of encoding dog type using the unary encoding scheme.

```
0000000001
0000000011
0000000111
...
1111111111
```

To retrieve all dog types that are larger than type 3 requires to examine only bit b3 (if it is 1 or not). Similarly for all dog types that are below type 3 requires to examine only bit b2 (if it is 0 or not). Range queries in the form of (a,b) can be expressed as two inequality queries in the form of $< a$ and $> b$. For example, to find all dog types between 3 and 8 requires examining only bits b2 (greater than 2) and b8 (less than 9). Similarly queries such as $\sim=a$ can be expressed as $< a$ or $> a$. For example, to find all dog types not equal to dog type 3 requires examining bits b2 (less than 3) and b3 (greater than 3).

2.4. Superimposed Encoding

Superimposed encoding scheme ([Knuth73]) is important for SSDs which contain large volume of bibliographical data or property data ([Shoshani, Oiken, Wong84]). To use superimposed encoding for an attribute, a hashing function is first

defined which maps each desired keyword in the attribute into a bit pattern of N bits. Given an attribute value (text with keywords), the collection of bit patterns of all the keywords are superimposed (logically ANDed together) and the resulting bit pattern is the encoded value. This scheme supports partial match queries. Given a list of keywords to be searched, the keywords are hashed, superimposed onto a bit vector and the resulting bit pattern is matched against the superimposed codes of the attribute. Because of the possible "false drops", this scheme can only be used as a "filter" in the sense that only some records not qualifying are eliminated but of the selected ones, a search for the keywords is still required to reject those that were selected because their codes coincide with the superimposed code of the query.

2.5. Composite Encoding

Each of the four encoding schemes mentioned above can be made "composite". Given an encoding scheme E and a bit vector with length N, a composite encoding scheme for E of D fields is the concatenation of D groups of bit vectors, each of which is encoded using E and with length N. For example, suppose there are 1000 possible values for the attribute dosage in our experiment database. An 1-of-1000 encoding would require 1000 bits for each value. A composite 1-of-10 encoding with 3 fields, which involves the concatenation of three 1-of-10 fields together, can be used. To find a particular dosage value, only 3 bits have to be examined, 1 from each field. Composite k-of-n encoding with d fields can be viewed as a n-bit radix number with d digits. It is not required for the fields of a composite encoding scheme to have the same length. For the example above, we could have the first field encoded as 2-of-5 and the last two as 1-of-10.

Given an attribute encoded in a particular scheme, to find the correspondence between a value of the attribute and its bit pattern is done by a code table lookup. The major advantage of the composite encoding scheme is the reduction of the code table size. The reason is that the number of possible encoded values of a composite

encoding scheme is the *product* of the number of possible encoded values of its fields, but the size of its code table is just the *sum* of the size of the code tables of its fields. In fact, in the case that all fields have the same encoding, then the same code table can be used. Another advantage of composite encoding is that for attributes with large number of possible values, multiple levels of grouping can be made so that selection can be performed based on the desired level. For example, in the composite encoding of dosage above (three 1-of-10 fields), there are three levels of grouping of values, one at the hundreds, one at the tens, and one at the ones level. Selection performed at the hundreds, tens, or ones level involves respectively one, two, or three bits. For large SSDBs, having multiple levels of grouping of values is very important and composite encoding scheme is invaluable.

The following table summarizes the properties of the encoding schemes. The formulas are expressed in terms of d (the number of fields, in the case of non-composite encoding, $d=1$), n (the width of each field), and k (the number of bits to turn on in the case of k -of- n encoding).

	# values	exact match	>	partial match
binary	$2^{n \cdot d}$	nd	nd	No
k of n	$\binom{n}{k}^d$	kd	nd	No
binary	$(n + 1)^d$	d	$2d-1$	No
superimposed	$O(2^{nd})$	no	no	*

* depends on code density, typically is $1/2nd$.

3. Bit Transposition

In this section we will describe the file structure using some examples. The steps in obtaining the BTFs involve the following: first, the encoding schemes are decided for selected attributes; then the attributes are encoded for all records in the database; for each bit position of the encoded record, a file consisting of all the bits across the whole database is generated and stored; finally, the files are compressed.

The database of radiation experiment on dogs is used again here to illustrate these steps. The attributes of the database include the dog type, weight, age, dosage, location, observation, etc. Assume the following encoding schemes

attribute	# values	scheme
dog type	10	2-of-5
weight	8	unary (8 bits)
age	20	binary (5 bits)
dosage	200	composite unary (3 fields of 6 bits)
location	10	1-of-10
observation	1000 keywords	superimposed on 10 bits

Using these encoding schemes, the database is transformed into bit patterns. For each bit position, a bit vector is stored as a file. For the example above, the number of bit vectors files is as follows:

attribute	#bit vectors
dog type	5
weight	8
age	5
dosage	18
location	10
observation	10

These bit vectors are then subject to compression. The compression method we use is a variation of the header compression scheme proposed by [Eggers, Olken, Shoshani81], which in turn is a variation of the run length encoding scheme with efficient access to the compressed data. Because of space limitation, the reader is referred to the above paper for the details of the compression method. The BTF

compression scheme has the additional capability of suppressing the compression in the case where the overhead exceeds the gain of compression. This happens when there are a large number of short runs of 1's and 0's. The suppression algorithm involves look ahead and constant evaluation and balance of the cost of the overhead vs the storage gain from the compression.

4. Query Processing

4.1. Boolean Operators on Bit Vectors The primitive operators on bit vectors are the boolean operators AND, OR, and NOT. These operators can be efficiently implemented by breaking up the bit vectors into words and feed to the boolean operators of the CPU. More efficiency is gained when the compression rate of the bit vectors is large. In the case of computing the AND operator between two bit vectors, for example, the runs of 0's in one of the bit vectors can be "skipped", and the corresponding part of the other bit vector can also be skipped. For bit vectors with large compression rate (which is one of the dominant characteristics of SSDBs), this skipping action can be used to produce very fast boolean operators over bit vectors.

4.2. Query Language

The current BTF query language is a simple boolean expression language which allows range, exclusion, and set conditions. For example, to retrieve all female dog records between age 3 to 5 and weigh more than 10 lbs, the following query can be used.

```
sex[1] & age[3:5] & weight[>10]
```

The query "retrieve all dogs except German Sheppards (which has value 105) and dogs that have developed cancer in the brain", can be expressed as

```
dogtype[~105] & observation["cancer","brain"]
```

(Note that in the current implementation of the BTF there is actually a menu-driven user interface which alleviates the user from having to memorize the internal codes of the attributes.)

4.3. Decoding of queries

Given a query, a series of table lookup has to be performed to translate the query into boolean expression of bit vectors. The first table is the attribute index encoding table which records the encoding scheme for each attribute and contains pointer to the attribute's bit assignment table. The bit assignment table records the bit pattern for each attribute value. In the case of composite encoding, there can be up to d value decode tables where d is the number of fields of the composite encoding scheme.

Given the bit assignments for each attribute in the query, the next step is to generate boolean expression on bit vectors. The generation procedure examines both the encoding scheme and the condition in the query for each attribute in order to generate the shortest boolean expression. Below, we will illustrate this step by some examples.

1. Simple exact match queries.

(a) find all German Shepherds

From table lookup, value 105 is found to have bit assignment

01100. The query

`dogtype[101]`

is translated to

`dogtype (b3 & b2).`

and can now be evaluated. (Remember that the bits are named from right to left.)

(b) find all 5-year-old dogs.

Age 5 is encoded as 00101 in a binary encoding scheme, so the following expression is generated

$\text{age} (\sim b_4 \ \& \ \sim b_3 \ \& \ b_2 \ \& \ \sim b_1 \ \& \ b_0)$

(c) find all 5-year-old German Shepherds.

is translated to

$\text{dogtype}(b_3 \ \& \ b_2) \ \& \ \text{age} (\sim b_4 \ \& \ \sim b_3 \ \& \ b_2 \ \& \ \sim b_1 \ \& \ b_0).$

2. Queries with set conditions

find all dogs that have been radiated on locations 1, 4, or 7.

The query is expressed as

$\text{location}[1,4,7]$

Since location is encoded as a 1-of-10, the query is translated to

$\text{location} (b_0 \ | \ b_3 \ | \ b_6).$

3. Queries with range conditions

(a) find all dogs lighter than weight class 7.

Recall that attribute weight is encoded as unary, the above query is translated simply to

$\text{weight} (\sim b_6).$

(b) find all dogs receiving more than 30 dosage units.

Attribute dosage is encoded as a Composite unary with 3 fields of 6 bits. Assume dosage 30 is encoded as 000111,000011,011111. The query can be translated to

$\text{dosage} ((b_{14} \ \& \ b_7 \ \& \ b_4) \ | \ (b_{14} \ \& \ b_8) \ | \ b_{15})$

4.4. Order of Evaluating Bit Vectors

After the boolean expression on bit vectors is obtained, an order of execution is determined which will minimize the running time. The optimal order of execution is to evaluate the bit vectors in the descending order of their compression rates. This is because the skipping action mentioned earlier is maximized. The rearrangement is performed by an algorithm that walks through the boolean expression to produce a new (but equivalent) expression where the order of the bit vectors appearance correspond to the descending order of their compression rates. The new expression is then evaluated from left to right.

5. Index Encoding Optimization

In this section, we would like to consider automating the optimal index encoding for one encoding scheme, the k-of-n. Future work will attempt to extend this approach to incorporate the rest of the encoding schemes.

Given an attribute A with v possible values, the k-of-x encoding method stores each value as a binary number with x digits. Exactly k digits are 1's and the other $x - k$ are 0's. Clearly we can represent at most $\binom{x}{k}$ (the number of combinations of x objects taken k at a time) different values for the attribute using this method and therefore we have the constraint that $\binom{x}{k}$ must be at least v . To meet this constraint we can choose to increase both x and k , increase only x while keeping k small, or increase only k . In any case k will not exceed $\frac{x}{2}$ since $\binom{x}{k}$ is maximized at either $k = \frac{x}{2}$ or $k = \frac{x-1}{2}$ and we will show that increasing k means more boolean operations to answer a query. On the other hand, a large x means that more storage will be required to store the bit vectors. Hence we have a time space tradeoff problem. In this section we address the following problem: Given a certain amount of space to store the bit vectors, what is the optimal partitioning of this space among m attributes such that the expected query processing time is minimized. A more formal

definition of the model and a dynamic programming solution to this problem is now given.

Given a database of N records on m attributes A_1, A_2, \dots, A_m , we would like to store the records as a set of bit vectors. The total number of bits reserved for encoding all attributes is C , so that the total storage requirement is $C \cdot N$. We assume that attribute A_j has v_j possible values and appears in a query with probability p_j . Our problem is to find for each attribute A_j , a k_j and a x_j such that the values for A_j will be encoded in a k_j -of- x_j encoding. We assume that when a value for attribute A_j is mentioned in a query, the amount of boolean operations required to find the appropriate records will be proportional to k_j because this is the number of columns we have to AND / OR in this case. Therefore, minimizing the expected time to answer a query amounts to minimizing

$$\sum_{i=1}^m p_i k_i.$$

The constraints are

$$\sum_{i=1}^m x_i \leq C.$$

$$\begin{pmatrix} x_i \\ k_i \end{pmatrix} \geq v_i.$$

We observe that the minimum value for any x_j is $\log_2(v_j)$, by information theoretic arguments and also the maximum value for k_j that we will consider is $\log_2(v_j)$ because otherwise we can use the usual binary encoding with this cost for query processing. The above optimization problem can be solved by dynamic programming techniques by using the following principal of optimality. Let us denote by $OPT_y(1, 2, \dots, j)$ the optimal expected query cost for the above problem where we only consider attributes A_1, A_2, \dots, A_j and allow these attributes to use a total of y bits. We observe that

$$OPT_w(1, 2, \dots, j+1) = \text{minimum}_y \{ OPT_y(1, 2, \dots, j) + OPT_{w-y}(j+1) \}.$$

In words, every partitioning of w bits for the first $j+1$ attributes is achieved by

finding some y where $y < w$ such that the first j attributes use y bits and the attribute A_{j+1} uses the remaining $w - y$ bits. Among all such feasible partitionings, we have to find the value for y which minimizes the sum of these costs. This provides us with an iterative approach where at each iteration we add one more attribute into consideration until we finally find $OPT_C(1,2,\dots,m)$ which is the optimal way of partitioning C bits among m attributes. A program which implements this idea was written in PASCAL and it took a very short time to compute optimal allocations for all practical size databases that we are currently using in our experiments. The details of the testing of the algorithm appear in Appendix A.

6. Implementation

A prototype of the BTF structure has been implemented in a VAX/VMS environment using mainly C with some assembler coding. The physical level of the prototype includes a compression package, an index encoder, a bit vector bulk loader, a set of boolean operators on compressed bit vectors. At the logical level, we have an user interface module, and a query processor. The user interface component is part of another experimental system called MICSUM, which uses the BTF structure and will be presented in a separate paper.

The largest database we have running using the bit transposed file is a 110,000 records cancer incidents database available from the National Institute of Health. Some performance experiments were performed comparing the retrieval time of the BTF with Datatrieve, a DEC relational DBMS, against the cancer data. The result is that BTF incur much smaller overhead (up to 10 times) and the retrieval time is consistently 10 times or more faster than Datatrieve. More details of some of experiments can be found in Appendix B. Besides the space and retrieval time, the loading time of the data is also of interest. We selected four attributes of the cancer database to have transposed bit vectors. Indices for the same attributes were generated in Datatrieve for a fair comparison. The transposition of the records

into bit vectors took about half an hour on our VAX, but it took Datatrieve 5 days to create two indices and 9 days for 4 indices. In fact, only about 75% of the database was loaded because of the excessive CPU time.

7. Related Work

As we mentioned in the Motivation Section, the basis of our approach is the transposed file, which is popular among SSDB implementors ([Turner et al79]). The BTF can be thought of as an extreme version of the transposed file. In addition to the advantages associated with the transposed file for SSDBs, the bit transposed file offers three potential benefits: indexing capability with minimum of overhead because bit vectors are data *and* indices; better compression rate because of the front compression opportunity (such as a telephone book) and the lack of word, or even byte boundary; and the inherent parallelism (and hence efficiency) associated with the boolean logic on bit vectors.

Two early versions of the BTF appear in [Brill & Tolken 77] and [Kiyoki, Tanaka, Aiso81]. The former only has the binary encoding scheme whereas the latter only the 1-of-n scheme. Neither consider other encoding schemes for different query types, compression of bit vectors, or optimization problems.

8. Current Work

We are concentrating our effort on three major areas: experimentation and development; optimization problems; and special parallel hardware.

Our current development on BTF includes the aggregation operators as well as other relational operators such as join. The aggregation operators will allow summary databases to be generated from BTFs, which in turn can be subject to further manipulation. We are also planning to experiment with more large SSDBs.

The first optimization problem we are working on is the generalization of the optimal index encoding algorithm presented earlier. We are interested in the optimal

index encoding assignment for attribute values considering any of the encoding schemes or their combinations and the values of an attribute may be encoded using more than one encoding scheme to optimize the access requirements. The second optimization problem is the aggregation operation. The problem is to find an optimal order to perform the aggregation among the attributes so that the number of passes over the bit vectors is minimized and the different compression rates associated with the attributes are exploited.

From our experience of implementing the BTF, it is apparent that simple yet powerful multiprocessor hardware can be built to support the file structure. We have a preliminary design for a transposer and a vlsi design for a boolean logic machine. The transposer consists of a 32 by 32 register matrix. 32 words (32 bits each) are read in at a time and the bits are slices into the matrix horizontally. The transposition is done by reading the data vertically from the top 32 registers. The entire database can be transposed using this matrix. The same transposer can also be used to convert from the bit transposed form to record format. The boolean logic machine is organized as a tree where each node is a simple processor with only AND, OR, and NOT operations built in. Given a query, the "tree machine" is dynamically reconfigured to correspond to the parse tree of the query. The data, which is in the form of bit vectors, is fed to the tree machine from the leafs. The result is propagated upward in a pipeline manner towards the root, which produces the result. A prototype 8-processor chip has been designed. The processors are connected in a full crossbar which has the necessary logic to make it dynamically reconfigurable.

9. Summary and Conclusion

The motivation of our research began with the examination of why current access methods are not in use for large SSDB processing. We will review our observations and examine whether our proposal provides part of the solution.

The first characteristic of SSDBs is that attributes tend to have small cardinality. As a result, most current access methods would add limited selective power yet incur large overhead. The BTF takes advantage of this property because small cardinality of attributes implies that it is possible to have small number of bit vectors, hence values can be efficiently retrieved. Also, there is minimal overhead associated with bit vectors because bit vectors are data *and* indices.

The second characteristic of SSDBs is the clustering effect of attribute values. The BTF takes advantage of this property by compressing the bit vectors. Unlike traditional compressed data, however, there is no need to uncompress in order to use the data. Instead the compressed bit vectors are used to implement efficient boolean operators.

The third characteristic is the static (or append only) property of SSDBs that tend to underuse the dynamic mechanism of most access methods. This property justifies the lack of update facilities of the BTF which only has the append operation.

The fourth characteristic of SSDBs is that queries tend to access many records but only on a few attributes. This property is the basic motivation of the transposed files. The BTF can be thought of as a transposed file with a built-in "generalized" indexing mechanism which incurs minimal overhead. Generalized indices because the elaborate index encoding schemes provide a continuum of indexing levels based on access requirements and storage considerations.

We envision the BTF to be used in coexistence with other access methods, especially in situations where efficient index encoding is difficult to obtain. Examples include attributes with continuous domains and very large cardinality. Our current implementation of the BTF, in fact, accommodates other file structures such as sequential files, and transposed files.

In conclusion, we believe that the BTF offers an interesting approach to SSDBs because of its simplicity, low overhead, inherent efficiency due to the parallel bit

operations in computers, the optimization opportunities, and amenability to parallel hardware implementation.

Acknowledgements

We would like to thank Arie Shoshani for his valuable comments. Credits are to Michael Ger for implementing the index encoding algorithm and providing the test data. We would also like to acknowledge the text editing help from Carole Agazzi.

Appendix A Index Encoding Optimization Algorithm Result

This appendix lists the test runs and the CPU time it took the optimization algorithm to obtain the optimal results. The first table contains the input and output of the test runs. For each test run, each attribute has two pairs of numbers. The left number of the upper pair represents the number of possible values for the attributes and the right number is the frequency of the attribute being accessed. The lower pair of numbers (a, b) represents the result of the optimal bit assignment.

The second table lists the CPU time comparison of the exhaustive search method and our dynamic programming method. In some instances, the latter's running time is less than 1% of the brute force method. As can be seen, this method is effi-

Run	Max # Bits	Attr. 1	Attr. 2	Attr. 3	Attr. 4	Attr. 5	Attr. 6	Attr. 7	Attr. 8	Attr. 9	Attr. 10
1	70	89000,30 (19,9)	2567,20 (14,8)	780,30 (12,5)	1000,2 (10,18)	5,6000 (5,1)	40,60 (10,2)				
2	80	800,20 (10,10)	56,400 (8,3)	70,3000 (13,2)	687,30 (10,10)	20,400 (6,3)	789,60 (12,5)	2,1000 (1,1)	38,600 (10,2)	50,300 (8,3)	4,200 (2,2)
3	70	5670,30 (15,7)	456,20 (9,9)	900,70 (13,5)	690,200 (14,4)	456,30 (9,9)	590,20 (10,10)				
4	80	8400,30 (14,14)	600,60 (12,5)	56,400 (8,3)	70,20 (7,7)	700,1000 (13,4)	60,800 (12,2)	9,100 (4,4)	567,30 (10,10)		
5	80	6790,30 (13,13)	89000,200 (23,6)	34567,90 (19,7)	23,1000 (8,7)	560,20 (10,10)	30,30 (7,7)				
6	90	500,20 (13,4)	600,30 (13,4)	700,10 (12,5)	60,100 (12,2)	30,200 (9,2)	6,100 (6,1)	36,9 (6,6)	25,10 (5,5)	46,100 (11,2)	3,1000 (3,1)

cient enough for most practical databases.

Run	Exhaustive Search Method *	Dynamic Programming Method *
1	1335	172
2	39603	411
3	3318	176
4	19643	320
5	2138	195
6	250576	595

* measured in CPU milliseconds in a CDC CYBER-170/730

Appendix B Performance Comparison

The database is a real cancer incidents records. It contains information such as the patient's sex, age, cancer site, type of cancer cells, year, etc.

The first table lists the size of the test database in Datatrieve and BTF. The overhead column of BTF is the size (in number of 512-byte pages) of the bit vectors. The overhead for Datatrieve is the size of the indices.

The list of queries contains twenty queries, ten in BTF syntax, and ten in Data-trieve syntax.

The second table lists the running time of the listed queries (in terms of minutes, seconds and fractions of seconds).

	# records	DB size (in pages)	Overhead (in pages)
BTF	110,000	6,974 ⁺	1,332
DATATRIEVE	83,729 [*]	8,100	10,134

DB Sizes

- + The size of the DB after four attributes are index encoded.
- * Only about 75% of the original DB is loaded because of excessive CPU time.

QUERY	BTF	DATATRIEVE
1	00:04.03	00:43.06
2	00:24.92	05:22.03
3	00:10.84	04:43.45
4	00:06.96	02:11.59
5	00:26.98	06:50.20
6	00:02.18	00:56.60
7	00:07.24	00:19.47
8	00:11.77	03:18.08
9	00:02.68	03:12.91
10	00:02.35	02:22.01

List of Queries

1. B: year[75]
D: find r01key4 with year = 75
2. B: year[73:78]
D: find r01key4 with year bt 73 and 78
3. B: year[73:77] & racerea[2]
D: find r01key4 with year bt 73 and 77 and racere = 2
4. B: year[75,77] & sexre[1]
D: find r01key4 with (year = 75,77) and (sexre = 1)
5. B: sexre[1] & racerea[1]
D: find r01key4 with sexre = 1 and racere = 1
6. B: year[74] & agere[10:12]
D: find r01key4 with year = 74 and agere bt 10 and 12
7. B: site[570:579] & sexre[1]
D: find r01key4 with site bt 570 and 579 and sexre = 1
8. B: year[76:78] & sexre[2]
D: find r01key4 with (year bt 76 and 78) and sexre = 2
9. B: year[73,75,77] & site[859]
D: find r01key4 with year = 73, 75, 77 and site = 859
10. B: year[76,78] & histolog[9730,9731]
D: find r01key4 with (year = 76,78) and (site = 9730, 9731)

References

[Shoshani,Olken,Wong 84]

Shoshani, A., Olken, F., Wong, H.K.T., "Characteristics of Scientific Databases", Proc. 1984 VLDB, Singapore, 1984.

[Wong84]

Wong, H.K.T., "Micro/Macro Statistical/Scientific Database Management", The First IEEE International Conference on Data Engineering, Los Angeles, March, 1984.

[Eggers, Olken, Shoshani 81]

Eggers, S., Olken, F., Shoshani, A., "A Compression Technique for Large Statistical Databases", in Proc. 1981 VLDB, Cannes, France, Sept, 1981.

[Turner et al79]

Turner, M., Hammond, R., Cotton, F., "A DBMS for Large Statistical Databases," Proc. 1979 VLDB, Rio de Janeiro, 1979.

[Brill & Tolken 77]

Brill, R.C, Tolken, S.E., "Subset Selection by Boolean Calculation", Unpublished memo, 1977.

[Knuth 73]

Knuth, D.E., *The Art of Computer Programming*, Volume 3, Addison Wesley, 1983.

[Batory79]

Batory, D.S., "On Searching Transposed Files," ACM TODS, Vol. 4, no. 4, Dec., 1979, 531-544.

[Wiederhold 83]

Wiederhold, G., *Database Design*, McGraw-Hill, 2nd Edition, 1983.

[Kiyoki, Tanaka, Aiso 81] Kiyoki, Y., Tanaka, K., and Aiso, H., "Design and Evaluation of a Relational Data Base Machine Employing Advanced Data Structures and Algorithms", in The 8th Annual Symposium on Computer Architecture, MAY 12-14, 1981, Minneapolis, Minn.

*LAWRENCE BERKELEY LABORATORY
TECHNICAL INFORMATION DEPARTMENT
UNIVERSITY OF CALIFORNIA
BERKELEY, CALIFORNIA 94720*