

Compressed Binary Bit Trees: A New Data Structure For Accelerating Database Searching

Andrew Smellie*

CambridgeSoft Corporation, 100 CambridgePark Drive, Cambridge, Massachusetts 02104

Received September 8, 2008

Molecules are often represented as bit string fingerprints in databases. These bit strings are used for similarity searching using the Tanimoto coefficient and rapid indexing. A new data structure is introduced, the compressed bit binary tree, that rapidly increases search and indexing times by up to a factor of 30. Results will be shown for databases of up to 1 M compounds with a variety of search parameters.

1. INTRODUCTION

Commonly, molecular structure is encoded in databases as a bit string or *fingerprint*. The bits in the fingerprint encode for the presence of certain functional groups or substructures in the molecule, as shown in Figure 1. Typically, fragments of molecules are identified and hashed to a short integer that is the index of a bit in the fingerprint to which the fragments map.^{1,2} Fingerprints tend to be short (<2048 bits) so they can be stored compactly in a computer.³ Because there are many, many more possible substructures or fragments than available bits in the fingerprint, each bit typically represents more than one possible substructure.

Thus, each molecule in a database has an associated fingerprint of fixed length L with one or more bits set to TRUE (or ON) representing the presence of one or more functional groups or fragments. A common metric used to compare the fingerprints of two molecules is the Tanimoto coefficient.⁴ Though this a well documented and widely used metric for fingerprint distance,^{5,6} it will be described in detail here because the data structures described in this paper were developed exclusively for acceleration of calculations in fingerprint space using the special properties of the Tanimoto distance metric.

Let:

M_1, M_2 = molecules 1 and 2

$F(M_1), F(M_2)$ = bit string fingerprints of M_1 and M_2 , respectively.

$F_i(M_1)$ = i th bit of the fingerprint of M_1

Then:

$$T(M_1, M_2) = \frac{F(M_1) \cap F(M_2)}{F(M_1) \cup F(M_2)} \quad (1)$$

representing the Tanimoto distance between the fingerprints of molecules 1 and 2.

There are many databases that support exact or substructure searching using bit string fingerprints. The earliest method used brute force searching,⁷ where a query (consisting of a bit string fingerprint representing some query structure) was compared serially against a set of molecules in the database. This is the so-called *brute force* method.

Later developments introduced powerful pruning techniques (Baldi et al.⁸) which meant that a large fraction of the database could be ignored during a search resulting in faster searching. Fingerprint scaling⁹ has also been developed to boost the performance of bit string searching, principally in terms of the quality of hits returned with a beneficial side effect of improved search performance. The techniques described in this paper belong to the class of methods that eliminate a large fraction of the database to be searched that result in an increase of speed. It is not concerned with improving the quality of fingerprint searching. But in contrast to the Baldi method, the bit binary tree is a totally different representation of the fingerprint in memory.

2. THE BIT BINARY TREE

Recognizing that molecules are stored as short bit strings, the bit binary tree is a simple data structure that can encode the fingerprints of every molecule in a database simultaneously. The data structure is very simple to construct. Starting at the first bit in the bit string of the molecule, the tree is traversed left if the bit is 0 and traversed right if the bit is 1. The method for inserting a molecule M in the tree is shown in Algorithm 1 (see Chart 1).

An example of molecules inserted into a bit binary tree is shown in Figure 2. The tree shows three molecules (A, B, and C), each with a fingerprint of 7 bits. Inserting molecule A into the tree results in the tree path illustrated by dotted lines. At the first level in the tree (corresponding to the first bit in the fingerprint of A), the tree is traversed to the left. At the second level (corresponding to the second bit in the fingerprint of A), the tree is also traversed left. At the third bit, the tree is traversed right, corresponding to the third bit in the fingerprint of A. The bitstring is traversed serially according to each bit in the fingerprint until the last bit is reached. A leaf is inserted into the tree, and an identifier for molecule A is stored at that leaf. The same procedure is followed for molecules B and C. Note that in this case, the three molecules share the first two bits in common (each with value 0), resulting in a similar path through the tree.

* Corresponding author e-mail: asmellie@cambridgesoft.com.

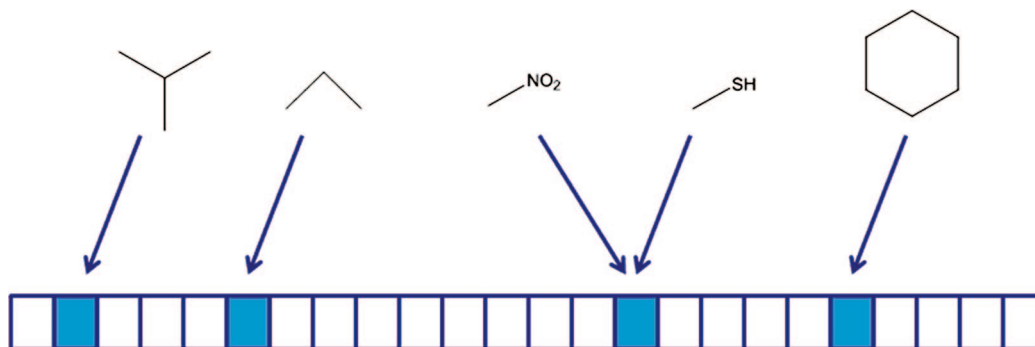


Figure 1. Fragments mapped into fingerprints.

Chart 1. Algorithm 1: Insert Molecule M into the Bit Binary Tree

- a) Compute the fingerprint $F(M)$ of the molecule M
- b) Set $currentNode = Root$ node of the tree
- c) Set $maxBitPos = Total$ number of bits in $F(M)$
- d) Set $bitpos = 0$
- e) While($bitpos < maxBitPos$) {
- f) if ($F_{bitpos}(M)$ is TRUE) { // test "bitpos"th bit of the fingerprint
- g) $currentNode = currentNode \rightarrow RightChild$ // traverse right in the tree
- h) }
- i) Else {
- j) $currentNode = currentNode \rightarrow LeftChild$ // traverse left in the tree
- k) }
- l) $bitpos = bitpos + 1$
- m) }
- n) Store molecule M at $currentNode$

3. STANDARD DATA SETS AND EXPERIMENTAL CONDITIONS

All results in this paper were generated using 4 sets of molecules containing 10000, 100000, 450477, and 967749 molecules. These molecules were randomly selected from the collection of databases CambridgeSoft currently distributes including ChemACX,¹⁰ Merck Index 14.0,¹¹ R&D Insight for Chemists,¹² ChemINDEX,¹³ NCI,¹⁴ TCM,¹⁵ Nanogen,¹⁶ and Ashgate.¹⁷ Details of these molecule sets, including some general properties, are shown in Table 1.

Table 1. Data Sets Used in the Study

set size	mean molecular weight	mean number of rotors
10000	257.5	3.84
100000	284.28	4.81
450477	316.14	5.29
967749	286.5	4.63

All calculations were done on an 4 CPU Xeon 3.0 GHz machine with 2GB of memory and 232GB hard disk. The calculations were performed in single CPU mode.

	1	2	3	4	5	6	7
A	0	0	1	0	1	1	0
B	0	0	0	1	1	0	0
C	0	0	0	0	0	1	0

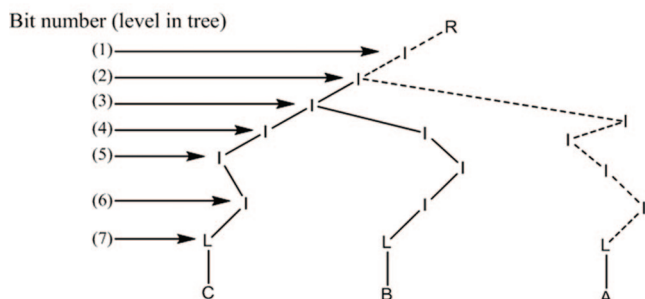


Figure 2. Inserting a molecule into the bit binary tree.

4. APPLICATION TO DATABASE SEARCHING-INVERTED SCREENING

Molecular fingerprints were originally devised to accelerate database searching.¹⁸ Given a query molecule Q, a database is searched to find target molecules that contain Q as a substructure (substructure search mode) or as a complete structure (full structure search mode). The fingerprints of each target molecule are precomputed and stored in the database. The fingerprint of the target query is computed at run time when the database search is initiated. As a *screening* step, target molecules are retrieved that share all bits in common with the query.

Historically, two types of screening have been used in database searching. In *regular screening*, the fingerprint is stored with each molecule. When a query is presented to the search, each target is examined individually, and the bits

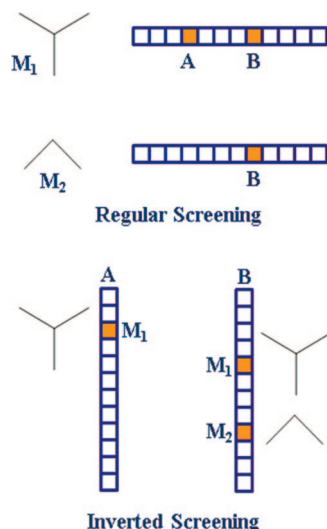


Figure 3. Regular and inverted screening.

in each target are checked against the corresponding bits of the query. Only targets that have all bits in common with the query are returned from the indexing search. The complexity of the search is $N_T \times L$, where N_T is the number of targets and L is the fingerprint width. In *inverted screening*, a list is maintained for each possible bit in the fingerprint consisting of molecules that contain that bit. When a query is presented to this search, each bit in the query that is set TRUE has its list of molecules retrieved. A join of these molecule lists yields a list of molecules that share all bits in common with the query.

Regular and inverted screening is shown in Figure 3. Molecules M_1 and M_2 are shown in regular screening. Molecule M_1 has bits A and B set in its fingerprint. Molecule M_2 has bit B set in its fingerprint. To search a database in regular screening mode, each target's fingerprint must be retrieved and compared against a query to determine if a target is a hit in screening. In inverted screening, lists of molecules are maintained for bits A and B. Molecules M_1 and M_2 are in B's list, but only molecule M_1 is in A's list. If a query is used in a search in inverted index mode, each bit of the query that is TRUE has its bits molecule list retrieved and joined with the other lists from other TRUE bits in the query. For large databases and many TRUE query bits the joins can be expensive. Many systems do a reduced number of joins to save time.¹⁹

The bit binary tree can be used for screening simulating an inverted screen in which the molecule lists for *every* query bit which is ON is effectively joined. To perform this screen, the bit binary tree is exhaustively traversed starting from the root node. The process is illustrated in Figure 4. The fingerprints of four molecules (A, B, C, and D) are shown, with each fingerprint having 7 bits. A query Q is used to search the tree. In this case, a perfect inverted screen would return all molecules that have the same bits on as Q (i.e., bits 4 and 5). The bit binary tree is traversed. If at depth d in the tree, the d th bit of the query is ON, the tree is only traversed to the right (because moving right in the tree corresponds to the presence of an ON bit in the target inserted into the tree). If the d th query bit is OFF, the tree must be traversed both right and left at depth d . This situation corresponds to the case where a bit is OFF in the query and may be ON or OFF in any target it may hit during inverted

	1	2	3	4	5	6	7
A	0	0	1	0	1	1	0
B	0	0	0	1	1	0	0
C	0	0	0	0	0	1	0
D	0	0	0	1	0	0	0
Q	0	0	0	1	1	0	0

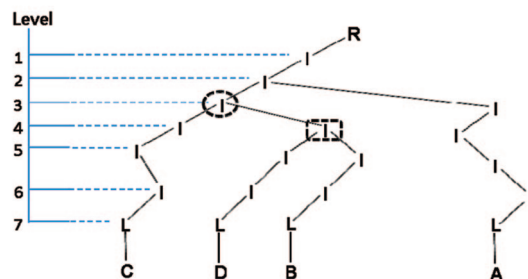


Figure 4. Inverted screening with the bit binary tree.

screening. So in Figure 4, when the tree is traversed, all nodes encountered at depth 4 or 5 must be traversed only to the right (because bits 4 and 5 in the query are ON). This results in only the leaf containing molecule B being reached in the search.

The performance of the bit binary tree as an inverted indexing tool is shown in Table 2. The search times reported are for the bit binary tree and a brute force algorithm in which the molecule joins are performed using precomputed molecule lists (stored in RAM) for each available bit in the query. For every bit that is ON in the query, the molecule lists are joined exhaustively in memory. For the purpose of this test, both the bit binary tree and all the molecular fingerprints are resident in memory to remove any issue of disk access when measuring run times.

It can be seen from the table that as expected as the size of the database being searched increases the number of hits returned also increases approximately linearly. The brute force search time (where each bit string is examined in memory) also increases linearly. It can be seen that the bit binary tree is approximately 4–5× times faster in runtime. Interestingly the percentage of nodes visited in the tree search remains approximately constant independent at the size of the database. One result is not reported: for 967,749 molecules. This was due to the bit binary tree occupying too much memory—more on this later.

5. APPLICATION TO DATABASE SEARCHING-SIMILARITY SEARCHING

Despite the fact that bit strings were originally devised for database screening they have also been used extensively in molecular similarity searching.^{20,21} As previously mentioned, one of the most common metrics of molecular similarity using bit strings is the Tanimoto coefficient. Typically, a user will define a similarity threshold (T_{min}). Molecules are deemed to be similar to the query if they have a Tanimoto distance to the query greater than or equal to the user-defined threshold.

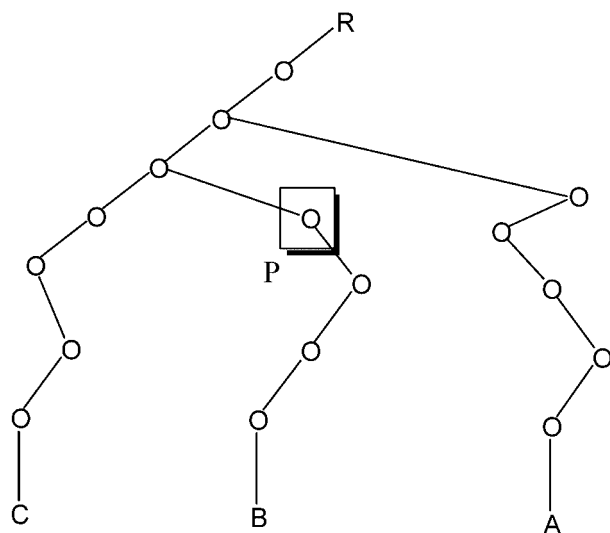
The bit binary tree can take advantage of the particular nature of this coefficient to achieve dramatic improvements in the speed of similarity searching using the *Tanimoto lookahead condition*.

Table 2. Inverted Screening Performance with a Regular Bit Binary Tree

set size	total number of tree nodes	number of tree nodes per molecule	mean number of inverted screen hits ^a	mean brute force search time (ms)	mean bit tree search time (ms)	mean percent tree nodes visited
10000	2,471,816	247.18	110.6	5.86	1.54	3.42
100000	23,160,348	231.60	1151.72	66.21	14.55	3.13
450477	78,105,106	173.38	5640.04	283.90	57.74	3.64
967749 ^b	-	-	-	-	-	-

^a Average using the first 2000 molecules as queries against the entire set of molecules. ^b Test not possible due to tree memory growth.

Bit number	1	2	3	4	5	6	7
Molecule A	0	0	1	0	1	1	0
Molecule B	0	0	0	1	1	0	0
Molecule C	0	0	0	0	0	1	0
Query	0	0	1	1	1	0	0

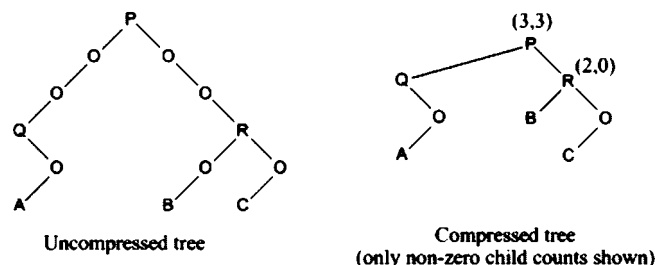
**Figure 5.** The Tanimoto lookahead condition.**Table 3.** Similarity Search with a Regular Bit Binary Tree

set size	user-defined similarity threshold ^a	mean number of hits ^b	mean brute force search time (ms)	mean bit tree search time (ms)	mean percent tree nodes visited
10000	0.6	41.1	8.73	38.57	15.06
100000	0.6	242.21	90.34	336.49	13.63
450477	0.6	686.51	372.91	1069.35	12.77
10000	0.7	18.86	8.69	21.6	8.44
100000	0.7	91.56	90.27	178.36	7.11
450477	0.7	234.04	373.34	559.59	6.57
10000	0.8	9.47	8.11	10.2	3.9
100000	0.8	37.66	90.44	75.25	2.92
450477	0.8	92.5	373.48	229.1	2.61
10000	0.9	3.86	8.50	2.48	0.95
100000	0.9	11.9	90.24	13.22	0.49
450477	0.9	28.3	373.07	35.94	0.38

^a Using each test molecule as a query, all molecules in the set are retrieved within this Tanimoto distance of the query. ^b All values are averaged using the first 2000 molecules as queries against the entire set of molecules.

The Tanimoto lookahead condition is shown in Figure 5 where a bit binary tree is constructed from three molecules A, B, and C. A query is also shown with bit string 0011100. Naively, all molecules that are most similar to Q can be found

Bit Number	1	2	3	4	5
Molecule A	0	0	0	1	0
Molecule B	1	1	1	0	0
Molecule C	1	1	1	1	0

**Figure 6.** The compressed bit binary tree.

by computing the Tanimoto distance (from eq 1) from the query to each molecule. All molecules with $T(M,Q) > T_{\min}$ are deemed to be “similar” to Q. To perform a similarity search with the Tanimoto lookahead and the bit binary tree, the query is marched through the tree subject to the following traversal algorithm.

Any node in the tree (for example, node P at depth d in Figure 5) defines a partial bit string. Any leaf node (at which a molecule is stored) which is reached via node P will have the first d bits in common in their bit strings.

Defining:

- $B(n,d)$ = partial bit string formed from the path traversed through the bit binary tree from the root node to node n at depth d

- $Q(d)$ = partial bit string formed from the first d bits of the query

- $QL(d)$ = the total number of ON bits in the query after position d in the bitstring

Then we can define the terms

$$RUC(n,d) = B(n,d) \cup Q(d)$$

This is the *running union count* and computes the total number of ON bits between the first d bits of the query and the partial bit string from the tree at node n.

$$RUC(n,d) = B(n,d) \cap Q(d)$$

This is the *running intersection count* and computes the total number of ON bits in common between the first d bits of the query and the partial bit string from the tree at node n.

As an example, the partial bit string from the tree traversal at node P is 0001 (three left traversals and one right traversal). The partial bit string from the query is 0011. Thus the running union count at P using query Q is $10001 \cup 00111 = 2$, and the running intersection count is $10001 \cap 00111 = 1$.

Table 4. Inverted Screening with a Compressed Bit Binary Tree

set size	total number of tree nodes	number of tree nodes per molecule	mean number of inverted screen hits	mean brute force search time (ms)	mean bit tree search time (ms)	mean percent tree nodes visited
10000	455,171	45.5	110.6	7.12	0.66	5.36
100000	4,948,210	49.48	1151.72	79.36	7.0	4.32
450477	20,740,319	46.0	5640.04	340.6	29.65	4.51
967749	45,557,420	47.1	14573.32	699.7	73.57	5.22

Table 5. Similarity Search with a Compressed Bit Binary Tree

set size	search radius ^a	mean number of hits ^b	mean brute force search time (ms)	mean bit tree search time (ms)	mean percent tree nodes visited
10000	0.6	41.1	10.1	4.5	21.95
100000	0.6	242.21	101.89	48.58	18.8
450477	0.6	686.51	434.58	168.08	15.70
967749	0.6	2040.56	898.61	363.97	15.84
10000	0.7	18.86	10.0	3.0	13.75
100000	0.7	91.56	101.73	31.06	10.99
450477	0.7	234.04	434.29	104.31	8.8
967749	0.7	610.98	896.96	227.16	9.01
10000	0.8	9.47	9.55	1.67	7.0
100000	0.8	37.66	101.63	16.33	4.96
450477	0.8	92.5	434.62	52.8	3.77
967749	0.8	192.92	892.07	115.74	3.95
10000	0.9	3.86	9.54	0.47	1.84
100000	0.9	11.9	101.71	4.13	0.96
450477	0.9	28.3	434.75	11.79	0.64
967749	0.9	53.65	892.52	27.95	0.71

^a Using each test molecule as a query, all molecules in the set are retrieved within this Tanimoto distance of the query. ^b All values are averaged using the first 2000 molecules as queries against the entire set of molecules.

Recalling the definition of Tanimoto distance from eq 1, the distance between the query molecule and any target molecule is defined as the number of ON bits the molecules have in common divided by the total number of ON bits between the target and the molecule. As we traverse the bit binary tree guided by the query we can assess at any node whether molecules reached at the leaves of the tree by traversing through that node can possibly have a Tanimoto score that is greater than or equal than the user-defined value.

At any node P, the maximum value the Tanimoto distance can possibly be between the query Q and *any* molecule that is reached by traversing the bit binary tree through node P is given by

$$T(P, Q) = \frac{RIC(p, d) + QL(d)}{RUC(p, d) + QL(d)} \quad (2)$$

The Tanimoto lookahead condition is now trivially stated

$$T(P, Q) < T_{\min} \quad (3)$$

The tree traversal can be terminated at node P if eq 3 is TRUE. Intuitively, the tree pruning will be more effective at high user-defined similarity thresholds.

The performance of molecular similarity searching in the bit binary tree can be seen in Table 3, where molecular databases of different sizes are searched using the first 2000 molecules in each database as test queries. Searches were performed at different user-defined similarity thresholds of 0.6, 0.7, 0.8, and 0.9, respectively. Tree search times were compared against brute force search times (where the bitstrings of all molecules were held in memory and the

Tanimoto distance was measured between each molecule and the target query). At lower similarity thresholds, brute force searching was faster than tree searching. This is because more tree nodes are traversed for lower values of the user defined similarity threshold (eq 3). For a threshold of 0.6, between 12% and 15% of tree nodes were traversed for different database sizes. Tree searching begins to get faster than brute force searching at higher user similarity thresholds (>0.7). At high similarity thresholds (0.9) the tree search is significantly faster as only a small fraction of the tree nodes are being traversed (<1%).

This observation led to an improvement in the algorithm. Table 2 lists the total number of tree nodes for each database that was searched. Each bitstring had a maximum length of 384 bits. When molecules are encoded in the tree they share a path of length d in the tree if they have first d bits in common. Thus the average number of nodes per molecule in tree varies between 247 and 173. As expected as the database size increases, the average falls since molecules will share some number of bits in common at the beginning of their bit strings. However, even with this commonality property, the number of nodes in the tree rises very quickly. Even with careful coding a tree node occupies about 16 bytes, so 100,000,000 nodes consumes about 1.6GB of memory. However, most molecules do not switch on very many bits in the bitstring, and run length encoded²² versions of the bit string can be very terse because there are long series of 0's in the bit string. In the uncompressed version of the tree, long series of 0's equates to a succession of left children in the tree.

6. THE COMPRESSED BIT BINARY TREE

The compressed bit binary tree records at each node the number of successive left and/or right children (with no other branching). Each node records the count of these successive left and right children. Figure 6 shows the construction of a simple compressed bit binary tree. From the uncompressed tree, there is a root node P that has 3 successive left children (ending at node Q) and 3 successive right children (ending at node R). In the compressed version of the tree, the nodes between P and Q are removed, and a "left" count is introduced on P to denote that there are 3 successive left children terminating at node Q. The *actual* left child of P is set to Q, and the left count on P is set to 3.

In a similar fashion, the nodes between P and R are removed, the actual right child of P is set to R, and the right count on node P is set to 3. Finally, the nodes between R and B are removed, the actual left child of R is set to B, and the left count of R is set to 2.

The efficacy of this compression can be seen from examining the total number of tree nodes in the uncompressed and compressed version of the tree built on the same

data (listed in Table 2 and Table 4, respectively). For example, with 100,000 molecules the uncompressed tree had 23,160,348 nodes and the compressed tree had 4,948,210 nodes, a factor of 4.7 less. In the test performed, the machine used was not able to handle a data set of 967,749 molecules because of lack of memory. However, with the compressed tree this database was easily searched. The number of new nodes per molecule dropped from 247–173 to 45–50.

7. APPLICATION TO DATABASE SEARCHING-SCREENING REVISITED

Even with added complexity of traversing the compressed tree, it can be seen from Table 4 that the search times for inverted screening were *always* faster for the compressed tree (by about a factor of 10) for any size of database searched. Interestingly, unlike uncompressed trees, the percentage of nodes visited in the compressed tree for searching any size database was about the same (4–5%). Comparing the inverted screening search times between the compressed and uncompressed versions of the tree, the compressed version was about 2× faster.

8. APPLICATION TO DATABASE SEARCHING-MOLECULAR SIMILARITY

The results for molecular similarity searching are even more striking. At *any* user defined threshold tested (0.6–0.9) and any size database searched, the tree search was always faster. For small databases and low similarity thresholds (where more of the tree had to be searched), the speed increase was on the order of 2.5. However for large databases and high similarity thresholds, the speed increase was a factor of 32 faster (27.95 ms vs 892.52 ms). The increase in performance is a direct consequence of the compressed tree having fewer nodes to traverse. The algorithm calls for the evaluation of the Tanimoto lookahead condition (eq 3) at every node visited in the search, so visiting fewer nodes in the compressed tree results in better performance.

9. CONCLUSIONS

A new memory-resident data structure has been described that greatly increased the speed of molecular similarity search in molecular databases, over a wide range of similarity thresholds (0.6–0.9). The *same* data structure can also be used to accelerate inverted screening in those databases. This is very valuable; the same data structure can serve more than one purpose on the same data. Work is going on to apply this data structure to other areas of molecular structure calculations.

ACKNOWLEDGMENT

Mark Olson, Harold Helson, and Andras Furst of CambridgeSoft Inc. are acknowledged for their many valuable discussions of this work.

REFERENCES AND NOTES

- (1) Willett, P.; Barnard, J. M.; Downs, G. M. Chemical Similarity Searching. *J. Chem. Inf. Comput. Sci.* **1998**, 38 (6), 983–996.
- (2) Flower, D. R. On the Properties of Bit String-Based Measures of Chemical Similarity. *J. Chem. Inf. Comput. Sci.* **1998**, 38 (3), 379–386.
- (3) Daylight Chemical Information Systems, 120 Vantis-Suite 550-Aliso Viejo, CA 92656. <http://www.daylight.com/dayhtml/doc/theory/theory.finger.html> (accessed Jan 2009).
- (4) Tanimoto, T. T. *IBM Internal Report*; Nov 17, 1957.
- (5) As a representative of many example papers, see: (a) Bender, A.; Mussa, H. Y.; Glen, R. C.; Reiling, S. Molecular Similarity Searching using Atom Environments, Information-Based Feature Selection and a Naïve Bayesian Classifier. *J. Chem. Inf. Comput. Sci.* **2004**, 44, 170–178.
- (6) Dunkel, M.; Günther, S.; Ahmed, J.; Wittig, B.; Preissner, R. SuperPred: drug classification and target prediction. *Nucleic Acids Res.* **2008**, Jul 1; 36(Web Server issue):W55–9.
- (7) Daylight Chemical Information Systems, 120 Vantis-Suite 550-Aliso Viejo, CA 92656. <http://www.daylight.com/dayhtml/doc/theory/theory.merlin.html> (accessed Jan 2009).
- (8) Swadimass, S. J.; Baldi, P. Bounds and Algorithms for Fast Exact Searches of Chemical Fingerprints in Linear and Sublinear Time. *J. Chem. Inf. Model.* **2007**, 47, 302–317.
- (9) Xue, L.; Stahura, F.; Bajorath, J. Similarity Search Profiling Reveals Effects of Fingerprint Scaling in Virtual Screening. *J. Chem. Inf. Model.* **2004**, 44, 2032–2039.
- (10) Solomon, B. *ChemBioNews*; 2006;16.2. <http://chembionews.cambridgesoft.com/Articles/Default.aspx?articleID=378> (accessed Jan 2009).
- (11) <http://www.merckbooks.com/mindex/index.html> (accessed Jan 2009).
- (12) Schreiman, I. *ChemBioNews*; 2005;15.3. <http://www.cambridgesoft.com/webinars/info/webinarid=80> (accessed Jan 2009).
- (13) Schreiman, I. *ChemBioNews*; 2005;. 13.3. <http://chembionews.cambridgesoft.com/Articles/Default.aspx?articleID=343> (accessed Jan 2009).
- (14) The NCI Database distributed by CambridgeSoft Inc. <http://nci.cambridgesoft.com/> (accessed Jan 2009).
- (15) The TCM Database distributed by CambridgeSoft Inc. <http://www.cambridgesoft.com/databases/details/?db=3> (accessed Jan 2009).
- (16) The Nanogen Database distributed by CambridgeSoft Inc. <http://www.cambridgesoft.com/databases/details/?db=22> (accessed Jan 2009).
- (17) The Ashgate Database distributed by CambridgeSoft Inc. <http://www.cambridgesoft.com/databases/details/?db=2> (accessed Jan 2009).
- (18) <http://www.daylight.com/dayhtml/doc/theory/theory.finger.html> (accessed Jan 2009).
- (19) Furst, A. *Oracle Cartridge 11.0 internal documentation*; CambridgeSoft Inc.
- (20) Bender, A.; Glen, R. C. Molecular similarity: a key technique in molecular informatics. *Org. Biomol. Chem.* **2004**, 2, 3204–3218.
- (21) Bender, A.; Mussa, H. Y.; Glen, R. C.; Reiling, S. Similarity searching of chemical databases using atom environment descriptors: evaluation of performance (MOLPRINT 2D). *J. Chem. Inf. Comput. Sci.* **2004**, 44, 1708–1718.
- (22) Held, G. *Data Compression: Techniques and Applications, Hardware and Software Considerations*; 2nd ed.; John Wiley & Sons: New York, NY, 1987.

CI800325V