

Nearest Neighbor Search in General Metric Spaces Using a Tree Data Structure with a Simple Heuristic

Huafeng Xu* and Dimitris K. Agrafiotis

3-Dimensional Pharmaceuticals, Inc., 665 Stockton Drive, Exton, Pennsylvania 19341

Received July 21, 2003

We present a new algorithm for nearest neighbor search in general metric spaces. The algorithm organizes the database into recursively partitioned Voronoi regions and represents these partitions in a tree. The separations between the Voronoi regions as well as the radius of each region are used with triangular inequality to derive the minimum possible distance between any point in a region and the query and to discard the region from further search if a smaller distance has already been found. The algorithm also orders the search sequence of the tree branches using the estimate of the minimum possible distance. This simple heuristic proves to considerably enhance the pruning of the search tree. The efficiency of the algorithm is demonstrated on several artificial data sets and real problems in computational chemistry.

I. INTRODUCTION

Similarity search is one of the most common tasks in computing.¹ It finds, in a collection, the objects that are most similar to a given object of interest. The collection is often referred to as the database and the given object of interest as the query. Similarity, or rather dissimilarity, is measured by a distance function: two objects are considered more similar if their distance is smaller. Two forms of similarity search exist. The first, referred to as range search, is to find objects whose distances to the query are smaller than a user-specified threshold. The second, referred to as nearest-neighbor search, is to find the k objects whose distances to the query are smaller than any other object in the database. Because the techniques involved in addressing these two types of queries are identical, in this article we will focus only on nearest neighbor search.

Similarity search has found numerous applications in chemical informatics and computational chemistry. For instance, medicinal chemists often want to retrieve from a chemical library molecules that resemble some biologically active lead compound, in the hope that they will exhibit similar biological activity. The inverse of this problem is the design of diverse chemical libraries comprising compounds that are as dissimilar to each other as possible. In this case, nearest neighbor search is used to ensure maximum separation between the members of the library.^{2,3} In some quantitative structure–activity relationship studies, the unknown activity of a compound is estimated based on those of the most similar compounds of known activity.^{4,5} These three types of studies require nearest neighbor searching based on some measure of similarity between chemical compounds. A rather different example is the search for low-energy molecular conformations. To identify unique conformations, a new conformation generated during the sam-

pling process is compared to previously generated conformations and is kept only if it is sufficiently different from its most similar predecessors.

Although the distance functions are defined differently for different problems, they are normally designed to be positive, symmetric, reflexive and satisfy the triangular inequality. Explicitly, if we denote S to be the universe of all valid objects, and the distance to be the function $d: S \times S \rightarrow R$, where R is the set of real numbers, the distance function satisfies the following:

$$\text{Positiveness: } \forall x, y \in S, d(x, y) \geq 0 \quad (1)$$

$$\text{Reflexivity: } \forall x, y \in S, d(x, y) = 0 \leftrightarrow x = y \quad (2)$$

$$\text{Symmetry: } \forall x, y \in S, d(x, y) = d(y, x) \quad (3)$$

$$\text{Triangular inequality: } \forall x, y, z \in S, d(x, y) + d(y, z) \geq d(x, z) \quad (4)$$

A distance function with the above properties is called a metric. The set S furnished with such a distance function is referred to as a metric space and is usually denoted as $(S; d)$. Most similarity searches are carried out in some metric space. To facilitate the discussion and introduce the notations, we also formally define the k -nearest neighbor search problem. Suppose we have a database X consisting of n points from the space S , $X = \{x_i, i = 1, 2, \dots, n \mid x_i \in S\}$. The k -nearest neighbor search problem is defined as follows: for any query point $q \in S$ find the k points from X , $Y = \{y_j, j = 1, 2, \dots, k \mid y_j \in X\}$ such that for any point $y_j \in Y$ and any other point $x \in X$ and $x \notin Y$, we have $d(x, q) \geq d(y_j, q)$. In plain words, we look for the k points in the database X that are closest to the query point q .

The straightforward solution to the nearest neighbor search problem is to compute the distances of the query point to all the points in the database and select the k points that are closest to the query. This solution scales linearly with the number of points in the database and becomes inefficient for large databases. Moreover, there is often a need to

* Corresponding author, current address: Department of Pharmaceutical Chemistry, University of California San Francisco, P.O. Box 2240, San Francisco, CA 94143; phone: (415) 476-6875; e-mail: huafeng@maxwell.ucsf.edu.

repeatedly query the database for the nearest neighbors to a large number of queries. In this case, it would be wasteful to scan the entire database for every query. A more efficient solution would be to organize the database into local neighborhoods so that only a fraction of the database needs to be considered for any single query.

The k -dimensional tree (kd-tree)^{6,7} provides a very efficient solution, if the objects are represented in a low-dimensional Euclidean space. kd-tree recursively partitions the data set into binary subsets along individual dimensions using the coordinates of the points as discriminators. However, this effectively prevents the kd-tree from being used in a general metric space where dissimilarity is not measured by the Euclidean (or generalized L^n) distance between vectors. In addition, the performance of kd-tree rapidly deteriorates as the number of components, or the dimensionality, of the vectors increases. Nevertheless, the partitioning concept of kd-tree also prevails in search algorithms for general metric spaces. These algorithms invariably partition the database into subsets and exploit triangular inequality to eliminate some of the subsets during the search. Most of these solutions were recently encompassed in a unifying model.⁸

Here we use the vantage point tree (vp-tree)⁹ as a prototypical example to illustrate nearest neighbor searching in a general metric space. Like kd-tree, vp-tree recursively partitions the database into subsets. At each level of the tree, a vantage point is selected, and the distances to this vantage point from all the points to be partitioned are computed. The median of these distances is then determined and used to partition the data points into an “internal” and an “external” subset depending on whether their distances to the vantage point are smaller or greater than the median distance, respectively. The partitioning is carried out until each subset contains only a small number of points. To search for the nearest neighbors to a query point, the distance from the query point to the vantage point at the root of the tree is first computed. Using triangular inequality, it is decided if each of the internal and external subsets needs to be searched. In an ideal situation, where only one branch of the tree needs to be searched at every level, the search can be completed in $O(\log N)$ time, which represents a drastic reduction from the $O(N)$ time required by the naive exhaustive search. At the other extreme, if both branches of the tree must be searched at every level, the number of distance computations is N , and vp-tree brings out not savings but additional overhead in constructing and traversing the tree. Ordinarily the performance of vp-tree lies between these two opposite ends, depending on the distribution of the data points, the location of the query, and the intrinsic dimensionality of the database.

Apparently the partitioning need not be restricted to binary divisions. For instance, the data points can be clustered according to their proximity to a few selected pivot points, so that a data point belongs to a cluster if it is closer to the corresponding pivot than to any of the other pivots. This m -ary tree is utilized in the geometric near-neighbor access tree (GNAT).¹⁰ The maximum distance from any point in a cluster to the corresponding pivot, or the radius of the cluster, is stored. The distance from the query point to any point in the cluster is greater or equal to the distance between the query and the pivot minus the radius of the cluster. If this minimum possible distance is larger than the nearest neighbor

distance encountered so far, the respective cluster need not be examined further. Partitioning the space into more clusters reduces the overlap boundary volumes between neighboring clusters and therefore increases, on average, the number of branches that can be pruned during the search. Of course, the increase in the “arity” of the tree is accompanied by an increase in the storage space as well as the number of mandatory distance computations at each level of the tree. Therefore, the arity has to be chosen judiciously to achieve an optimal balance.

In all the tree-based search methods, the efficiency of the search is determined by how often a branch can be pruned from further consideration. In a metric space, the pruning decision is invariably made using the minimum possible distance of the query to any point in the branch under consideration. This minimum possible distance or lower bound is usually estimated using triangular inequality. Clearly, we hope to tighten this bound at our best effort. The larger we can raise this lower bound, the more branches we are likely to prune during the search. In this work, we show, through a data structure similar to GNAT, how to utilize all available information to increase the estimate of the minimum possible distance. We also propose a simple heuristic for ordering the search sequence that further enhances search performance.

II. METHODS

Similar to other search methods for metric spaces, our method organizes the database points into a tree structure based on their mutual proximity. In want of a descriptive and still available name, we call our data structure μ -tree. Each node of the tree, except for the terminal nodes (leaves), has μ child nodes under it. The constant μ , termed the “arity” of the tree, is selected based on the intrinsic dimensionality, as well as the size, of the database. The μ nodes that share the same parent node will be referred to as sibling nodes henceforth. Except for the root node, every node in the tree is associated with a point in the database and stores some auxiliary data, described below, to expedite the search. At each level of the tree, the μ database points associated with the respective μ sibling nodes serve as pivots to partition the space into μ Voronoi regions.¹¹ A point is assigned to a Voronoi region if its distance to the respective pivot is smaller than to any of the other pivots. μ -tree thereby recursively divides the data space into Voronoi regions and clusters the database points based on their neighborhoods. The essential idea of μ -tree is to maximally utilize neighborhood information to derive lower bounds of the minimum possible distance of the query to any data point in a subtree and order the search sequence based on a simple heuristic, so that as many branches of the search tree can be pruned as early as possible. The construction and the search of a μ -tree are described in the following paragraphs.

Construction. We first select μ points from the database, $\{x_j, j = 1, 2, \dots, \mu \mid x_j \in X\}$, and call them the pivot points. We divide all the other points in the database into μ clusters, C_j , so that any point in C_j has x_j as its nearest pivot, i.e., $C_j = \{y \in X \mid \forall i = 1, 2, \dots, \mu \wedge i \neq j, d(y, x_j) \leq d(y, x_i)\}$. We call x_j the pivot point of cluster C_j and write it as x_{C_j} . As we assign each point into its corresponding cluster, we compute its distance to all μ pivot points. We register the smallest

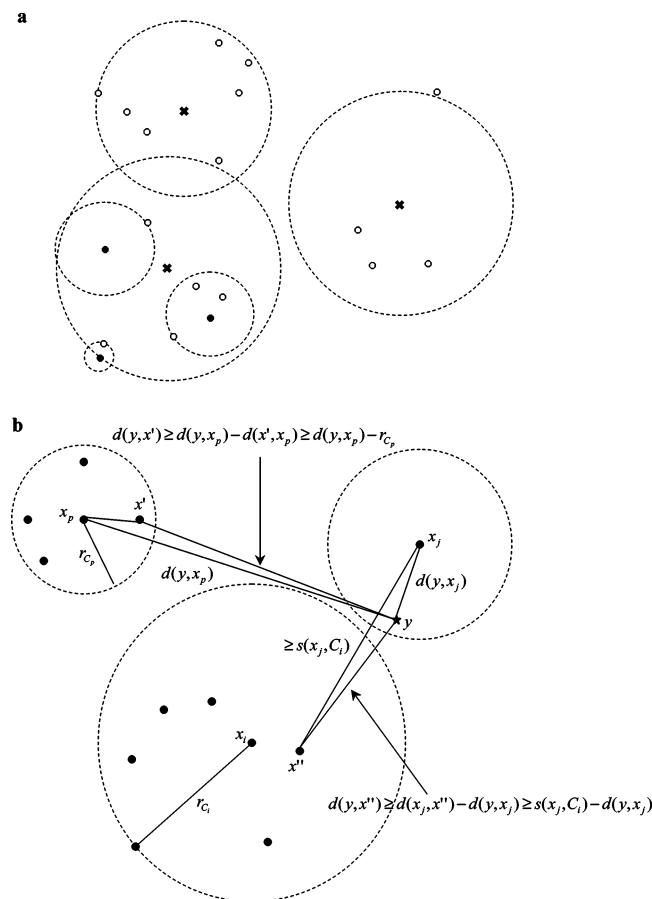


Figure 1. Illustration of the construction and the search of μ -tree. (a) The two-dimensional data set is first partitioned into 3 clusters, using pivots that are accentuated by \times marks. The dashed circles show the radii of the respective clusters. Each of these clusters is then further partitioned into 3 subclusters (only one subdivision is shown), using pivots highlighted by \bullet marks. (b) The minimum possible distance of the query y to any point in a cluster is estimated using triangular inequalities. Depending on whichever is smaller, the minimum possible distance can be estimated using the radius of the cluster (as in the case between y and cluster C_p), or using the separation between the clusters (as in the case between y and cluster C_i).

distance of pivot point x_j to the points in each cluster C_i , $i \neq j$, call that the separation between the reference point x_j and cluster C_i , and denote it as

$$s(x_j, C_i) \equiv \min_{y \in C_i} (d(y, x_j))$$

We also register the largest distance of pivot point x_j to the points in cluster C_j , call that the radius of the cluster, and denote it as

$$r_{C_j} \equiv \max_{y \in C_j} (d(y, x_j))$$

We then repeat the above procedure recursively. For each cluster C , we select μ pivot points, $\{x_j^C\}$, other than the pivot point of the cluster x_C (see above), and divide the other points in the cluster into μ subclusters based on the nearest pivots. A simple illustration of the partitioning process is shown in Figure 1a. We call the resulting μ subclusters child clusters of C and sibling clusters of each other. We keep the radius of each cluster as well as the separations between each pivot point and the sibling clusters. We proceed until no cluster

contains more than μ database points. In the end we obtain a tree structure. Each node in the tree corresponds to a cluster C and contains: (1) the reference point x_C of the cluster, (2) the radius of the cluster r_C , (3) the separation of the reference point from the sibling clusters $s(x_C, C_i)$, $C_i \neq C$, and (4) pointers to the μ child clusters under C .

Slight variants can be devised in the construction of the tree. For example, the arity of the tree need not remain constant—a cluster of large radius containing many points may be divided into more subclusters than a cluster of small radius containing a few points; the recursive division may be terminated by different criteria, such as when a cluster has a very small radius, and therefore its division is unlikely to cause any elimination of its subtrees in the search. These variants will be pursued in future refinements of the algorithm.

Search. To search for the nearest neighbor in the database to a given query point y , we first compute its distances to the μ pivot points at the top level of the tree. The smallest of these serves as an upper bound of the nearest neighbor distance. We then use triangular inequality extensively to derive rigorous lower bounds for the distances of the query point to the points in each cluster and eliminate the clusters with lower bounds bigger than the current nearest neighbor distance from further consideration. Suppose the distance between the query point and the pivot x_j is $d(y, x_j)$. The distance between any point x' in cluster C_j and x_j is no greater than the radius of the cluster, i.e., $d(x', x_j) \leq r_{C_j}$. From the triangular inequality (see Figure 1b for an illustration), we have $d(y, x') \geq d(y, x_j) - d(x', x_j) \geq d(y, x_j) - r_{C_j}$. Therefore the minimum distance between any point in cluster C_j and the query y cannot be smaller than $d(y, x_j) - r_{C_j}$. This is what vp-tree and GNAT use to estimate the lower bounds of the distances and eliminate branches. We point out, however, that further information on the separation between sibling clusters can be used to derive additional and potentially greater lower bounds. Consider the pivot points x_i 's of the sibling clusters. In constructing the tree, we have computed and stored the minimum separation between each pivot point x_j and any point in cluster C_j . For any point $x'' \in C_j$, we have $d(x'', x_i) \geq s(x_i, C_j)$. Using triangular inequality on the triangle of y , x_i , and x'' , we have $d(y, x'') \geq d(x'', x_i) - d(y, x_i) \geq s(x_i, C_j) - d(y, x_i)$ (see Figure 1b for an illustration). Therefore for each $i = 1, 2, \dots, \mu \wedge i \neq j$, triangular inequality yields a lower bound for the distance between the query point and any point in cluster C_j . Including the lower bound derived above from the radius of the cluster, there are μ lower bound estimates in total. The largest one of these is the best lower bound. If this is larger than the nearest neighbor distance found so far, cluster C_j is omitted from further consideration.

From the above procedure of deriving the lower bounds, it is apparent that more clusters are likely to be eliminated from the search if (1) each cluster has smaller radius, and (2) the separation between clusters is larger. As soon as we select μ reference points, assigning other points to the cluster corresponding to the nearest reference point guarantees that the radius of the largest cluster is minimum, and the smallest separation between any two clusters is maximum.

In addition to using the intercluster separations to tighten the estimate of the minimum possible distances, we also realize that by carefully selecting the order in which the remaining branches of the tree are searched, we can increase

Chart 1. μ -Tree**Construction**

```

function makeTree ( database, root )

  select  $\mu$  pivots from database;

  for i = 1; i  $\leq$   $\mu$ ; i++;

    child[i] = new Node(parent = root, point = i'th pivot);
    child[i].radius = 0;

    for j = 1; j  $\leq$   $\mu$ ; j++;

      child[i].sep[j] = distance(pivot i, pivot j);

    end

  end

  foreach x in database AND (x NOT a pivot)

    p = the pivot that is closest to x;
    insert x in subset[p];
    child[p].radius = max(distance(x, pivot p), child[p].radius);

    for i = 1; i  $\leq$   $\mu$ ; i++;

      if i == p continue;

      child[i].sep[p] = min(distance(x, pivot i), child[i].sep[p]);

    end

  end

  for i = 1; i  $\leq$   $\mu$ ; i++;

    makeTree(subset[i], child[i]);

  end

end function

```

Search

```

function [nearest, d] = search ( query, root of  $\mu$ -tree )

  priority_queue <by min_pos_dist> queue;

  // Find the minimum distance to the query among the K pivots.

  foreach c in children of root;

    if distance(query, c.point) < d

      d = distance(query, c.point);
      nearest = c.point;

    end

  end

```

the probability that a smaller distance to the query point is discovered early on in the search. The earlier a small distance is discovered, the more branches will be pruned from the search. A very intuitive heuristic for ordering the branches for search arises from the following reasoning: the lower bound of the minimum distance between the query and any point in the cluster estimated by triangular inequality approaches the actual minimum distance as the points in the database become infinitely dense. This is true because when the space is densely populated, there is always a point on the perimeter of the cluster that lies near the line joining the query point and the pivot of the cluster. Our simple heuristic is thus to first search the branch with the smallest minimum possible distance estimated by the triangular inequalities. This can be trivially implemented using a priority queue. We will demonstrate that this ordering considerably improves the efficiency of the search.

```

end if

sn = new SearchNode(c);

sn.min_pos_dist = distance(query, c.point) - c.radius;
sn.min_pos_dist = max(0, sn.min_pos_dist);

for i = 1; i  $\leq$   $\mu$ ; i++;

  if root.child[i] == c continue;

  est = c.sep[i] - distance(query, root.child[i].point);
  sn.min_pos_dist = max(est, sn.min_pos_dist);

end

if sn.min_pos_dist > d continue;

insert sn in queue;

end

while queue NOT empty

  // Get the node with the smallest lower bound of distance.
  pop sn from queue with minimum min_pos_dist;

  if sn.min_pos_dist > d break;

  foreach c in children of sn

    if distance(query, c.point) < d

      d = distance(query, c.point);
      nearest = c.point;

    end if

    csn = new SearchNode(c);

    csn.min_pos_dist = distance(query, c.point) - c.radius;
    csn.min_pos_dist = max(0, csn.min_pos_dist);

    for i = 1; i  $\leq$   $\mu$ ; i++;

      if sn.child[i] == c continue;

      est = c.sep[i] - distance(query, sn.child[i].point);
      csn.min_pos_dist = max(est, csn.min_pos_dist);

    end

    if csn.min_pos_dist > d continue;

    insert csn in queue;

  end

end while

return [nearest, d];

end function

```

The optimal value of μ is determined empirically by trial and error. With increasing μ , the height of the tree as well as the radius of each cluster decreases, and the separations between the clusters increase, both contributing to a smaller number of distance computations. On the other hand, at least μ distance computations must be performed for every subtree, which makes a large μ value counterproductive. The optimal μ value must come between the two extremes and apparently depends on the intrinsic dimensionality, the distribution of data points, and the size of the database. However, a mathematical expression for the optimum value of μ appears to be too difficult to obtain, and even a simple estimate based on scaling arguments leads to a complicated form hardly of any practical use.

To assist the readers implementing our algorithm, we present the pseudocode for the construction and search of μ -tree in Chart 1.

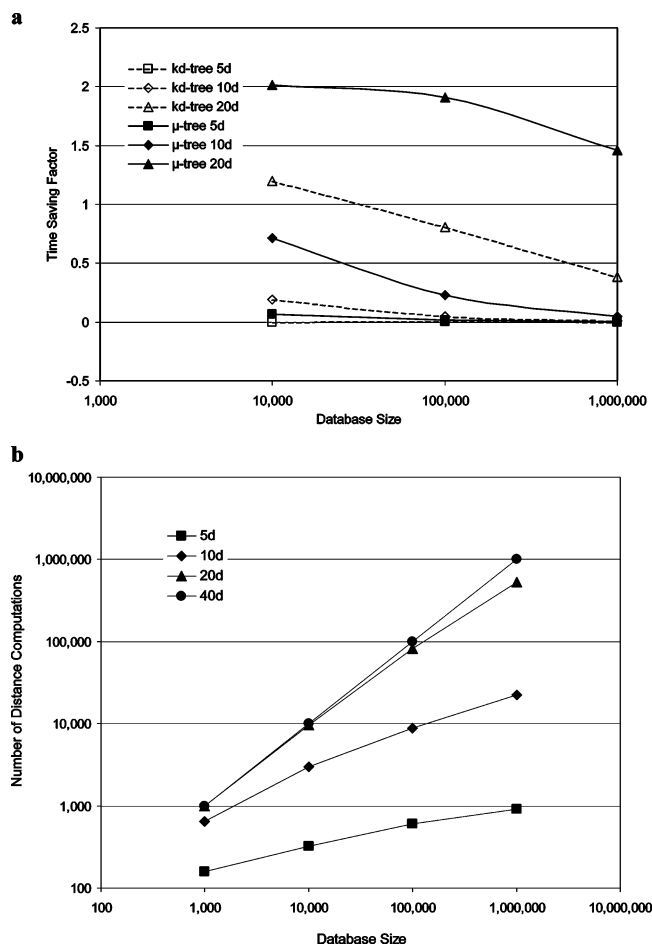


Figure 2. Efficiency of μ -tree in searching for the first nearest neighbor ($k = 1$) of 1000 random queries in a database of random points uniformly distributed in a unit hypercube of 5, 10, 20, and 40 dimensions. The database size and the number of distance computations are displayed on the logarithmic scale for easy viewing. (a) Time saving factors of the search methods as a function of the database size. The time saving factor is defined as the ratio of the retrieval time of the tree-based method to that of a naïve exhaustive search. Both kd-tree and μ -tree achieve drastic savings for low-dimensional data, and the savings get better with increasing database size. Kd-tree outperforms μ -tree significantly for low-dimensional vectorial data. The advantages of the tree-based methods vanish when handling high-dimensional data. (b) Average number of distance computations μ -tree performs in order to search for the nearest neighbor in a database of varying sizes. Sublinear growth is evident for low-dimensional data. But for high-dimensional data, μ -tree is little different from exhaustive search.

III. RESULTS AND DISCUSSION

By trial and error, we determined $m = 5$ to be a good value for all the data sets studied in this work and report results with this value throughout this article. $m = 3, 4, 6, 7$ gave very similar results. The best search efficiency was achieved with m taking one of the above values depending on the data set.

We have previously implemented the kd-tree algorithm for nearest neighbor search.¹² Naturally, we compare the performance of the new data structure, μ -tree, against that of kd-tree. Because kd-tree is only applicable to vectorial data, we first compare the two algorithms using artificially synthesized vectorial data. A set of random data points uniformly distributed in a d -dimensional hypercube was constructed and was used to construct kd-trees and μ -trees

based on the Euclidean distance. Query points were also randomly generated in the same d -dimensional hypercube. The performances of kd-tree and μ -tree are compared in Figure 2a. Here we present the time saving factor for each method, as defined by the ratio of the search time of the method to that of naïve exhaustive search. It is seen that, for uniform vectorial data, kd-tree has better search efficiency than μ -tree. Still, μ -tree also drastically reduces the search time when dealing with low-dimensional data. As expected, both methods lose their advantage in a high-dimensional database. The raw search time, however, can be biased by the implementation and can be affected by the computational cost of the distance function. To investigate how μ -tree does in eliminating search branches, we plot, in Figure 2b, the number of distance computations performed in μ -tree search as a function of database size. It is clear that for low-dimensional data, μ -tree cuts down a tremendous number of distance computations, calculating the distance of the query to only a tiny fraction of the entire database (for the five-dimensional database of 1 million members, μ -tree calculates the distance to less than 0.1% of the data points). This result suggests that for data sets where the cost of distance computation overwhelms the tree operations, the savings from μ -tree can be dramatic.

Because kd-tree partitions the space along alternating input dimensions, its efficiency worsens with increasing *apparent* dimensionality of the data set (i.e., the number of components in each vector). μ -tree, on the other hand, works exclusively in metric space, and thus its performance depends only on the *intrinsic* dimensionality, which is typically lower than the *apparent* one. Simplistically, the intrinsic dimension is one less than the largest number of points that can be placed in the metric space so that they are all equidistant to each other. For example, we can have an equilateral triangle in a plane, but not four points with all six of their pairwise distances equal, therefore the intrinsic dimension of the plane is $3 - 1 = 2$. From this perspective, it is obvious why nearest neighbor search requires more distance computations as the intrinsic dimensionality grows. To illustrate the advantage of μ -tree over kd-tree in this respect, we synthesized a low dimensional data set embedded in a high dimensional vector space. This data set was constructed by generating a set of random points uniformly distributed in a d -dimensional vector space, $\mathbf{x}_i = (x_i^1, x_i^2, \dots, x_i^d)$, and augmenting them with $D - d$ additional dimensions using linear combinations of their first d components

$$X_i^k = \sum_{t=1}^d c_{kt} x_i^t$$

for $k > d$ and

$$X_i^k = x_i^k$$

for $k \leq d$. The coefficients c_{kt} were chosen at random $\{-1 < c_{kt} < 1\}$, $k = d + 1, d + 2, \dots, D$, $t = 1, 2, \dots, d$, where d is the preset intrinsic dimension and D is the apparent vectorial dimension. Apparently, the vectors X_i are in a d -dimensional hypercube embedded in a D -dimensional vectorial space. Once again, Euclidean distances in the D -dimensional space were used as a measure of dissimilarity. The results for kd-tree and μ -tree are compared in Figure 3.

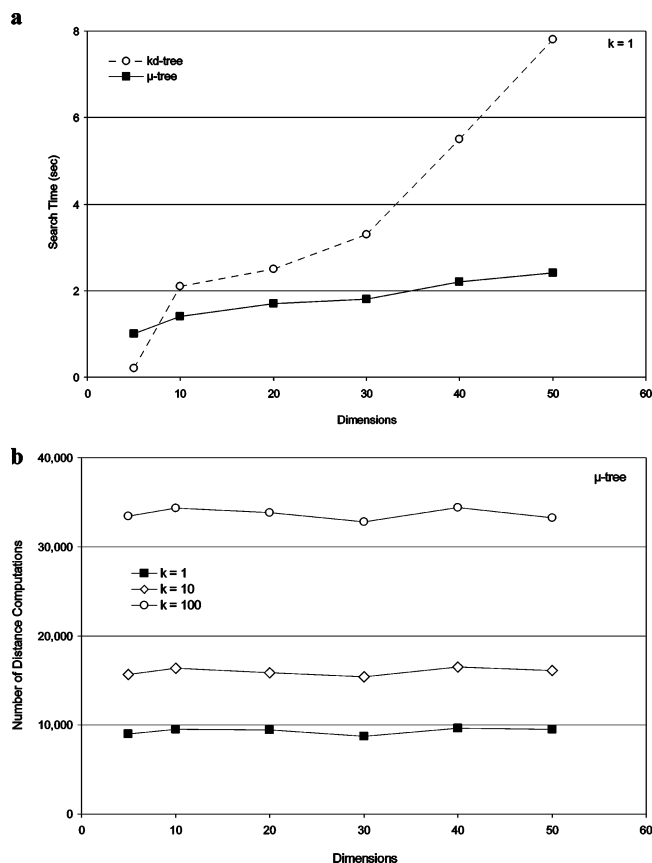


Figure 3. (a) Average retrieval time for the nearest neighbors of 1000 random queries in a database of 1 million random points lying on a five-dimensional manifold embedded in 5, 10, 20, 30, 40, and 50 apparent dimensions, as a function of the apparent dimensions. The limitations of kd-tree and the advantage of μ -tree are apparent. (b) The average number of distance computations μ -tree performs to retrieve the k -nearest neighbors for the above random queries. They remain constant with respect to the apparent dimensions, aside from some statistical fluctuations.

As expected, μ -tree's performance remains the same, whereas kd-tree's rapidly deteriorates with increasing apparent dimensions (the slight increase in the search time for μ -tree is due to the increased cost of computing distances in higher dimensions). This is further illustrated in Figure 3b, which shows that the number of distance computations involved in the k -nearest neighbor retrieval by μ -tree does not change with the apparent dimension.

We next applied μ -tree to some nearest neighbor search problems commonly encountered in chemical informatics and compared its performance to that of kd-tree whenever applicable. In computer-aided design, we often need to select from a large collection of compounds, a small number of candidates that resemble some bioactive lead. As a first example, we conducted nearest neighbor searches of lead compounds in a 1-million-member virtual combinatorial library based on the reductive amination reaction.¹³ Each compound was characterized by 117 topological descriptors,¹⁴ including molecular connectivity indices, kappa shape indices, subgraph counts, information-theoretic indices, Bonchev-Trinajstić indices, and topological state indices. These descriptors were reduced to 26 principal components to eliminate strong linear correlations between them. The distance between any two compounds was defined as the Euclidean distance between the principal components of the

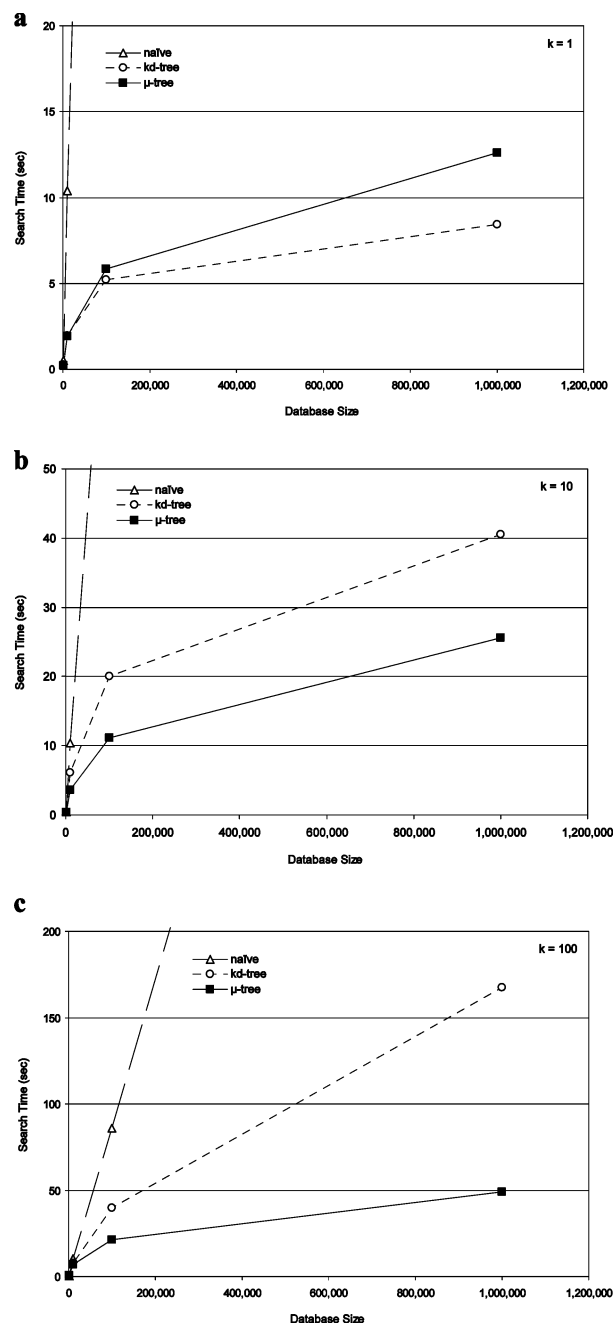


Figure 4. Average retrieval time required by the naive, k-tree, and μ -tree algorithms to retrieve the k nearest neighbors ($k = 1, 10$, and 100) of 1000 random queries in the diamine database, as a function of the database size (number of compounds). Euclidean distance between the vectors formed by the first 26 principal components of the chemical descriptors is used. (a) $k = 1$. (b) $k = 10$. (c) $k = 100$. μ -tree demonstrates a clear advantage over kd-tree as more nearest neighbors are sought.

respective descriptors. This vectorial representation makes it possible to apply the kd-tree algorithm and to compare the performances of kd-tree and μ -tree. The results are presented in Figure 4. It shows that even though kd-tree slightly outperforms μ -tree when only the single nearest neighbor is sought ($k = 1$), μ -tree exhibits higher efficiency for $k > 10$, and its advantage increases with increasing k . We do not fully understand this behavior, and we continue to investigate this issue. Figure 5 shows that μ -tree drastically reduces the number of distance computations involved in the nearest neighbor search.

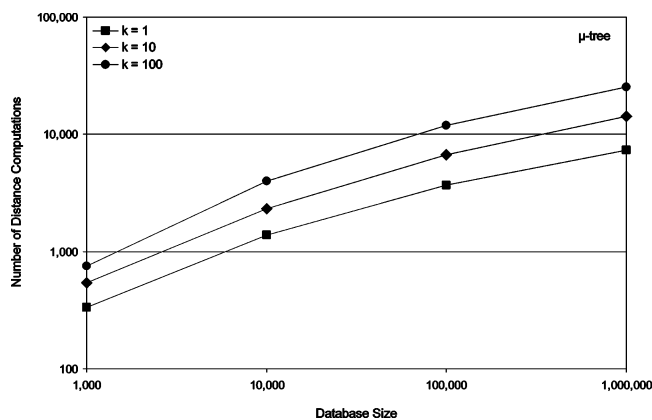


Figure 5. Average number of distance comparisons required by the μ -tree algorithm to retrieve the k nearest neighbors ($k = 1, 10$, and 100) of 1000 random queries in the diamine database, as a function of the database size. Logarithmic scales are used for easier viewing.

A rather different but equally popular set of descriptors are the substructure keys or molecular fingerprints.³ These descriptors are binary vectors, where each bit represents the presence or absence of a predefined structural feature. Here, we used 166-dimensional bitsets, based on the public fragment definitions in the ISIS chemical database management system.¹⁵ The distance between two compounds was measured using the Tanimoto coefficient

$$r_{ij} = 1 - \frac{|(\mathbf{x}_i \& \mathbf{x}_j)|}{|\mathbf{x}_i| |\mathbf{x}_j|}$$

where $\mathbf{x} \& \mathbf{y}$ represents the bitwise AND, $\mathbf{x}|\mathbf{y}$ represents the bitwise OR, and $|\mathbf{x}|$ gives the number of 1 or "on" bits in the bit-vector. In Appendix A we show that the Tanimoto distance is indeed a valid metric. A second, 4-component virtual combinatorial library containing 1 million compounds was constructed according to the Ugi reaction.¹³ The fingerprint for each compound was computed, each comprising 166 bits. Since the fingerprints are not vectorial data, the kd-tree algorithm is no longer applicable. We simply characterized the search efficiency of μ -tree by comparing it to naïve search. The savings that μ -tree brings about are clearly demonstrated in Figure 6. The efficiency of μ -tree becomes increasingly manifest as the database size increases. To retrieve the nearest neighbor from a data set of 1 million members (a typical size of a compound library in a pharmaceutical company), μ -tree computes the distances of the query to only 0.25% of the database members, and is more than 100 times faster than the exhaustive search, as measured by CPU time (naïve search time is off the scale and not shown). The small number of distance computations again suggests that the time savings can still be drastically improved by optimizing the implementation to reduce the cost of tree operations.

The problem of conformational analysis provides another example of nearest neighbor search. A flexible molecule can adopt a large number of conformations. A number of computer programs have been created to sample conformational space and to identify distinct low-energy conformations that are physicochemically important.¹⁶ A common problem in conformational analysis is to identify, in a set of sampled conformations, the one that is most similar to a query

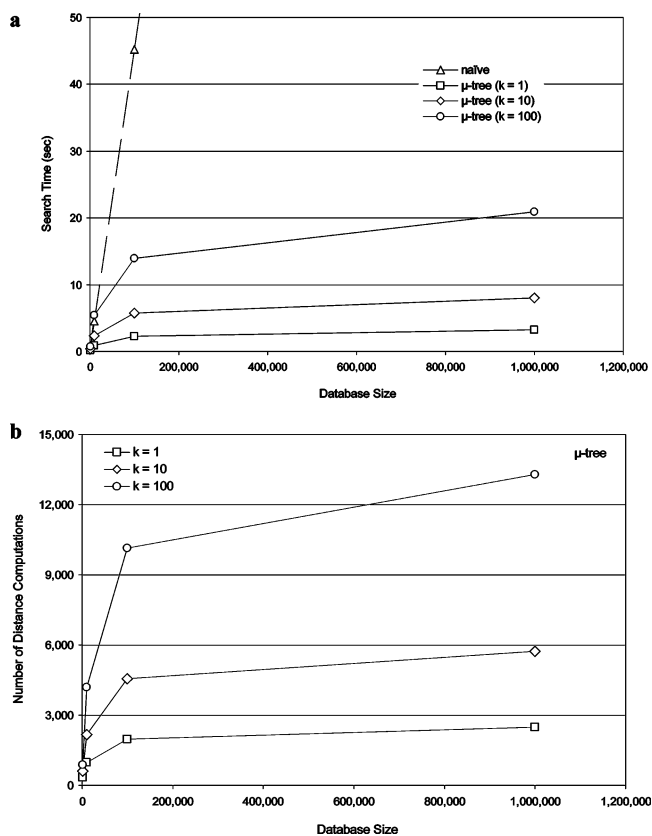


Figure 6. (a) Average CPU time of the naïve and μ -tree algorithms to retrieve the k nearest neighbors ($k = 1, 10$, and 100) of 1000 random queries in the Ugi database, as a function of the database size. Distance based on the Tanimoto coefficients between structural fingerprints is used. Naïve search is only used to retrieve $k = 1$ nearest neighbors. (b) Average number of distance computations required by the μ -tree algorithm to retrieve the k nearest neighbors ($k = 1, 10$, and 100) of 1000 random queries in the Ugi database, as a function of the database size.

structure. The simplest example is to test whether a new conformation is so similar to a previously generated conformation that it can be discarded as a duplicate. The distance between two molecular conformations is conveniently defined as the root-mean-square deviation (RMSD), which we have shown to be a valid metric in a previous publication.¹⁷ Again, since kd-tree is not applicable to this problem, μ -tree is compared to naïve search. The performance of the μ -tree depends on the intrinsic dimensionality of the data set. The intrinsic dimensionality of the conformational space increases with the flexibility of the molecule, which is determined by the number of rotatable bonds, ring puckering, and other factors such as steric hindrance. Four molecules were used in this study: raloxifene, the free base of Gleevec (imatinib mesylate), cycloheptadecane, and [Met⁵]-enkephalin (sequence YGGFM). Eight thousand conformations were generated for each molecule using a self-organizing method¹⁸ and stored in data files. These conformations were *not* locally minimized in energy and therefore were distributed continuously in conformational space. Another 100 conformations were independently generated for each molecule and used as queries to search for the nearest neighbors in the previously saved conformations. The results are summarized in Figure 7. For the modestly flexible molecules raloxifene and Gleevec, μ -tree effects huge savings in the search time. μ -tree computes distances between the query and only about

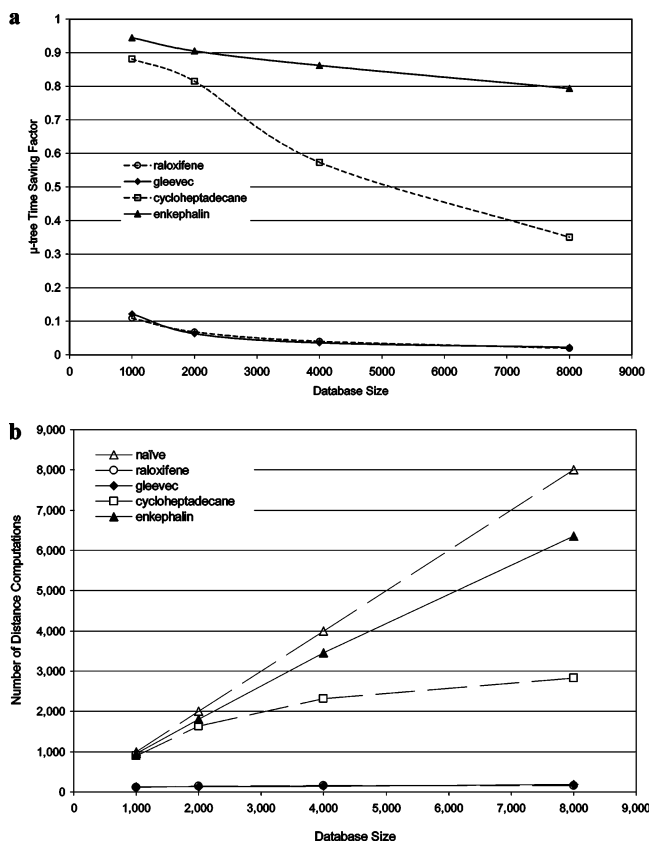


Figure 7. (a) The time saving factors of μ -tree in search of the nearest neighbors in the conformational database of raloxifene, Gleevec, cycloheptadecane, and [Met⁵]-enkephalin. (b) Average number of distance computations required by the naïve and μ -tree algorithms to retrieve the nearest neighbor ($k = 1$) of a random query in the conformational databases, as a function of the database size.

2% of the database conformations. As the flexibility of the molecules, and concomitantly the intrinsic dimensionality of the search space, increases, the advantage of μ -tree diminishes and all but vanishes for a highly flexible molecule such as [Met⁵]-enkephalin. Fortunately, the compounds that are of most interest to medicinal chemists are relatively rigid, usually containing fewer than 10 rotatable bonds. We therefore expect μ -tree to be widely applicable in conformational analysis of small molecules.

Having demonstrated the potential of μ -tree in several computational chemistry applications, we now discuss the merits of μ -tree attributable to the use of the intercluster separations and the simple heuristic introduced to order the search sequence. These are the two main improvements in μ -tree compared to GNAT. We carried out experiments to quantitatively assess each enhancement. We performed nearest neighbor searches in the random data set uniformly distributed in the unit hypercube and in the conformational data sets of raloxifene and cycloheptadecane. We used μ -tree, with either enhancements disabled, to search for the nearest neighbors. The efficiency of the search is measured by the fraction of distances computed, defined as the number of distance computations involved in one search divided by the size of the database. As shown in Table 1, the simple heuristic to first search the cluster with the smallest estimate of minimum possible distance improves the search efficiency by as much as 50%. In comparison, using intercluster separations does not improve the search as much but still

Table 1. Effect of Using the Intercluster Separations To Improve the Estimate of Lower Bounds of Minimum Possible Distances and the Effect of Using the Simple Heuristic To Prioritize the Search Sequence^a

system	<i>N</i>	w/o separation	w/o heuristic	μ -tree
5-d hypercube	1000	0.1716	0.2085	0.1511
	10000	0.0427	0.0624	0.035
	100000	0.0074	0.0122	0.0063
	1000000	0.0011	0.0021	0.0009
10-d hypercube	1000	0.6481	0.7260	0.6685
	10000	0.3024	0.3837	0.2992
	100000	0.0960	0.1316	0.0906
	1000000	0.0213	0.0394	0.0221
raloxifene	1000	0.1227	0.1345	0.116
	2000	0.0778	0.0947	0.0738
	4000	0.0468	0.0641	0.0436
cycloheptadecane	8000	0.0252	0.0443	0.0231
	1000	0.8678	0.8900	0.8678
	2000	0.8362	0.8659	0.8362
	4000	0.7812	0.8299	0.7812
	8000	0.6403	0.7331	0.6402

^a The efficiency of the search is measured by the fraction of distances computed, defined as the number of distance computations involved in searching for the nearest neighbor divided by the database size. The numbers in the column under w/o separation show the fraction of distances computed when μ -tree searches without using the intercluster separations to tighten the estimate of minimum possible distances. The numbers under w/o heuristic show the fraction of distance computed when μ -tree searches the remaining branches using plain breadth-first search. The numbers under μ -tree show the fraction of distance computed when μ -tree searches with both improvements. The searches are carried out in random data sets in 5- and 10-dimensional hypercubes and in the conformational data sets of raloxifene and cycloheptadecane.

generates a consistent and noteworthy reduction of distance computations. Just like the performance of the search method itself, the improvement induced by either enhancement diminishes with increasing dimensionality. The rather small gain in efficiency brought about by using intercluster separations can be understood for the reasons that follow. The lower distance estimate is based on the difference between a cluster's separation from another and the distance of the query to the pivot of the latter cluster. For this to be large enough to eliminate the former cluster, the intercluster separation has to be large and the query has to be close to the pivot of the latter cluster. Due to our ad hoc selection of the pivots, there is significant overlap between the clusters, leading to small intercluster separations. This problem can be remedied by judiciously selecting well-separated pivots. But inherent is the problem that the query is often far away from any pivot points, even more so at high dimensions, consequently yielding a small estimate of the lower distance bound. This limits the benefits of using intercluster separations to database of moderate intrinsic dimensionality. Moreover, computing the intercluster separations increases the construction time of the tree by a factor of μ . Therefore we advocate the use of this technique only for relatively static databases of low intrinsic dimensionality.

Finally, we would like to point out an obvious improvement to μ -tree. Currently, all pivots are randomly selected when constructing the tree. This is done to minimize the construction time of the tree. For many database applications, however, the tree need only be rebuilt infrequently when major modifications to the database occur. Therefore the construction time is usually not of paramount concern. Under such circumstances, the pivots can be selected by sampling

the database to maximize their mutual separations. This will lead to smaller and well-separated clusters, which in turn will help prune the branches more extensively during the search. A number of diversity sampling techniques can be helpful in this regard.^{2,3} This obvious improvement has been incorporated in other tree-based methods. Since the main purpose of this article is to introduce the better estimate of the lower bound of distances and the simple heuristic to order the search sequence and to demonstrate the applicability of μ -tree to genuine problems in computational chemistry, we will not dwell on the issue of the selection of pivots, which is itself an interesting problem worth exploring.

IV. CONCLUSIONS

We presented a new data structure, μ -tree, for fast nearest neighbor search of large databases in general metric spaces. μ -tree utilizes the separation between clusters to derive tighter estimates of minimum possible distances, which improve the pruning of the search tree. It also prioritizes the search sequence using the estimates of minimum possible distances, which accelerates the discovery of small distances. We demonstrated the efficiency of μ -tree using both artificial data and real problems encountered in computational chemistry. We believe μ -tree can be a very useful addition to the toolkit for nearest neighbor searching. Moreover, the new concept of ordering the search sequence based on a simple heuristic may open up new possibilities for developing better nearest neighbor search algorithms. We will actively pursue such possibilities in the future.

ACKNOWLEDGMENT

We thank Dr. F. Raymond Salemme, Victor S. Lobanov, and Sergei Izrailev of 3-Dimensional Pharmaceuticals for many useful discussions.

APPENDIX A

Here we show that the Tanimoto distance is indeed a valid metric. Positiveness, reflexivity, and symmetry are obvious. We will prove that it also satisfies the triangular inequality and therefore is a valid metric. First let us introduce the notations for bitwise operations. In the following, $\&$ stands for bitwise AND, $|$ stands for bitwise OR, \sim stands for bitwise NOT, and $|x|$ is the number of 1 (or 'on') bits in the bit-vector x . The Tanimoto distance between x and y is

$$r(x,y) = \frac{|x|y| - |x \& y|}{|x|y|} = \frac{|x| + |y| - 2|x \& y|}{|x|y|}$$

Consider the "triangle" formed by three bit-vectors x , y , and z . We divide the 1 bits of y into three categories: (1) the 1 bits that overlap with those of both x and z ; (2) the 1 bits that overlap with either those of x or z but not both; (3) the 1 bits that do not overlap with either those of x or z . We therefore write

$$|x|y| = |x| + |y \& z \& \sim x| + |y \& \sim x \& \sim z|$$

$$|y| = |y \& x \& \sim z| + |y \& z \& \sim x| + |y \& x \& z| + |y \& \sim x \& \sim z|$$

and

$$|x \& y| = |x \& y \& \sim z| + |x \& y \& z|$$

We have

$$\begin{aligned} r(x,y) &= \frac{|x| - |y \& x \& \sim z| + |y \& z \& \sim x| - |x \& y \& z| + |y \& \sim x \& \sim z|}{|x| + |y \& z \& \sim x| + |y \& \sim x \& \sim z|} \\ &\geq \frac{|x| - |y \& x \& \sim z| + |y \& z \& \sim x| - |x \& y \& z|}{|x| + |y \& z \& \sim x|} \\ &\geq \frac{|x| - |y \& x \& \sim z| + |y \& z \& \sim x| - |x \& y \& z|}{|x| + |z \& \sim x|} \\ &= \frac{|x| - |y \& x \& \sim z| + |y \& z \& \sim x| - |x \& y \& z|}{|x|z|} \end{aligned}$$

The first inequality is true because $r(x,y) \leq 1$, and therefore the fraction decreases when both the denominator and the numerator are subtracted by the same nonnegative $|y \& \sim x \& \sim z|$.

We thus have

$$\begin{aligned} r(x,y) + r(y,z) &\geq \frac{|x| - |y \& x \& \sim z| + |y \& z \& \sim x| - |x \& y \& z|}{|x|z|} + \frac{|z| - |y \& z \& \sim x| + |y \& x \& \sim z| - |x \& y \& z|}{|x|z|} \\ &= \frac{|x| + |z| - 2|x \& y \& z|}{|x|z|} \\ &\geq \frac{|x| + |z| - 2|x \& z|}{|x|z|} = r(x,z) \end{aligned}$$

REFERENCES AND NOTES

- (1) Apers, P. M. G.; Houtsma, M. A.; Blanken, H. T. *Multimedia Databases in Perspective*; Springer-Verlag: New York, 1997.
- (2) Agrafiotis, D. K. *IBM J. Res. Dev.* **2001**, 45, 545.
- (3) Agrafiotis, D. K. The diversity of chemical libraries. In *Encyclopedia of computational chemistry*; Schleyer, P. v. R. et al., Ed.; John Wiley & Sons: Chichester, 1998; p 742.
- (4) Parzen, E. *Ann. Math. Stat.* **1962**, 33, 1065.
- (5) Cedeno, W.; Agrafiotis, D. K. *Int. J. Comput. Res.* **2002**, 11, 443.
- (6) Bentley, J. L. *Comm. ACM* **1975**, 18, 509.
- (7) Friedman, J. H.; Bentley, J. L.; Finkel, R. A. *ACM Trans. Math. Soft.* **1977**, 3, 209.
- (8) Chávez, E.; Navarro, G.; Baeza-Yates, R.; Marroquín, J. L. *ACM Computing Surveys* **2001**, 33, 273.
- (9) Yianilos, P. N. Data structures and algorithms for nearest neighbor search in general metric spaces; Proceedings of the Fifth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), 1993.
- (10) Brin, S. Near neighbor search in large metric spaces; Proceedings of the 21st conference on very large databases, 1995.
- (11) de Berg, M.; van Kreveld, M.; Overmars, M.; Schwartz, O. *Computational Geometry: Algorithms and Applications*; Springer-Verlag: New York, 2000.
- (12) Agrafiotis, D. K.; Lobanov, V. S. *J. Chem. Inf. Comput. Sci.* **1999**, 39, 51.
- (13) Agrafiotis, D. K.; Lobanov, V. S. *J. Comput. Chem.* **2001**, 22, 1712.
- (14) Kier, L. B.; Hall, L. H. *Molecular connectivity in structure-activity analysis*; Wiley: New York, 1986.
- (15) MDL Information Systems. MDL Information Systems; MDL Information Systems.
- (16) Leach, A. R. A survey of methods for searching the conformational space of small and medium-sized molecules. In *Reviews in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH: New York, 1991; Vol. 2.
- (17) Agrafiotis, D. K.; Xu, H. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, 99, 15869.
- (18) Xu, H.; Izrailev, S.; Agrafiotis, D. K. *J. Chem. Inf. Comput. Sci.* **2003**, 43, 1186.