# An Intersection Inequality Sharper than the Tanimoto Triangle Inequality for Efficiently Searching Large Databases

Pierre Baldi and Daniel S. Hirschberg*

School of Information and Computer Sciences, Institute for Genomics and Bioinformatics, University of California, Irvine, Irvine, California 92697-3435

Bounds on distances or similarity measures can be useful to help search large databases efficiently. Here we consider the case of large databases of small molecules represented by molecular fingerprint vectors with the Tanimoto similarity measure. We derive a new intersection inequality which provides a bound on the Tanimoto similarity between two fingerprint vectors and show that this bound is considerably sharper than the bound associated with the triangle inequality of the Tanimoto distance. The inequality can be applied to other intersection-based similarity measures. We introduce a new integer representation which relies on partitioning the fingerprint components, for instance by taking components modulo some integer $M$ and reporting the total number of 1-bits falling in each partition. We show how the intersection inequality can be generalized immediately to these integer representations and used to search large databases of binary fingerprint vectors efficiently.

## INTRODUCTION

Bounds on distances or similarity measures can be useful for searching large databases.[1−5] In a fairly typical similarity search situation, if item $\mathcal{A}$ is close to item $\mathcal{B}$ and item $\mathcal{C}$ is far from item $\mathcal{A}$, then one would like to infer that $\mathcal{C}$ is also far from $\mathcal{B}$, thereby allowing one to prune the search space and ignore $\mathcal{B}$. In particular, we consider the typical chemoinformatics situation where a large database of molecules represented by binary fingerprint vectors of fixed length $N$ must be searched using the Tanimoto similarity measure.[6−9] Given a molecule $\mathcal{A}$, we denote its fingerprint vector by $\vec{A}$. The Tanimoto similarity measure is given by

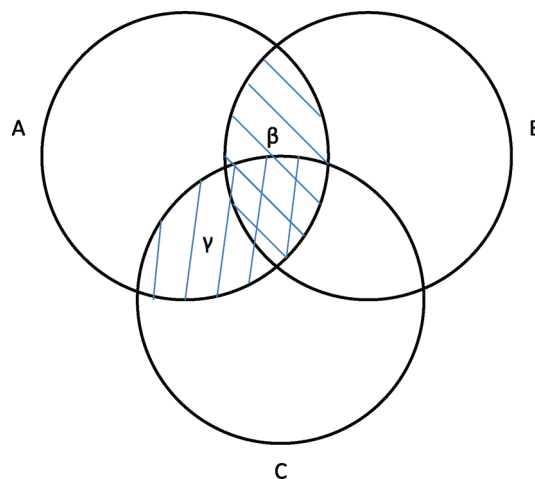$$T(\mathcal{A}, \mathcal{B}) = T(\vec{A}, \vec{B}) = \frac{A \cap B}{A \cup B} \qquad (1)$$

Here and everywhere else we use $A$, $A{\cap}B$, $A{\cup}B$ ⋯ to denote the number of 1-bits in the vectors (or sets) $\vec{A}$, $\vec{A}{\cap}\vec{B}$, $\vec{A}{\cup}\vec{B}$ ⋯. It is well-known that $D(\vec{A}, \vec{B}) = 1 - T(\vec{A}, \vec{B})$ is a distance.[10,11] Thus, using the triangle inequality we can bound the distance $D(\vec{B}, \vec{C})$ by

$$|D(\vec{A}, \vec{B}) - D(\vec{A}, \vec{C})| \le D(\vec{B}, \vec{C}) \le D(\vec{A}, \vec{C}) + D(\vec{A}, \vec{B}) \qquad (2)$$

from which one obtains the following bound on the similarity $T(\vec{B}, \vec{C})$

$$T(\vec{B}, \vec{C}) \le 1 - |T(\vec{A}, \vec{B}) - T(\vec{A}, \vec{C})| \qquad (3)$$

Our goal here is to first derive another general bound on $T(\vec{B}, \vec{C})$ and then study its relationship to the triangle inequality. We then further generalize this bound and show how it can be applied to efficiently search large databases of chemical fingerprint vectors.



**Figure 1.** Venn diagram. The three sets represent the 1-bit components of three fingerprint vectors associated with three molecules $\mathcal{A}$, $\mathcal{B}$, and $\mathcal{C}$. $\beta = A{\cap}B$ and $\gamma = A{\cap}C$. We are interested in estimating or bounding the Tanimoto similarity between $\vec{B}$ and $\vec{C}$, as a function of $A$ and its similarity to $\vec{B}$ and $\vec{C}$.

## THE INTERSECTION BOUND

As can be seen in the Venn diagram of Figure 1, the following relation holds

$$\vec{B} \cap \vec{C} = (\vec{B} \cap \vec{C} \cap \vec{A}) \cup [(\vec{B} - \vec{A}) \cap (\vec{C} - \vec{A})] \qquad (4)$$

By letting $A{\cap}B = \beta$ and $A{\cap}C = \gamma$

$$|\vec{B} \cap \vec{C} \cap \vec{A}| \le \min(A \cap B, A \cap C) = \min(\beta, \gamma) \qquad (5)$$

and

* Corresponding author e-mail: pfbaldi@ics.uci.edu.

$$|(\vec{B} - \vec{A}) \cap (\vec{C} - \vec{A})| \leq \min[B - (A \cap B), C - (A \cap C)] = \min(B - \beta, C - \gamma) \quad (6)$$

Combining these inequalities, we get

$$B \cap C \leq \min(A \cap B, A \cap C) + \min[C - (A \cap C), B - (A \cap B)] = \min(\beta, \gamma) + \min(B - \beta, C - \gamma) \quad (7)$$

Since $T(\vec{B}, \vec{C}) = (B \cap C)/(B \cup C) = (B \cap C)/(B + C - (B \cap C))$, we finally have

$$T(\vec{B}, \vec{C}) \leq \frac{\min(\beta, \gamma) + \min(B - \beta, C - \gamma)}{B + C - \min(\beta, \gamma) - \min(B - \beta, C - \gamma)} \quad (8)$$

which we call the intersection bound or intersection inequality.

## THE INTERSECTION INEQUALITY IS SHARPER THAN THE TRIANGLE INEQUALITY

In this section, we prove that the intersection inequality (eq 8) is sharper than the triangle inequality (eq 3). More specifically we have the following theorem.

**Theorem.** For any fingerprints $\vec{A}$, $\vec{B}$, and $\vec{C}$, we have the inequality

$$T(\vec{B}, \vec{C}) \leq \frac{\min(\beta, \gamma) + \min(B - \beta, C - \gamma)}{B + C - \min(\beta, \gamma) - \min(B - \beta, C - \gamma)} \leq 1 - \left| \frac{\beta}{A \cup B} - \frac{\gamma}{A \cup C} \right| \quad (9)$$

Furthermore the inequality is strict, except in degenerate cases detailed in the proof where (1) at least one of the fingerprints is equal to $\vec{0}$; or (2) $\vec{A} \subset \vec{B}$ or $\vec{A} \subset \vec{C}$; or (3) $\vec{B}$ and $\vec{C}$ are symmetric with respect to $\vec{A}$ (i.e., $B = C$ and $\beta = \gamma$). Furthermore, if two of the fingerprints are identical ($\vec{A} = \vec{B}$ or $\vec{B} = \vec{C}$ or $\vec{A} = \vec{C}$), then the two bounds are equal and exact.

**Proof.** First note that $B$ and $C$ play entirely symmetric roles, and this symmetry can be used to reduce the number of cases that need to be examined. Second consider situations where some of the denominators in eq 9 are equal to 0. If $A \cup B = 0$, then $A = B = 0$ and similarly for $A \cup C = 0$. The denominator on the left-hand side is 0 if and only if $B = 0$ or $C = 0$. Cases where $A = 0$, or $B = 0$, or $C = 0$ can be examined directly and are summarized in Tables 1 and 2.

Thus in the rest of the proof we can assume that $A \neq 0$, $B \neq 0$, and $C \neq 0$ and that none of the denominators in eq 9 is equal to 0. To aggregate the fractions in the right-hand side of eq 9, it is convenient to parametrize the problem by assuming that $A \cup B = k(A \cup C)$ where $k$ is a strictly positive real number. Without any loss of generality, we can assume that $k \geq 1$. By the definition of $k$, we have $A + B - \beta = k(A + C - \gamma)$. This can be rewritten as

$$B - C = (k - 1)(A + C) + \beta - k\gamma \quad (10)$$

and also as

$$B - \beta = (k - 1)(A \cup C) + C - \gamma \quad (11)$$

Equation 11 shows that under our assumptions $B - \beta \geq C - \gamma$, with equality possible if and only if $k = 1$. By noticing that $A + C \geq 2\gamma$, we also have

$$B - C \geq k\gamma - 2\gamma + \beta \quad (12)$$

**Table 1.** Table of Possibilities with $A = 0^a$

|  | $B = 0$ | $B \neq 0$ |
|---|---|---|
| $C = 0$ | (1,1) | (0,0) |
| $C \neq 0$ | (0,0) | ((min $(B,C)$)/(max $(B,C)$),1) |

$^a$ In each case, the left entry is the value of the bound given by the left-hand side of eq 9, and the right-hand side is the bound resulting from the triangle inequality (right-hand side of eq 9). Note that when $B \neq 0$ and $C \neq 0$, the two bounds are identical if and only if $B = C = 0$. Otherwise, if $B \neq C$, then the intersection bound is strictly better than the triangle inequality.

**Table 2.** Table of Possibilities with $A \neq 0^a$

|  | $B = 0$ | $B \neq 0$ |
|---|---|---|
| $C = 0$ | (1,1) | (0,1 − (β/(A∪B))) |
| $C \neq 0$ | (0,1 − (γ/(A∪C))) | main case |

$^a$ In each case, the left entry is the value of the bound given by the left-hand side of eq 9, and the right-hand side is the bound resulting from the triangle inequality (right-hand side of eq 9). Note that when $C = 0$ and $B \neq 0$, the bounds are identical if and only if $\vec{A} = \vec{B}$. Otherwise, if $\vec{A} \neq \vec{B}$, then the intersection bound is strictly better than the triangle inequality. Similar considerations hold for the symmetric case ($B = 0$ and $C \neq 0$).

Now we can write eq 9 as

$$\frac{\min(\beta, \gamma) + \min(B - \beta, C - \gamma)}{B + C - \min(\beta, \gamma) - \min(B - \beta, C - \gamma)} \leq 1 - \frac{|\beta - k\gamma|}{A \cup B} = 1 - \frac{\max(\beta, k\gamma) - \min(\beta, k\gamma)}{A \cup B} \quad (13)$$

By assumption here none of the denominators is equal to 0, and we have seen that $\min(B - \beta, C - \gamma) = C - \gamma$. Thus after reducing to common denominators and rearranging the terms, we need to prove the equivalent relation given by

$$[B - \min(\beta, \gamma) + \gamma][\max(\beta, k\gamma) - \min(\beta, k\gamma)] \leq [A \cup B][B - C + 2\gamma - 2\min(\beta, \gamma)] \quad (14)$$

Three cases can be distinguished depending on the size of $\beta$ relative to $\gamma$ and $k\gamma$. It is convenient to treat the special symmetric case $k = 1$ separately, and we begin with that.

**Case 1: $k = 1$.** In this case, eq 10 yields $B - \beta = C - \gamma$. Two subcases can be considered.

• In the most symmetric subcase, $\gamma = \beta$ which implies that $B = C$, and both bounds are equal to 1.

• Alternatively, without any loss of generality, assume that $\gamma < \beta$. Then eq 14 becomes $B(\beta - \gamma) \leq (A \cup B)(B - C)$. Since here $B - C = \beta - \gamma$, the inequality simplifies to $B \leq A \cup B$ which is always true. Furthermore, equality is achieved only when $\vec{A} \subset \vec{B}$ in which case both bounds are equal to $C/B$ and are not necessarily exact.
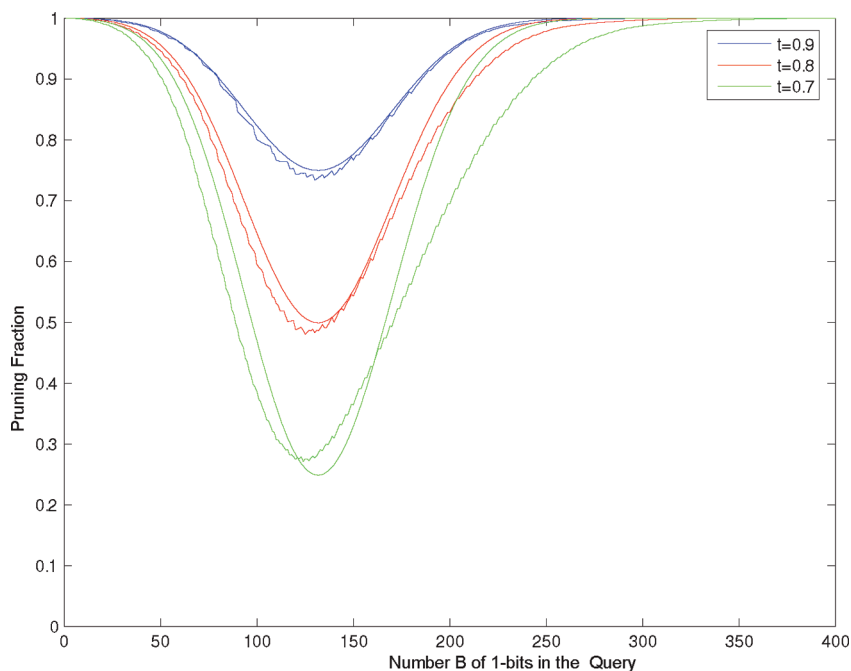
Thus in the rest of the proof, we can assume that $k > 1$ (in addition to $A \neq 0$, $B \neq 0$, and $C \neq 0$).

**Case 2: $k\gamma \leq \beta$.** In this case, eq 14 yields

$$B[\beta - k\gamma] \leq [A \cup B][B - C] \quad (15)$$

This inequality is true because $B \leq A \cup B$, and $\beta - k\gamma \leq B - C$ by eq 10. Under the current assumptions, equality in eq 15 cannot be achieved since equality could occur only in the following two situations:

• $\beta = k\gamma$ and $B = C$. But then eq 10 would imply $k = 1$, contradicting the assumption $k > 1$.

**Figure 2.** Average database pruning for the case of $M = 1$. The $y$-axis represents the fraction of the database being pruned. The $x$-axis represents the number $B$ of bits in the query molecule. Curves correspond to three different Tanimoto similarity threshold. Blue corresponds to $t = 0.9$, red to $t = 0.8$, and green to $t = 0.7$. Results are computed using binary fingerprints of length $N = 512$ using lossy OR-folding compression and a subset of 50,000 molecules extracted randomly from the ChemDB database. Smooth curves correspond to theoretical predictions. Rough curves correspond to simulations.

• $B = A \cup B$ and $\beta - k\gamma = B - C$, which also implies $k = 1$.

**Case 3:** $\gamma \leq \beta \leq k\gamma$. In this case, eq 14 yields

$$B[k\gamma - \beta] \leq [A \cup B][B - C] \qquad (16)$$

Again this inequality is true because $B \leq A \cup B$, and $k\gamma - \beta \leq B - C$ by eqs 10 and 12. Equality between the bounds in eq 16 can only be achieved when $\vec{A} = \vec{C} \subset \vec{B}$. In this case, both bounds are exact and equal to $C/B$. To see this, under the current assumptions, equality could occur only in the following two situations:

• $\beta = k\gamma$ and $B = C$. But then eq 10 would imply $k = 1$, contradicting the assumption $k > 1$.

• $B = A \cup B$ and $k\gamma - \beta = B - C$. But then $\vec{A} \subset \vec{B}$ and eq 12 would imply $\gamma \geq \beta$, hence $\beta = \gamma = A$. In turn, this would yield $B = A \cup B = k(A \cup C) = kC$ which combined with $k\gamma - \beta = B - C$ forces $C = \gamma$ and therefore $\vec{C} = \vec{A} \subset \vec{B}$. In this case, we have $T(\vec{B}, \vec{C}) = C/B$, and this is the value assumed by both the intersection inequality and the triangle inequality and they are both exact.

**Case 4:** $\beta \leq \gamma$. In this case, eq 14 yields

$$[B - \beta + \gamma][k\gamma - \beta] \leq [A \cup B][B - C + 2\gamma - 2\beta] \qquad (17)$$

Again this inequality is true because $B - \beta + \gamma \leq A \cup B$, and $k\gamma - \beta \leq B - C + 2\gamma - 2\beta$ by eqs 10 and 12. Equality between the bounds in eq 16 can only be achieved when $\vec{A} = \vec{C}$. In this case, both bounds are exact and equal to $T(\vec{B}, \vec{C}) = T(\vec{A}, \vec{B}) = \beta/(A \cup B) = \beta/(A \cup C)$. To see this, under the current assumptions, equality could occur only in the following two situations:

• $B - C = 2\beta - 2\gamma$. Using eq 12, this would imply $\beta = k\gamma$ which is possible here only if $\beta = \gamma = 0$ and $B = C$. But this implies $A \cup B = A \cup C$ which contradicts $k > 1$.

• $A \cup B = B - \beta + \gamma$ and $k\gamma - \beta = B - C + 2\gamma - 2\beta$. The first equality yields $A = \gamma$, hence $\vec{A} \subset \vec{C}$. Combining this with $A \cup B = k(A \cup C)$ gives $C = \gamma$, and so $\vec{A} = \vec{C}$. When $\vec{A} = \vec{C}$, the intersection and triangle inequality bounds are identical and exact and equal to $\beta/(A \cup B)$.

## ADDITIONAL PROPERTIES AND GENERALIZATION OF THE INTERSECTION INEQUALITY

The following properties can be derived immediately.

**Property 1.** The intersection bound does not depend on any bits in $\vec{A}$ that are not in common with $\vec{B}$ or $\vec{C}$.

**Property 2.** The intersection bound can be applied to many other similarity measures $S(\vec{A}, \vec{B})$ for binary fingerprints. Although the Tanimoto measure is by far the most widely used measure, two dozen or so other measures have been described in the literature.[12] All these measures involve algebraic expressions that explicitly depend on the intersection $A \cap B$. Thus, with the proper adjustments, bounds on the intersection naturally lead to bounds on these other similarity measures.

**Property 3.** In general, the intersection inequality is considerably sharper than the triangle inequality. We have seen in the theorem above that the intersection bound is always better than the triangle inequality, except in some degenerate cases. But is it much better and worth using? The answer to this question is clearly yes as illustrated by the following two examples.

**Example 1.** Consider the case where the three sets $\vec{A}$, $\vec{B}$, and $\vec{C}$ are disjoint, and $A \neq 0$. Then eq 9 gives

$$T(\vec{B}, \vec{C}) \leq \frac{\min(B, C)}{\max(B, C)} \leq 1 \qquad (18)$$

In other words, in this case the triangle inequality bound is useless and equal to 1. The intersection bound is much

EFFICIENT SEARCHING OF LARGE DATABASES

*J. Chem. Inf. Model., Vol. 49, No. 8, 2009* **1869**

sharper and can diverge maximally from the triangle inequality bound by being as small as 0 when $\min(B,C) = 0$ and $\max(B,C) \neq 0$.

**Example 2.** Consider now a typical case (e.g., Daylight fingerprint system with OR-folding lossy compression) where fingerprints of length $N = 512$ are used. Depending on various implementation details, the average number of 1-bits per fingerprint vector could be something like 250. Consider then two fairly typical fingerprints with $B = 150$ and $C = 350$. Note that to derive a bound on $T(\vec{B},\vec{C})$, we are free to choose the vector $\vec{A}$ (and only the vectors of all 0-bits or all 1-bits can be used to compute bounds, without making additional assumptions on $\vec{B}$ and $\vec{C}$). Consider the choice $\vec{A} = \vec{0}$. Then eq 9 yields

$$T(\vec{B}, \vec{C}) \le \frac{150}{350} \le 1 \tag{19}$$

Thus in this fairly typical case the triangle inequality gives a useless bound equal to 1, whereas the intersection inequality gives a much sharper bound equal to 3/7. If, for instance, $B$ was the query and we were interested in finding molecules with a Tanimoto similarity of at least 0.5, the intersection inequality would allow us to immediately discard the molecule $\vec{C}$ from the search, without having to compute precisely the value of $T(\vec{B},\vec{C})$. This is of course also the basic idea used in ref 5 and generalized below in the Section on Applications to Efficient Database Searches.

GENERALIZED TRIANGLE INEQUALITY

Let us use the notation $\vec{A}^C$ to denote the complement of $\vec{A}$. i.e. $\vec{A}^C$ is the vector obtained by replacing the 0-bits of $\vec{A}$ with 1-bits and the 1-bits of $\vec{A}$ with 0-bits. Then, with the usual notation, eq 7 can also be written as

$$B \cap C \le \min(A \cap B, A \cap C) + \min(A^C \cap B, A^C \cap C) \tag{20}$$

More generally, if $A_1, \cdots, A_M$ denotes any family of disjoint subsets of components ($A_i \cap A_j = 0$ for any $i$ and $j$), then we have immediately the generalized triangle inequality

$$B \cap C \le \sum_{i=1}^{M} \min(A_i \cap B, A_i \cap C) \tag{21}$$

In particular, this is true for any partition $A_1, \cdots, A_M$ of the entire set of components. Thus one important question is how to choose the partition in order to get good bounds. Below we further investigate the case where the partition is obtained by considering the fingerprint components modulo some integer $M$.

APPLICATIONS TO EFFICIENT DATABASE SEARCHES: THE MODULO HASHING APPROACH

Let $M$ be any integer between 1 and $N$ and assume that to any molecule $\mathscr{A}$ and binary fingerprint $\vec{A}$, we also attach a new integer vector of counts $\vec{a} = (a_1, \cdots, a_M)$ of length $M$. The count $a_i$ represents how many 1-bits in $\vec{A}$ fall on components that are congruent to $i$ modulo $M$. For instance, if $\vec{A} = (1,0,1,1,0,0)$ and $M = 3$, then $\vec{a} = (2,0,1)$ because there are two 1-bits in positions congruent to 1 modulo 3 (first and fourth), zero 1-bits in positions congruent to 2 modulo 3 (second and fifth), and one 1-bit in positions congruent to 3 modulo 3 (third and sixth).

If we consider two fingerprint vectors $\vec{B}$ and $\vec{C}$, with their respective count vectors $\vec{b}$ and $\vec{c}$, then $\vec{B}$ and $\vec{C}$ can have at most $\min(b_i, c_i)$ bits in common at positions that are congruent to $i$ modulo $M$. Thus in this case the intersection inequality of eq 21 can be written

$$B \cap C \le \sum_{i=1}^{M} \min(b_i, c_i) \tag{22}$$

resulting in the following bound on the Tanimoto similarity

$$T(\vec{B}, \vec{C}) \le \frac{\sum_{i=1}^{M} \min(b_i, c_i)}{B + C - \sum_{i=1}^{M} \min(b_i, c_i)} \tag{23}$$

Suppose that $\vec{B}$ is the query fingerprint and that we are interested in retrieving from the database all the molecules $\vec{C}$ satisfying $T(\vec{B},\vec{C}) > t$, i.e. with a Tanimoto similarity of at least $t$ to the query. Then, using the intersection bound, we can discard from the search all the molecules $\vec{C}$ satisfying

$$\frac{\sum_{i=1}^{M} \min(b_i, c_i)}{B + C - \sum_{i=1}^{M} \min(b_i, c_i)} \le t \tag{24}$$

or equivalently

$$\sum_{i=1}^{M} \min(b_i, c_i) \le \frac{t}{1 + t}(B + C) \tag{25}$$

Thus, in summary, if we are searching a database of molecules with a query $\vec{B}$ and a Tanimoto similarity threshold $t$, we can prune the search by removing all the molecules that satisfy eq 25 and are thereby guaranteed to have a degree of similarity to $\vec{B}$ that is below the threshold $t$. To show that this corresponds to a very significant amount of pruning, we consider in more details the cases of $M = 1$ and $M = 2$.

**Case M = 1.** This is exactly the case described and used in ref 5. In this case, eq 25 implies that we can discard all the molecules $\vec{C}$ with

$$C \le tB \text{ or } C \ge B/t \tag{26}$$

Using the same example as above with fingerprints of length $N = 512$, consider a typical query $\vec{B}$ with $B = 300$ and a similarity threshold of $t = 0.8$. Then, using the intersection inequality, we can immediately discard from the search all the molecules $\vec{C}$ with $C \le 240 = 0.8 \times 300$ or $C \ge 375 = 300/0.8$.

More generally, Figure 2 shows the significant amount of database pruning achieved by the intersection inequality with $M = 1$ in a practical setting, for different Tanimoto similarity thresholds $t$ and different queries. The curves are derived using a random subset of 50,000 molecules from the ChemDB database[8,13] using Daylight-style fingerprints of length 512 (with folded-OR lossy compression). For example, for a similarity threshold of $t = 0.8$ and a query containing

$B = 200$ 1-bits, on average 80% of the database can be pruned from the search.

**Case M = 2.** The case $M = 2$ subsumes the case $M = 1$ and allows one to prune additional molecules. In this case, we have $\vec{b} = (b_1, b_2)$ and $\vec{c} = (c_1, c_2)$. Here $b_1$ represents the number of 1-bits in $\vec{B}$ falling on odd-numbered components, $b_2$ the number of 1-bits in $\vec{B}$ falling on even-numbered components, and similarly for $\vec{c}$. Obviously, $b_1 + b_2 = B$ and $c_1 + c_2 = C$. When $M = 2$, the intersection bound allows one to prune all the molecules $\vec{C}$ such that

$$\min(b_1, c_1) + \min(b_2, c_2) \le \frac{t}{1+t}(B + C) \quad (27)$$

This yields four subcases:

1. $b_1 + b_2 = B \le t(B + C)/(1 + t)$ which yields $C \ge B/t$ (as in the case of $M = 1$).

2. $c_1 + c_2 = C \le t(B + C)/(1 + t)$ which yields $C \le tB$ (as in the case of $M = 1$).

3. $b_1 + c_2 \le t(B + C)/(1 + t)$ which yields $c_2 \le -b_1 + t(B + C)/(1 + t)$.

4. $c_1 + b_2 \le t(B + C)/(1 + t)$ which yields $c_1 \le -b_2 + t(B + C)/(1 + t)$.

The first two subcases correspond to the pruning derived when $M = 1$. The last two subcases correspond to new additional pruning.

To see this, consider the example given above in the case of $M = 1$ with a query $\vec{B}$ with $B = 300$ and a similarity threshold of $t = 0.8$. By using the $M = 1$ level corresponding to the first two subcases, one needs to focus only on candidate molecules satisfying $240 < C < 375$. By using $M = 2$ and the last two subcases, one can prune additional molecules that fall within this interval. For instance, assume that $\vec{b} = (200, 100)$. Then even among the molecules $\vec{C}$ containing the same number of 1-bits as the query (i.e., $C = B = 300$), we can eliminate all of those associated with $\vec{c} = (c_1, c_2)$ for which $c_1 \le 166$ (subcase 4) or $c_2 \le 66$ (subcase 3).

In short, it should be clear that the greater the value of $M$, the greater the amount of pruning associated with the intersection bound and eq 25. While the amount of pruning is an essential consideration, other equally important factors must be taken into account in an actual database implementation. For instance, larger values of $M$ may require additional storage space and computing time. Furthermore, the data structure required to implement the pruning algorithm is also an important consideration. While a complete analysis of these issues is beyond the scope of this paper, it is worth noting that the cases of $M = 1$ and $M = 2$ suggest a natural and efficient data structure. First, organize the fingerprints $\vec{C}$ in the database into bins containing all the fingerprints that have the same number $C$ of 1-bits. Organize the bins by increasing values of $C$. Second, within any such bin,

organize the fingerprints by increasing values of, for instance, $c_1$ corresponding to the number of bits falling on odd components. For a given query $\vec{B}$ and threshold $t$, one can first use the intersection inequality with $M = 1$ to discard all the bins associated with values of $C$ that are too large or too small and focus on bins with $C$ in a limited range. Then, for these remaining bins, one can use the intersection inequality with $M = 2$ to discard additional molecules and focus on sub-bins of molecules with $c_1$ in a limited range. This procedure can be further repeated using, for instance, $M = 4$ raising again interesting questions of optimality and trade-offs.

Formulas for estimating the amount of pruning for different values of $M$ and different thresholds $t$ have been derived, but will be given elsewhere, together with an analysis of the corresponding trade-offs. Yet, even without those analyses, it is clear that the intersection bound is a simple but powerful tool for organizing and searching large databases of fingerprint vectors.

## REFERENCES AND NOTES

(1) Burkhard, W.; Keller, R. Some Approaches to Best-Match File Searching. *Commun. ACM* **1973**, *16*, 230–236.
(2) Shapiro, M. The Choice of Reference Points in Best-Match File Searching. *Commun. ACM* **1977**, *20*, 339–343.
(3) Shasha, D.; Wang, T.-L. New Techniques for Best-Match Retrieval. *ACM Trans. Inf. Syst.* **1990**, *8*, 140–158.
(4) Downs, G. M.; Willett, P. Similarity Searching and Clustering of Chemical-Structure Databases Using Molecular Property Data. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 1094–1102.
(5) Swamidass, S.; Baldi, P. Bounds and Algorithms for Exact Searches of Chemical Fingerprints in Linear and Sub-Linear Time. *J. Chem. Inf. Model.* **2007**, *47*, 302–317.
(6) Willett, P.; Barnard, J.; Downs, G. Chemical Similarity Searching. *J. Chem. Inf. Comp. Sci.* **1998**, *38*, 983–996.
(7) Leach, A. R.; Gillet, V. J. *An Introduction to Chemoinformatics*; Springer: Dordrecht, The Netherlands, 2005.
(8) Chen, J.; Swamidass, S. J.; Dou, Y.; Bruand, J.; Baldi, P. ChemDB: a Public Database of Small Molecules and Related Chemoinformatics Resources. *Bioinformatics* **2005**, *21*, 4133–4139.
(9) Baldi, P.; Hirschberg, D. S.; Nasr, R. J. Speeding Up Chemical Database Searches Using a Proximity Filter Based on the Logical Exclusive-OR. *J. Chem. Inf. Model.* **2008**, *48*, 1367–1378.
(10) Spath, H. *Cluster Analysis Algorithms*; Ellis Horwood: 1980.
(11) Lipkus, A. A Proof of the Triangle Inequality for the Tanimoto Distance. *J. Math. Chem.* **1999**, *26*, 263–265.
(12) Holliday, J. D.; Hu, C. Y.; Willett, P. Grouping of Coefficients for the Calculation of Inter-Molecular Similarity and Dissimilarity Using 2D Fragment Bit-Strings. *Comb. Chem. High Throughput Screening* **2002**, *5*, 155–166.
(13) Chen, J.; Linstead, E.; Swamidass, S. J.; Wang, D.; Baldi, P. ChemDB Update-Full Text Search and Virtual Chemical Space. *Bioinformatics* **2007**, *23*, 2348–2351.

CI900133J