

MultiMCS: A Fast Algorithm for the Maximum Common Substructure Problem on Multiple Molecules

Ramesh Hariharan,^{*,†} Anand Janakiraman,[†] Ramaswamy Nilakantan,[†] Bhupender Singh,[†] Sajith Varghese,[†] Gregory Landrum,[‡] and Ansgar Schuffenhauer[‡]

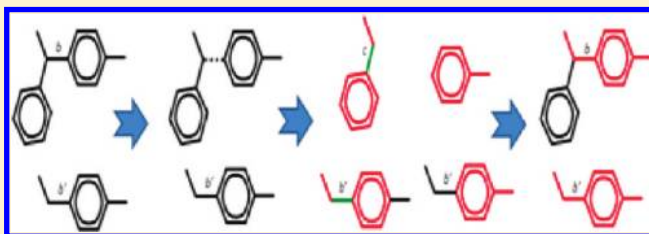
[†]Strand Life Sciences, Fifth Floor, Kirloskar Business Park, Bellary Road, Hebbal, Bangalore 560024, India

[‡]Novartis Institutes for BioMedical Research, CH-4002, Basel, Switzerland

 Supporting Information

ABSTRACT: Several efficient correspondence graph-based algorithms for determining the maximum common substructure (MCS) of a pair of molecules have been published in the literature. The extension of the problem to three or more molecules is however nontrivial; heuristics used to increase the efficiency in the two-molecule case are either inapplicable to the many-molecule case or do not provide significant speedups.

Our specific algorithmic contribution is two-fold. First, we show how the correspondence graph approach for the two-molecule case can be generalized to obtain an algorithm that is guaranteed to find the optimum connected MCS of multiple molecules, and that runs fast on most families of molecules using a new divide-and-conquer strategy that has hitherto not been reported in this context. Second, we provide a characterization of those compound families for which the algorithm might run slowly, along with a heuristic for speeding up computations on these families. We also extend the above algorithm to a heuristic algorithm to find the disconnected MCS of multiple molecules and to an algorithm for clustering molecules into groups, with each group sharing a substantial MCS. Our methods are flexible in that they provide exquisite control on various matching criteria used to define a common substructure.



INTRODUCTION

Medicinal and computational chemists are frequently in search of common structural patterns shared by sets of bioactive compounds. The maximum common substructure (MCS) is thus of great interest to them. Rather than being an abstraction consisting of fragments or descriptors, the MCS is very visual and easy to comprehend. However, while it is easy to define, the MCS is inherently hard to compute, and finding the MCS can be computationally very demanding when the data set consists of closely related complex structures.

Previous Work for the Two-Molecule Case. Efforts to develop MCS algorithms for chemical applications are not new. One of the first such attempts was that of Armitage et al.^{1,2} This was a brute-force approach to find the connected MCS between a pair of molecules, and it did not scale well to larger structures. Subsequently, Levi³ and Cone et al.⁴ published a deterministic algorithm to find the connected MCS shared by two molecules. Their method involves the use of a so-called compatibility table (or equivalently, a correspondence graph) that identifies compatible pairs of bonds/atoms in the two molecules. A path-tracing algorithm is then used to find paths connecting these atoms. Subsequent work by Raymond et al.^{5,6} placed the above approach of using a correspondence graph on a much more formal footing, first by noting that MCS computation reduces to computation of the largest clique in this graph, and second, by drawing upon a plethora of heuristics available in the literature on graph algorithms to compute the largest clique (for

instance, Pardalos and Xue⁷ and Wood⁸). This computation ends up using a Bron and Kerbosch⁹-like algorithm for finding all maximal cliques in graphs along with additional heuristics that allow avoiding parts of the search space if easily computable upper bounds on the clique size in these parts turn out to be smaller than the size of the largest clique found so far.

Previous Work for the Multiple-Molecule Case. The generic approach to this problem is to generate candidate substructures in one or more molecules and then check for the presence of these substructures in other molecules. Algorithms that implement this generic approach vary only in how candidate substructures are enumerated and how one checks for the presence of a given substructure in other molecules. Varkony et al.¹⁰ described the first such algorithm that enumerated all possible candidate substructures in increasing order of size in each of the molecules; these substructures were then converted to canonical keys using which pairs of substructures resulting from distinct molecules could be compared for identity. Keeping counts now enabled determination of the MCS. Efficiency was derived by generating candidate substructures of size m only by extending substructures of size $m - 1$ that were common to all the molecules. Accelrys Pipeline Pilot¹¹ has an MCS routine that is based on the same paradigm, except that canonical keys are generated via hashing to special types of fingerprints. Takahashi

Received: August 4, 2010

Published: March 29, 2011

et al.¹² gave an algorithm that used a substructure search instead of canonical keys or fingerprints to take candidate substructures in one molecule and determine the occurrence of these candidate substructures in other molecules. ChemAxon¹³ has developed another algorithm called Library MCS that uses the pairwise approach of the previous paragraph to determine the MCS for a given pair and uses substructure search to find occurrences of these candidate substructures in other molecules. Note that Library MCS aims not at finding the MCS of multiple molecules but at hierarchical clustering, where the MCS size is used as a measure of similarity between a given pair of molecules. Hierarchical clustering can often be performed using just pairwise MCS computations and does not need determination of the MCS of multiple molecules; indeed, as expected, there are cases where Library MCS does not find the correct MCS shared by a set of compounds (Figure 23). Also note that both Library MCS and Accelrys Pipeline Pilot seem to find only connected MCSs.

Comparing Algorithms for the Multiple-Molecule Case on Performance. It is challenging to compare running times of published algorithms because of differences in their associated computational and programming platforms. An alternative basis for comparison would be a measure such as the dependence of running time on the output MCS size. For instance, consider the algorithm due to Varkony et al.¹⁰ or Takahashi et al.;¹² it appears that these algorithms will need to consider every single substructure of the eventual MCS. Because the number of such substructures grows exponentially with the MCS size for most molecules, the running time of the algorithm will indeed be exponential in the output MCS size for most molecules. The algorithm implemented in Accelrys Pipeline Pilot¹¹ avoids this exponential behavior via optimizations that do not generate all possible substructures. Because these optimizations are not explicitly described to the best of our knowledge, it is hard to determine their impact on running time. In addition, a less than exhaustive generation of candidate substructures carries with it the risk of losing guarantees on correctness of the MCS found (we have found specific cases where Accelrys Pipeline Pilot fails to find the correct MCS shared by a set of compounds; see Figure 22 and the surrounding text). Indeed, we have found no previous work that is able to mute the exponential dependence of running time on the output MCS size (at least for MCS sizes of practical interest, say up to 50 or so), while simultaneously guaranteeing optimality for many if not all molecule families and characterizing exception molecule families if any.

MultiMCS: A New Algorithm for the Multiple-Molecule Case. In this paper, we present MultiMCS, a new algorithm for finding the MCS of multiple molecules. MultiMCS uses the approach of generating candidate substructures by enumerating all *maximal* common substructures for pairs of molecules (rather than enumerating *all* substructures in one or more molecules). In particular, we perform this computation for a chosen pivot molecule paired with each of the other molecules in turn; it turns out that the correct MCS for the connected case can then be obtained by taking intersections of these maximal common substructures. This approach poses the following key challenge: generating all maximal common substructures for a pair of molecules requires computing not just the maximum clique in the correspondence graph but all maximal cliques as well. This renders several heuristics for maximum clique computation invalid, leading to increased computational time.

Our solution to the above challenge is novel and not reported hitherto in literature. We use a new divide-and-conquer strategy

that provides substantial speedups in computing all maximal common substructures for a given pair of molecules. Indeed, we find empirically that the time taken by the algorithm applied without the divide-and-conquer strategy grows exponentially with the MCS size, while this exponential dependence gets substantially muted with the application of this divide-and-conquer strategy. We also show that the divide-and-conquer strategy guarantees the optimality of the MCS found for the connected case. However, even though the divide-and-conquer strategy provides substantial speedups for most molecules in practice, it does not provide speedups for all molecules. We provide a characterization of those molecule families for which the algorithm might run slowly and provide an extension of the divide-and-conquer approach to such molecules. While this extension does not guarantee the optimality of the MCS found, we observe that optimality indeed holds for several cases in practice.

We also extend the above algorithm to a heuristic algorithm for the disconnected case and to an algorithm for clustering molecules into families, with compounds in each family sharing a substantial common substructure.

The flexibility of our algorithm arises from the ability to specify the following:

- Whether the discovered MCS must be connected (single component) or may be disconnected (multicomponent).
- Whether or not ring bonds are allowed to match chain bonds in this MCS.
- Whether or not rings are allowed to match partially in this MCS.
- Whether or not the MCS found should include a specified seed substructure, i.e., a substructural fragment that is required by the user to be a part of the MCS.
- Different atom-typing schemes: a default scheme, an element-type-based scheme, or any other user-supplied scheme.

Finally, we illustrate our method by applying it to a series of test cases for the two main applications of interest, namely, (1) for tens or hundreds of molecules, finding the exact MCS of these molecules, and (2) for thousands of molecules, classifying these molecules into clusters, each of which has a substantial MCS (with molecules that cannot be clustered into an uninteresting miscellaneous set).

METHODS

We define the MCS problem for multiple molecules, recalling the correspondence graph framework, and describe why techniques used in previous work for speeding up the two-molecule case are not entirely applicable to the problem at hand. We then describe our techniques to solve this problem.

Definition. We are given molecules M_1, \dots, M_m , along with an atom-type label for each atom in each molecule and a bond-type label for each bond in each molecule. Typically, bond labels are single, double, triple, or aromatic. Consider the set \mathcal{S} of all tuples (b_1, \dots, b_n) , where bond $b_i \in M_i$. A common substructure (CS) between these molecules is defined as a subset $\mathcal{C} \subseteq \mathcal{S}$ with the following properties:

- For any two tuples $(b_1, \dots, b_n), (b'_1, \dots, b'_n) \in \mathcal{C}$ and all $i, 1 \leq i \leq n$, b_i and b'_i must be distinct bonds in M_i .
- If tuple $(b_1, \dots, b_n) \in \mathcal{C}$ then the bond types of bonds b_1, b_2, \dots, b_n are identical as are the atom-types of the end point atoms of these bonds.

- For any two tuples $(b_1, \dots, b_n), (b'_1, \dots, b'_n) \in \mathcal{G}$, one of the two conditions below must hold:
 - For every molecule M_i , $1 \leq i \leq n$, b_i and b'_i share an end point atom in M_i ; further, the atom type of this shared atom is the same in all molecules M_i , $1 \leq i \leq n$.
 - For every molecule M_i , $1 \leq i \leq n$, b_i and b'_i do not share an end point atom in M_i . In particular, note that the following two scenarios are ruled out by the above conditions. Consider the scenario where bonds b_i and b'_i share an end point atom in molecule M_i , but bonds b_j and b'_j do not share an end point atom in molecule M_j ; this is not allowed. Likewise, consider the scenario where bonds b_i and b'_i share an end point atom in molecule M_i with atom type, say x_1 , and bonds b_j and b'_j share an end point atom in molecule M_j with atom type, say $x_2 \neq x_1$; this again is not allowed.
- The number of distinct end point atoms over all bonds from molecule M_1 that appear in the various tuples in \mathcal{G} must be identical to the corresponding number of distinct end point atoms over all bonds from molecule M_i that appear in the various tuples in \mathcal{G} , for all i , $2 \leq i \leq n$.

A maximum common substructure (MCS) is defined as the CS with the largest number of tuples. Often additional constraints are placed on the nature of the MCS, as follows.

Connected/Disconnected MCS. When only connected CSs are of interest, then for any two bonds b_1, b_2 from molecule M_1 that appear in such a CS, there exists a path between b_1 and b_2 in molecule M_1 such that all bonds on this path also appear in this CS (the definition of a CS ensures that the same property will hold for all other molecules as well).

Ring Bonds Match Only Ring Bonds. When ring bonds must match only ring bonds, the following must hold for any CS of interest: for each tuple (b_1, \dots, b_n) in this CS, either all bonds in this tuple are chain bonds or all are ring bonds. This definition requires the identification of ring versus chain bonds in all molecules. A *ring bond* is defined as one that participates in a cycle, and all other bonds are called *chain bonds*.

Rings Must Match Fully. When rings must match fully, the following must hold for any CS of interest: if ring bond b from molecule M_1 is in the CS, then all bonds in some cycle comprising b in molecule M_1 must be in this CS as well. This option applies only when ring bonds must match ring bonds, in which case, by the definition of a CS, the above property will automatically hold not just for the first molecule but for every molecule.

Seeded CS. When a seed structure is specified, the following must hold for any CS of interest: in each molecule, the bonds participating in the CS must contain an occurrence of the seed structure as a substructure.

Of the cases above, the connected case with ring bond matching only ring bonds and rings matching fully is the most common use case in practice.

The Correspondence Graph Framework. The correspondence graph framework has been the topic of several previous publications (Raymond et al.,^{5,6} Stahl et al.¹⁴) for finding the MCS of two given molecules. For instance, given a pair of molecules M_1, M_2 , the RASCAL algorithm⁵ defines a correspondence graph CG as follows and casts MCS computation as a maximum clique determination algorithm on the product graph.

Correspondence Graph CG. The nodes of CG are bond pairs (b_1, b_2) , where $b_1 \in M_1$ and $b_2 \in M_2$, b_1, b_2 have the same bond

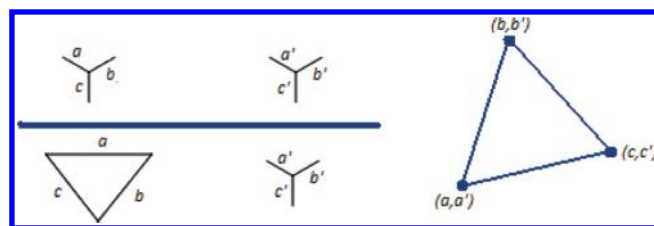


Figure 1. Two different graph pairs with the same correspondence graph. The first pair appears above the horizontal line, and the second appears below that line. In each pair, molecule 1 has bonds a, b, c , and molecule 2 has bonds a', b', c' . The correspondence graph in either case is the same and shown on the right. It has a clique of size 3, which corresponds to the MCS for the first pair of molecules but not for the second pair.

type, and the end point atoms of b_1, b_2 have the same atom types. In the ring-bonds-match-ring-bonds-only case, an additional criterion will be that either both b_1, b_2 are chain bonds or both are ring bonds. Pairs (b_1, b_2) and (b'_1, b'_2) are connected by an edge in CG if and only if either b_1, b'_1 do not share an atom in M_1 and likewise for b_2, b'_2 in M_2 (gray edges), or b_1, b'_1 do share an atom in M_1 and b_2, b'_2 also share an atom with the same atom-type label in M_2 (green edges).

Maximum Cliques. It is easy to see that the MCS of two molecules is now simply the maximum clique in CG with the following property: the number of distinct end point atoms over all bonds from M_1 that appear in the bond pairs in this clique must be identical to the corresponding number of distinct end point atoms over all bonds from M_2 that appear in bond pairs in this clique. This additional property is needed to guard against star–delta transforms where two different graph pairs (one isomorphic pair and one nonisomorphic pair, Figure 1) both yield the same correspondence graph (and as proved by Whitney,¹⁵ the star–delta transform is the only confounding case).

Optimizations for Clique Computation. Because maximum clique determination is a theoretically hard problem, i.e., NP-hard,¹⁶ the goal is to determine optimizations required to enable fast computation in practice and characterize input graph families to which these optimizations apply. There are two types of optimizations that appear in the literature.^{5,6}

Optimizations in Maximum Clique Determination. The standard maximum clique algorithm uses branch and bound to traverse the entire search space. At each branch point, a judgment has to be made whether or not to search further in the subspace implied by the choices made so far. For instance, suppose the algorithm is at a particular branch point where the vertex set V has been partitioned into U (included in the current clique), X (excluded), and R (the remaining vertices), and the algorithm seeks to extend U into a larger clique using vertices in R . At this point, it may be possible to compute upper bounds on the size of the clique that can be obtained by extending U . One such upper bound, though a very trivial one, will just be the number of vertices in R that are connected to all vertices in U . Finer bounds can be obtained by using the degrees of vertices in R or bounds on their chromatic number obtained via a greedy coloring, etc. If this upper bound turns out to be lower than the size of the largest clique already computed, the algorithm can backtrack without exploring the current subspace. The RASCAL algorithm⁵ suggested several such upper bounding mechanisms.

Optimizations in Correspondence Graph Definition. Maximum clique computation can be expensive in spite of the above

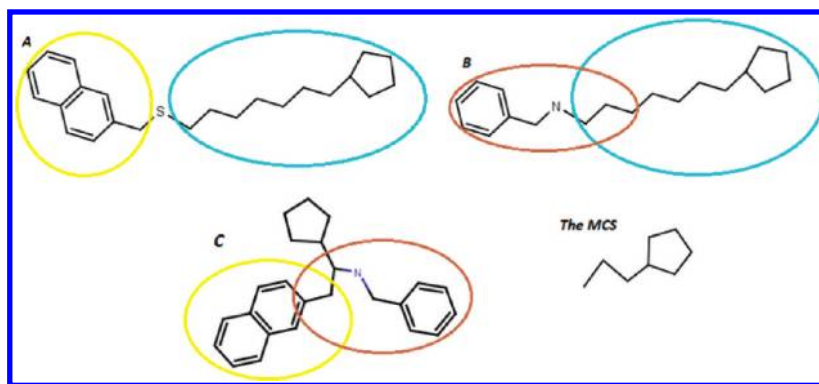


Figure 2. Performing pairwise MCS is not sufficient. Note that the MCS shown above is relative to an element-based atom-typing scheme, i.e., the atom type of an atom is its element type, and chain bonds are not allowed to match ring bonds.

optimizations, which begs the following question: can unnecessary nodes and edges in the correspondence graph be pruned off before invoking the maximum clique procedure? Some such optimizations, namely, symmetry and ring substitution alignment heuristics and weak and strong ring heuristics, appear in the RASCAL paper⁵ and subsequent papers.⁶

As described below, optimizations in maximum clique determination do not apply to our setting, and optimizations in correspondence graph definition do apply but provide limited speedup.

Challenges in Generalizing to Multiple Molecules. Can the above framework for two molecules be generalized to find the MCS of multiple molecules? For instance, can the overall MCS be determined from the individual pairwise MCSs? As shown in the example in Figure 2, the answer is no. In this example, the cyan ellipses show the connected MCS between molecules A and B, the yellow ellipses show the connected MCS between molecules A and C, and the brown ellipses show the connected MCS between molecules B and C. The actual MCS of the three molecules is also shown in the figure and is clearly distinct from all the above pairwise MCSs. Note that these results are relative to an *element-based atom-typing* scheme, i.e., the atom type of an atom is its element type, and chain bonds are not allowed to match ring bonds. So, an approach that finds these pairwise MCSs and uses a substructure search or other means to determine the presence of these substructures in other molecules (Library MCS¹³ would be an example) would not find the true MCS of these three molecules. It is also easy to see that even if we were to identify all pairwise disconnected MCSs, we would still not have enough information to obtain the actual MCS from these pairwise disconnected MCSs.

The above example does suggest an approach to generalization: the overall MCS can be obtained by computing all maximal common substructures and not just the MCS for each pairwise comparison (a maximal common substructure is one that cannot be extended further with more bonds; such substructures can be computed by determining all maximal cliques as opposed to just the maximum clique). However, this means that pairwise comparisons can no longer use optimizations in maximum clique computation described above, thus leading to decreased performance. Optimizations in correspondence graph definition as described above can indeed be applied, but achieve limited speedup. Hence, new ideas are needed to obtain substantial speedups in this setting.

Algorithm for Connected MCS Computation. The overall algorithm for the connected case is outlined in Algorithm 1. In

the clique finding phase, the algorithm performs pairwise comparisons between one specially chosen *pivot* molecule, say molecule M_1 , and every other molecule in turn. The comparison between molecule M_i , $2 \leq i \leq n$, and M_1 determines *all* maximal cliques corresponding to connected maximal common substructures between M_i and M_1 ; for each of these connected maximal common substructures found here, the corresponding substructure of M_1 is added to the set S_i . It now suffices to find the largest substructure with a presence in each of the S_i sets, $2 \leq i \leq n$. To this end, the clique intersections phase starts with S_2 , intersects each structure in S_3 with each structure in S_2 to determine all maximal intersections, then intersects each structure in S_4 with each of these maximal intersections, and so on. At the end, one is left with a collection of maximal intersections, each of which is a candidate common substructure over all of M_1, \dots, M_n , and the largest of these is the MCS.

Algorithm 1: The overall algorithm for the connected case for multiple molecules.

Input: Molecules M_1, \dots, M_n

Output: Connected MCS of these molecules

Pick a pivot molecule, say M_1 ;

//The Clique Finding Phase;

foreach Molecule M_i , $i = 2 \dots n$ **do**

Set $S_i \leftarrow \Phi$;

Compute the correspondence graph CG_i of molecules M_1 and M_i ;

Compute all maximal cliques in CG_i which correspond to connected substructures;

foreach maximal clique C_i , $i = 2 \dots k$ **do**

Project C_i on to M_1 to obtain substructure X ;

Add X to S_i ;

//The Clique Intersections Phase;

Set $S \leftarrow S_2$;

foreach $i = 3 \dots n$ **do**

$Z \leftarrow \{Y \cap X \mid \text{for all } Y \in S_i, X \in S\}$;

$S \leftarrow$ The set of maximal substructures in Z ;

Report the largest substructure in S as the MCS;

Optimality. It is easy to see that Algorithm 1 guarantees that the output is indeed an MCS of M_1, \dots, M_n . For instance, for the example in Figure 2, two pairwise computations are run in the clique finding phase, say A with B, and A with C. The first obtains

the substructure marked by the cyan ellipse as one of the maximal substructures, and the second obtains the actual MCS as one of the maximal substructures. In the clique intersections phase, the MCS survives, and being the largest such substructure, it is declared the MCS.

Pivot Molecule Selection. Selecting the smallest molecule as the pivot favors speed, but that is a question we will revisit in the context of further optimizations.

Finding All Maximal (Connected) Cliques. Given a graph G , we first outline a standard search algorithm (e.g., the Bron-Kerbosch algorithm⁹) to obtain all maximal cliques in G . The algorithm maintains three sets at each instant: set C comprising vertices in the current clique, set N comprising vertices that have edges to all vertices in C , and set X comprising vertices to be excluded when seeking to construct cliques that are extensions of C . It initializes C to Φ , N to the set of all vertices, and X to Φ . In general, given C, N, X , the algorithm first checks if N is empty, recording C if this is indeed the case; C is a maximal clique because there is no vertex that has edges to all vertices in C . Otherwise, if N is nonempty, it extends C by simply picking each vertex w in $N - X$ in turn. For each chosen w , C is then augmented with w , N is pruned down to only vertices that have edges to w (note that vertices in N are already guaranteed to have edges to other vertices in C), and the algorithm then recurses to compute all maximal cliques which are extensions of $C \cup \{w\}$. After this recursive call returns, the algorithm moves on to the next w in $N - X$. However, before doing this, it adds w to X so cliques containing $C \cup \{w\}$ are not generated again. It is easy to see that each maximal clique in G is output by Algorithm 2.

Algorithm 2: Find all maximal cliques in a graph G .

Input: Graph G

Output: All maximal cliques in G

begin

$C \leftarrow \Phi$;

$X \leftarrow \Phi$;

$N \leftarrow$ Set of all vertices in G ;

FindMaximalCliques(C, X, N);

Function FindMaximalCliques(C, X, N);

begin

if $N = \Phi$ **then**

Record C as a maximal clique;

else

foreach vertex w in N **do**

if $w \notin X$ **then**

FindMaximalCliques($C \cup \{w\}, X, N \cap \{\text{neighbors of } w\}$);

$X \leftarrow X \cup \{w\}$;

by vertices in this clique is connected on the basis of green edges alone (i.e., ignoring gray edges). We use the term green-connected for such cliques in the description below. Because we seek only green-connected cliques and because an overwhelming proportion of edges in CG are gray edges, we derive efficiency by treating gray edges implicitly and working only with green edges and red nonedges explicitly (i.e., if a pair of vertices in CG has no associated edge, we introduce a red edge between these vertices). This reduces wasteful time on repeatedly traversing gray edges. The appropriately modified Bron-Kerbosch variant is described in Algorithm 3. Note the key addition: we have two new sets NG and NR , comprising vertices that are connected to *at least one* of the vertices in C via a green edge and a red edge, respectively. So, effectively $NG - NR$ comprises vertices that are connected to each of the vertices in C via a green or a gray edge (with at least one of these connections being green).

A standard way to make the Bron-Kerbosch algorithm faster is to introduce the notion of pivoting. Consider the *for* loop that goes over all vertices $w \in N$; this loop is skipped for those vertices w that are neighbors of a chosen pivot vertex z (a different pivot is chosen in each recursive call). Skipping these vertices w in the *for* loop does not affect correctness because cliques that contain C and w together are found in the recursive call that adds z to C . Note, however, that this argument holds only for the standard Bron-Kerbosch algorithm⁹ and not the variant we use above, which drives efficiency by treating gray edges implicitly. Pivoting would indeed apply if we were to treat gray edges explicitly; however, we do not know how to apply pivoting if gray edges are treated implicitly.

Algorithm 3: Find all maximal green-connected cliques in CG .

Input: Correspondence graph CG for Molecules M_i, M_j

Output: All maximal green-connected cliques for these molecules

begin

$C \leftarrow \Phi$;

$X \leftarrow \Phi$;

$NG \leftarrow \Phi$;

$NR \leftarrow \Phi$;

$N \leftarrow$ Set of all vertices in CG ;

FindMaximalGreenConnectedCliques(C, X, NG, NR, N);

Function FindMaximalGreenConnectedCliques(C, X, NG, NR, N);

begin

if $N = \Phi$ **and** C satisfies the star-delta check **then**

Record C as a maximal green-connected clique;

else

foreach vertex w in N **do**

if $w \notin X$ **then**

FindMaximalGreenConnectedCliques($C \cup \{w\}, X, NG \cup$

$\{\text{green neighbors of } w\}, NR \cup \{\text{red neighbors of } w\}, NG - NR$;

$X \leftarrow X \cup \{w\}$;

We implement the Algorithm 2 on our correspondence graph framework, but with the following changes. Recall that edges in the correspondence graph are designated as green or gray, and note that we seek only maximal cliques corresponding to connected substructures. A clique corresponding to a connected substructure is a clique that satisfies the following additional property: the subgraph of the correspondence graph CG induced

Performance. Figure 3 shows timings for several three-molecule instances (each dot represents an MCS computation on a three-molecule instance, where the compounds are taken

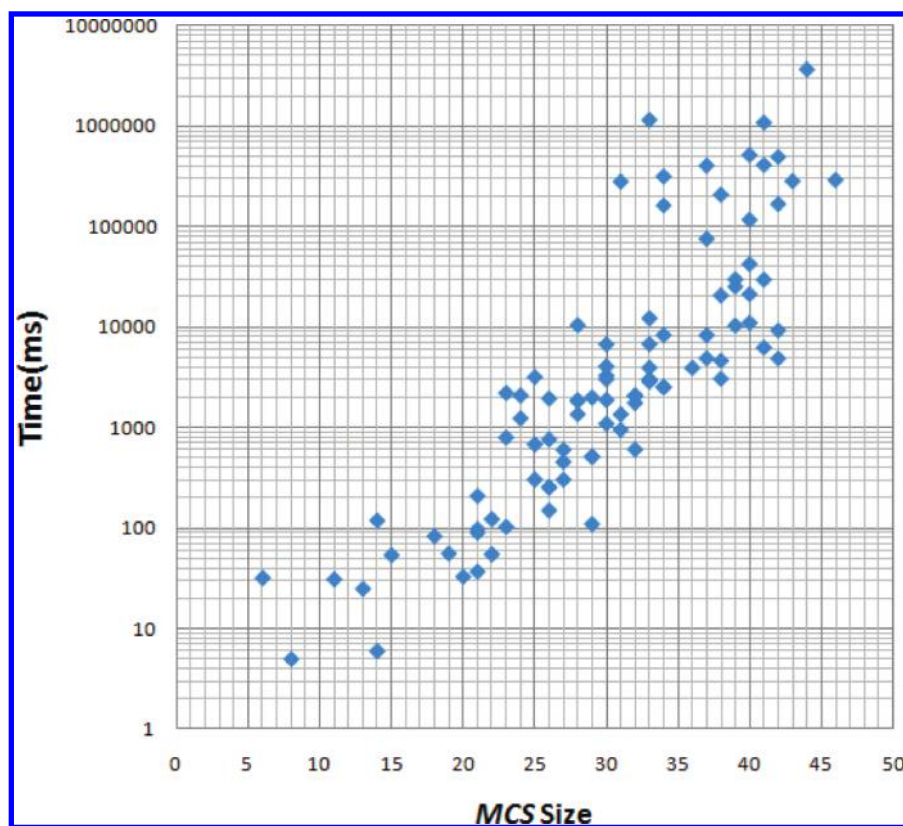


Figure 3. Time vs output size. Note that time is shown on the logarithmic scale. The R^2 value is 0.716.

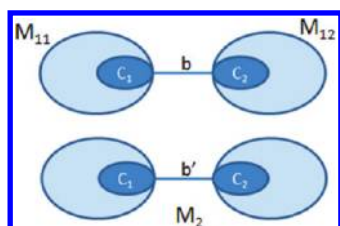


Figure 4. Setting for chain bond decomposition.

from a data set described in Lounkine et al.¹⁷). As shown in the figure, time grows exponentially with the size of the MCS being found. In particular, if the MCS has size beyond 30, the time starts getting unacceptably large.

The challenge now is the following: how do we speed up the computation without losing provable correctness for the vast majority of molecules?

Key Improvement: Chain Bond Decomposition. A natural approach is one of divide and conquer, i.e., breaking up a molecule into smaller molecules, solving the problem recursively on smaller molecules, and then assembling the final results from the results for the individual pieces. The challenge is to show that this approach does indeed find all maximal green-connected cliques. We show how this can be achieved for many molecule classes and precisely characterize those classes for which this approach might not work in terms of yielding speedup.

To describe the divide-and-conquer framework, consider the case when ring bonds are not allowed to match chain bonds and suppose we are trying to find all maximal green-connected cliques for molecules M_1, M_2 . Next, suppose we remove a chain bond b from molecule M_1 to obtain two smaller molecule fragments, M_{11}

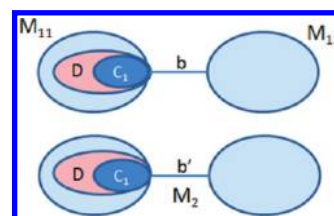


Figure 5. D does not cross bond b' in M_2 .

and M_{12} , and we run Algorithm 3 for each of the pairs M_{11}, M_2 and M_{12}, M_2 . From these two results, we can piece together the desired result for the pair M_1, M_2 , as described next.

Consider any green-connected maximal clique C that is output by Algorithm 3 run on molecules M_1, M_2 and in which bond b in M_1 matches say bond b' in M_2 . Let C_1 denote the subclique comprising only bonds in M_{11} , and let C_2 denote the subclique comprising only bonds in M_{12} (Figure 4).

If C_1 and C_2 are identified as maximal green-connected cliques when Algorithm 3 is run on M_{11}, M_2 and M_{12}, M_2 , respectively, then C can be obtained by simply combining C_1, C_2 and the pair (b, b') . Otherwise, suppose, without loss of generality, C_1 is not identified as a maximal green-connected clique when Algorithm 3 is run on M_{11}, M_2 . This can only happen because C_1 is not maximal, i.e., a larger green-connected clique D containing C_1 must have been identified by this run of the algorithm. We show that C_1 can be determined from D using a simple procedure. There are two cases here as described below.

First, suppose D does not contain any matching for bond b' (Figure 5; for simplicity, C_2 is not shown in this figure). Then, because both b and b' are chain bonds, it is easily seen that $D \cup C_2$

$\cup (b, b')$ is a bigger green-connected clique than $C = C_1 \cup C_2 \cup (b, b')$, and it contains C . C is therefore not maximal, a contradiction.

That leaves the case that D does contain a matching for bond b' (Figure 6). Let c be the bond in M_{11} matching b' . Because b' is a chain bond, it must be the case that removing bond pair (c, b') breaks D into two green-connected subcliques, one of which, say D' , completely contains C_1 . So, $D = E \cup D' \cup (c, b')$. If D' is strictly bigger than C_1 , then we can show that C is not maximal by considering $D' \cup C_2 \cup (b, b')$ as above. So, it must be the case that $D' = C_1$. So effectively, one can compute C_1 from D quite easily, and likewise for C_2 ; and from C_1 and C_2 , one can compute C .

A concrete example of the above logic is given in Figure 7. Part A shows the two molecules in question. In part B, a chain bond b is removed from M_1 to yield M_{11} and M_{12} , and we seek to identify the MCS where bond b in M_1 matches bond b' in M_2 (the algorithm will try all applicable bonds b' in M_2 in turn). Algorithm 3 is now run on the pairs M_{11}, M_2 and M_{12}, M_2 to identify all maximal green-connected cliques in each case. The largest such clique found for M_{11}, M_2 appears in part C and for M_{12}, M_2 in part D. These cliques are indicated in red, except for the bond pair (c, b') which is indicated in green (compare part C with Figure 6). Note that the largest such clique shown for M_{12}, M_2 in part D does not involve b' and hence has no green bonds. To combine results from the two parts, one simply drops the green bond c to obtain D' and E , and then combines D' from part C with the red structure C_2 in part D along with the bond b to obtain the final result, shown in part E.

We capture the above procedure in the Algorithm 4.

Algorithm 4: Find all maximal green-connected cliques for M_1, M_2 using Algorithm 3 on

M_{11}, M_2 and M_{12}, M_2 .

Input: All maximal green-connected cliques for M_{11}, M_2 and for M_{12}, M_2 , where M_1 has been broken up into M_{11}, M_{12} and bond b

Output: All maximal green-connected cliques for M_1, M_2

Record all maximal green-connected cliques computed for M_{11}, M_2 and for M_{12}, M_2 as maximal green-connected cliques for M_1, M_2 if they do not contain bonds that share an endpoint atom with b ;

foreach chain bond $b' \in M_2$ **do**

$X_1 \leftarrow$ Set comprising those maximal green-connected cliques for M_{11}, M_2 that contain a bond in M_{11} which shares an endpoint atom with b ;

$X_2 \leftarrow$ Set comprising those maximal green-connected cliques for M_{12}, M_2 that contain a bond in M_{12} which shares an endpoint atom with b ;

$X'_1 \leftarrow \Phi$;

$X'_2 \leftarrow \Phi$;

foreach $i = 1, 2$ **do**

foreach clique D in X_i **do**

if D contains a bond pair (c, b') **then**

 Split D into two green connected cliques D', E by removing this bond pair;

 Add to X'_i whichever of D', E contains a bond in M_{1i} that shares an endpoint atom with b ;

else

 Add D to X'_i

foreach maximal clique pair $C_1 \in X'_1, C_2 \in X'_2$ **do**

if $C_1 \cup C_2 \cup (b, b')$ is indeed a valid clique **then**

 Record $C_1 \cup C_2 \cup (b, b')$ as a maximal green-connected clique for M_1, M_2 ;

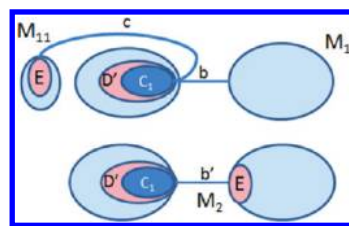


Figure 6. D straddles bond b' in M_2 .

removed chain bonds are added back one at a time, and the procedure in Algorithm 4 is then run each time. The goal of bond removal is to ensure that M_1 is decomposed into fragments that are small (via a customizable cutoff, we use a default of 20 atoms).

Pivot Molecule Selection. The pivot molecule (Algorithm 1) was earlier chosen as the smallest of the given molecules. Now we switch strategies so the pivot is chosen as any molecule which is indeed decomposable into small fragments as above, and if there are many such molecules, the smallest of those is chosen.

Rings Matching Chains Case. The above strategy only works when ring bonds are not allowed to match chain ring bonds. For the case when ring bonds are indeed allowed to match chain bonds, we do the following. The correspondence graph is constructed for M_1, M_2 so it has all relevant ring–chain bond pairs but no chain–chain bond pairs involving removed bonds; Algorithm 3 is used to find all green-connected maximal cliques. Then each of the removed bonds is put back and Algorithm 4 is run. It can be shown that this strategy indeed finds all relevant green-connected maximal cliques for M_1, M_2 .

Performance. The divide-and-conquer policy leads to very significant speedups for a variety of molecules *without* compromising on correctness. Figure 8 shows timings for the same set of instances represented in Figure 3. As shown, the exponential dependence of time on MCS size is significantly muted, and all timings are significantly better. This holds for both the ring-bonds-matching-ring-bonds-only and the ring-bonds-matching-chain-bonds cases.

Characterizing Slow Molecules and Heuristics for Speedup. Which molecules do not admit to speedup via the above divide-and-conquer strategy? These are molecules with substantial ring components that cannot be decomposed into smaller fragments via chain bond removal (an example would be staurosporine). Our strategy to speed up computations for these classes of molecules is analogous to the divide-and-conquer strategy used above, i.e., drop bonds, compute maximal cliques for the resulting molecules, and put back the removed bonds and update the cliques above using a modified version of Algorithm 4. However, this strategy is not provably guaranteed to find all green-connected maximal cliques; it is a heuristic that seems to work for most cases, as illustrated below.

Consider any pair of bonds $b \in M_1, b' \in M_2$ that are matched in a particular green-connected maximal clique C , and suppose we drop the pair (b, b') from the correspondence graph and find all green-connected maximal cliques via Algorithm 3. The question then is: can the resulting maximal cliques be combined in some way to result in the clique C ? It may happen that the only way to do this is by merging two resulting maximal cliques, along with the vertex (b, b') , that are mutually incompatible, i.e., their union is not a clique. The merging process will then need to drop some vertices. There may be many possible combinations of vertices that could be dropped to obtain a clique, and each such combination

Bond Removal Strategy. The above strategy can be generalized to the removal of multiple chain bonds from M_1 . The

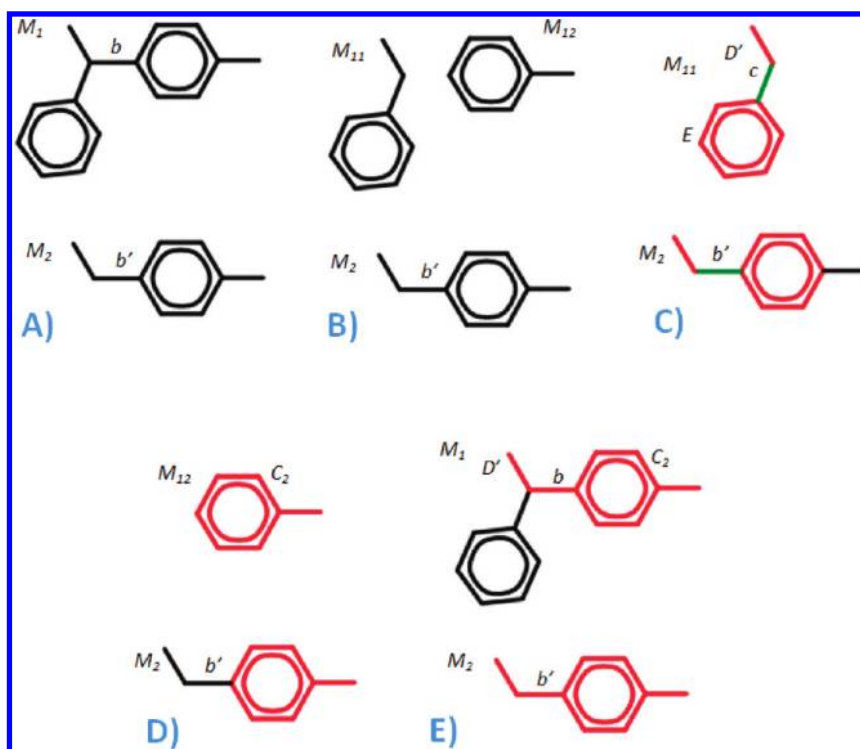


Figure 7. Chain bond decomposition on a real example.

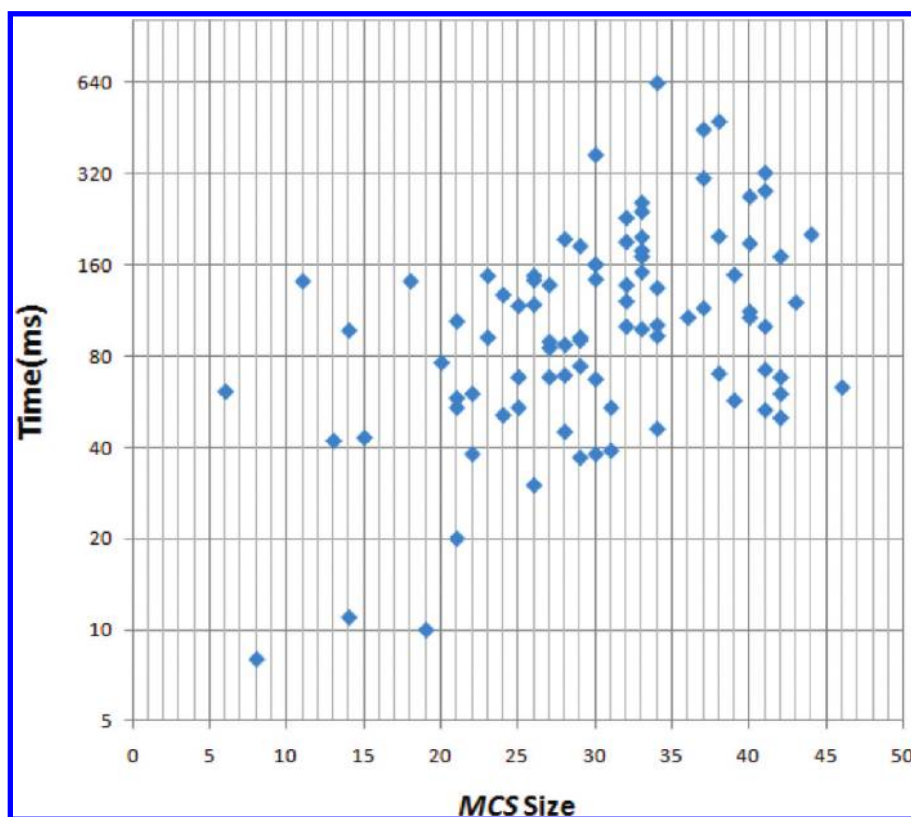


Figure 8. Time vs output size with the divide-and-conquer strategy. Note that time is shown on the logarithmic scale. The R^2 value is down to 0.203 (compare with Figure 3).

yields a maximal clique. However, identifying all such combinations is time consuming.

Instead, we obtain only those maximal cliques that can be obtained by merging compatible green-connected maximal

cliques resulting from running Algorithm 3 after dropping (b, b') . Because we do not obtain all possible maximal cliques, we cannot guarantee optimality of the MCS found at the end. However, we find that this algorithm yields incorrect results only in rare cases. An example where this algorithm yields correct results is shown in Figure 9. Note that in this instance, the steroid nucleus is not chain-bond decomposable. Also note that this example as well as the example in Figure 10 are run using the element-based atom-typing scheme, i.e., the atom type of an atom is its element type.

We illustrate a scenario that yields an incorrect result in Figure 10. The top row shows the two molecules in question with the size 20 MCS marked in red. Now consider the bottom row, and imagine that the pink bond labeled 10 is removed from the first molecule. Then all nodes in the correspondence graph comprising this pink bond are removed; in particular, the node $(10, 10)$ is removed. After this removal, we find all green-connected maximal cliques in the correspondence graph. Two such cliques are marked in the bottom row: clique A in green and clique B in blue/yellow. Let clique B' denote clique B restricted to just the blue bonds. Note that the MCS shown in the first row can be obtained by combining cliques A and B' along with the node $(10, 10)$. However, clique B is not green-connected maximal after removal of node $(10, 10)$ from the correspondence graph, clique B is. But cliques A and B and node $(10, 10)$ cannot be combined into a single clique because bonds 22, 8 share an atom in molecule

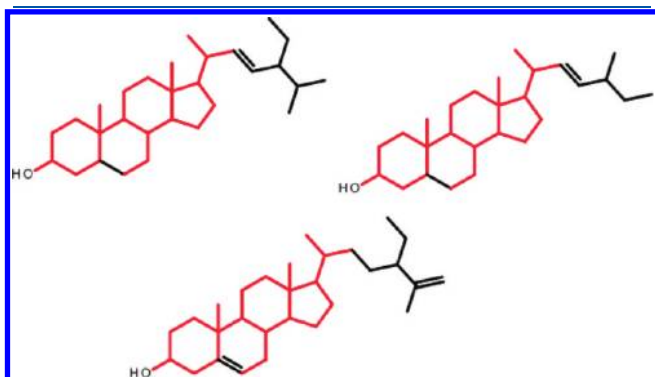


Figure 9. Example where the heuristic yields correct results for the connected case. The running time of the algorithm reduces from 38 s to 1 s when the heuristic is applied.

1 but not in molecule 2, and likewise bonds 24, 7 share an atom in molecule 2 but not in molecule 1. So, the algorithm ends up reporting a size 19 common substructure obtained by combining clique A with a completely different green-connected maximal clique B'' shown in blue in Figure 11.

Clearly, our approach could be extended in future work, so the violating bond pairs 22, 8 and 24, 7 are dropped while merging cliques A and B , in which case the optimum MCS would indeed be found.

Disconnected Case. We solve the disconnected case for a pair of molecules by combining together green-connected maximal cliques to generate candidate disconnected cliques. For the case of multiple molecules, these then go through the clique intersections phase in Algorithm 1. Candidate combinations of green-connected cliques are generated as in Algorithm 5, via a heuristic procedure.

Algorithm 5: Find candidate disconnected cliques.

Input: All green-connected maximal cliques for M_1, M_2

Output: Candidate disconnected cliques for M_1, M_2

Order green-connected maximal cliques in decreasing order of size $C_1 \dots C_n$;

foreach $i = 1$ to n **do**

FindGoodCombination($C_i, \{C_1 \dots C_n\}$);

Output this good combination;

Function FindGoodCombination(C, X);

begin

$C' \leftarrow C$;

repeat

Find that clique C'' in X which when combined with C' results in the largest increase in size;

Update C' to reflect this combination;

Remove C'' from X ;

until $X = \Phi$;

Return C' ;

Clustering Molecules. The goal here is to group a set of molecules into clusters, where each cluster has a prominent MCS. Some molecules that do not fall in any cluster are combined into a miscellaneous cluster. The clustering procedure proceeds as follows.

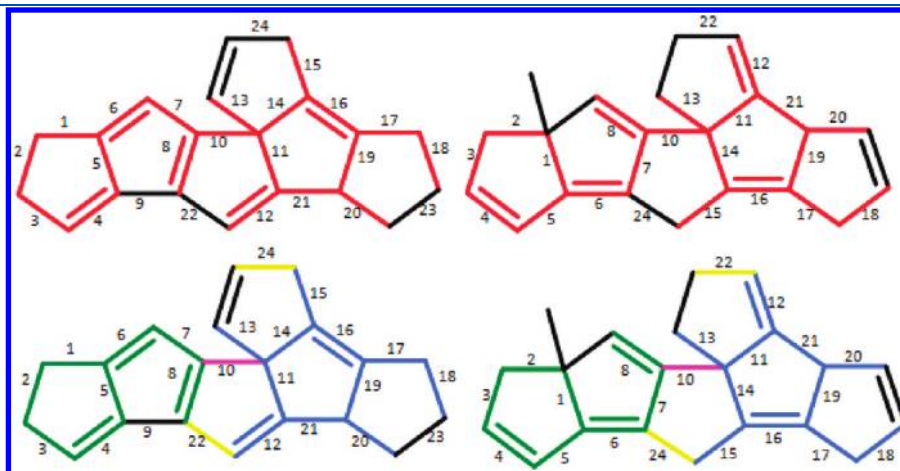


Figure 10. Top row shows the size 20 MCS between the two molecules in question. The bottom row shows green-connected maximal cliques found when the pair of pink bonds $(10, 10)$ is removed from the correspondence graph (see previous explanation).

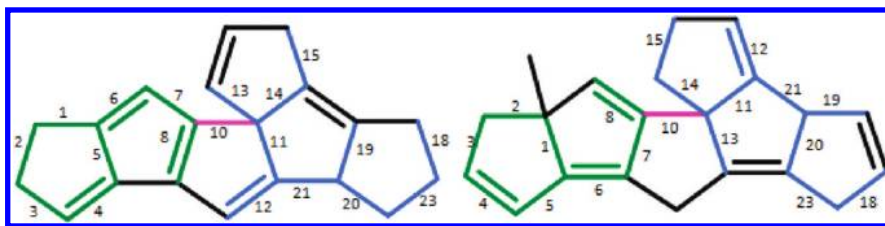


Figure 11. Size 19 common substructure found by the heuristic for the molecules in Figure 10. Note that the actual MCS has size 20.

First, we compute the Tanimoto similarity for each pair of compounds on the basis of fingerprints provided by the user. Next, we pick the most similar pair of compounds, find the MCS of this pair, and then use a substructure search method to find those molecules in the remaining set that contain this MCS, thus defining a cluster. We then keep the molecules in that cluster aside and recurse on the remaining molecules.

There is one problem with the previous approach that needs to be addressed: the MCS of the two molecules picked may be too specific to those two molecules. As a result, the cluster returned by the substructure search could be very small, causing heavily fragmented clustering. To circumvent this problem, we pick multiple molecules instead, find their MCS, and issue a substructure search with the result. Using multiple molecules ensures that the MCS found is common enough and not too specific. The number of molecules used in this process is specifiable by the user and drives the trade-off between specificity and sensitivity. The algorithm is outlined in Algorithm 6.

Algorithm 6: Cluster molecules into families sharing a substantial common substructure
Input: Set of Molecules \mathcal{M} along with associated Fingerprints for Similarity Computation
Output: Clusters of Molecules with the MCS of each cluster, and a miscellaneous set of unclustered molecules, or singletons
 Determine the Tanimoto similarity of all pairs of molecules;
 Mark all molecule pairs as valid;
repeat
 Pick a valid pair P that has the highest Tanimoto similarity;
 Mark the two molecules in P as invalid;
 If the similarity of P is below the *minTanimotoSimilarity* threshold, break;
 Pick all valid molecules M whose similarity to both molecules in P exceeds the *minTanimotoSimilarity* threshold;
 If $|M| < \text{minClusterSize} - 2$, continue;
 Find the MCS of the seed molecules $M \cup P$;
 if this MCS is smaller than minClusterMCSSize then
 Mark molecules in M as invalid;
 else
 Do substructure search using this MCS on the remaining molecules to identify say H hits;
 Combine M, P, H into a cluster and output this cluster;
 Mark all molecules in $M \cup P \cup H$ as invalid;
until there are no more valid molecules left;
 Output each invalid molecule as a singleton;

The *minClusterSize* threshold drives the number of seed molecules in a cluster. Large threshold values result in a small number of large clusters. Smaller thresholds produce more,

smaller, clusters. The *minClusterMCSSize* threshold is used to guarantee that the MCSs found are large enough to be meaningful. In general, large values of either threshold tend to produce more chemically meaningful results at the cost of singletons. These two parameters can be tweaked on any given data set to balance the numbers of clusters and singletons. An example of this tweaking can be seen in Figure 17.

Note that the computation of *minTanimotoSimilarity* can be achieved using any set of molecular fingerprints that satisfy the following property: two molecules with a substantial MCS should yield good similarity using these fingerprints. Our own experiments below are conducted with fingerprints that are extensions of the Avalon fingerprints.¹⁸ We have found that changes in the *minTanimotoSimilarity* parameter do not seem to change results very much, unless the parameter is increased from its default value of 0.57 to 0.75 or higher. This may be too high a value, i.e., two distinct molecules with significant MCSs may still not have such a high similarity value. Therefore, we recommend leaving this parameter as is.

Finally, note that our aim here is not to delve into the intricacies of clustering molecules; rather, our goal is to show how a routine for finding the MCS of multiple molecules can be used in a clustering procedure to drive the trade-off between cluster sizes and singletons. Also note that there are other approaches known for discovering fragments that discriminate between active and inactive compounds in a large structure databases.^{19,20}

RESULTS AND DISCUSSION

MultiMCS has been implemented as a Java program that executes using JRE 1.6. Various results obtained by MultiMCS are analyzed below. All timings have been obtained on a Windows XP system with 2GB RAM and an Intel Core 2 Duo P8400 processor running at 2.26 GHz. Our implementation uses only a single thread, though extension to multiple threads can easily be done. Unless specified otherwise, all timings below use the *default* atom typing scheme described in Figure 24.

Finding the MCS in Sets with Tens or Hundreds of Molecules. *A. Two-Molecule Case.* We illustrate the power of MultiMCS with an example where the structures are fairly complex and share a large common substructure. We used two staurosporine analogs. The calculations were done with the following options: ring bonds can only match ring bonds, partial ring matches are allowed, and disconnected MCS is allowed. Figure 12 shows the structures and the resulting 40 bond MCS. This calculation took 0.5 s.

B. Multiple Molecules: Simple Test Case. Next we give an example of what might be the standard use case for MultiMCS. We ran the MCS calculation on a set of 51 factor Xa inhibitors from Fontaine et al.²¹ The program was run in two ways: (a) not allowing disconnected MCS and (b) allowing disconnected

MCS. In both cases, ring bonds were allowed to match only ring bonds and partial ring matches were disallowed. The results are shown in Figure 13. The upper row of structures shows the connected case. Three different structural types are shown with the MCS colored in red. The bottom row shows the disconnected case for the same three compounds, with the MCS colored in red. It is clear that these 51 compounds share a significant common substructure: 17 bonds in the connected case and 26 bonds in the disconnected case. It can also be seen that this set of compounds all contain a biphenyl sulfonamide fragment and a benzamidine fragment connected to each other by a varying bridge fragment. This fact is buried in the data and not clearly apparent until the MCS is discovered. Once this pattern is discovered, it allows the scientist to focus on the variations in biological activity as a function of the bridging fragment. The run times are very fast: 5.0 s for the connected case and 6.6 s for the disconnected case.

C. Multiple Molecules: More Challenging Case. We tested MultiMCS on a set of 111 rapamycin analogs. These compounds are large macrocycles having between 63 and 79 non-hydrogen atoms. The calculations were done with the following options: ring bonds are allowed to match only ring bonds, partial ring

matches are allowed, and disconnected MCS is allowed. Rather than draw the 111 analogs and their MCS mappings, in Figure 14, we show one of the structures from the set with the MCS highlighted (the MCS contains 40 bonds and consists of two disconnected fragments). These fragments include parts of the macrocycle as well as the smaller rings. This calculation took about 40 s. The large size and the structural complexity of these structures are noteworthy. Without the heuristic described above, the computation time is in excess of 2000 s (the program timed out after 2000 s).

Multiple Molecules: Dependence of Running Time on the Number of Input Molecules. The earlier discussion focused on the dependence of run time on the output MCS size. Another important factor is the scaling of run time with the number of input molecules. Here, we show that run time scales linearly with the number of input molecules. To measure this, we created test sets containing multiple copies of a single molecule (a rapamycin

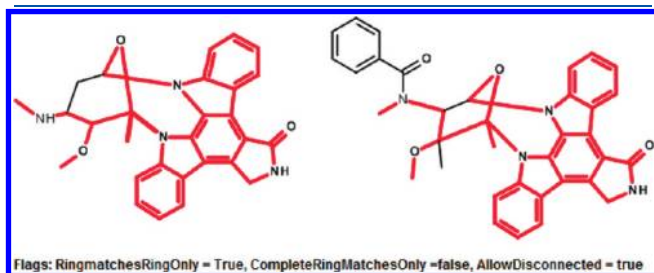


Figure 12. MCS calculation on two staurosporine analogs. Settings: ring bonds can only match ring bonds, partial ring matches are allowed, and disconnected MCS is allowed. The resulting MCS (colored red) has 40 bonds.

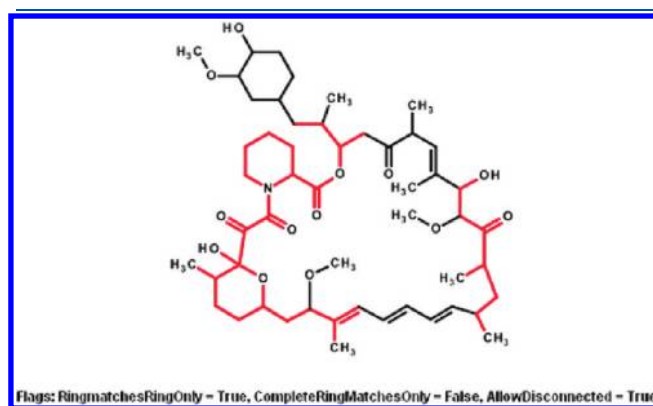


Figure 14. MCS calculation on a set of 111 rapamycin analogs. Settings: ring bonds can only match ring bonds, partial ring-matches are allowed, and disconnected MCS is allowed. The MCS (colored red) has 40 bonds and consists of two disconnected fragments.

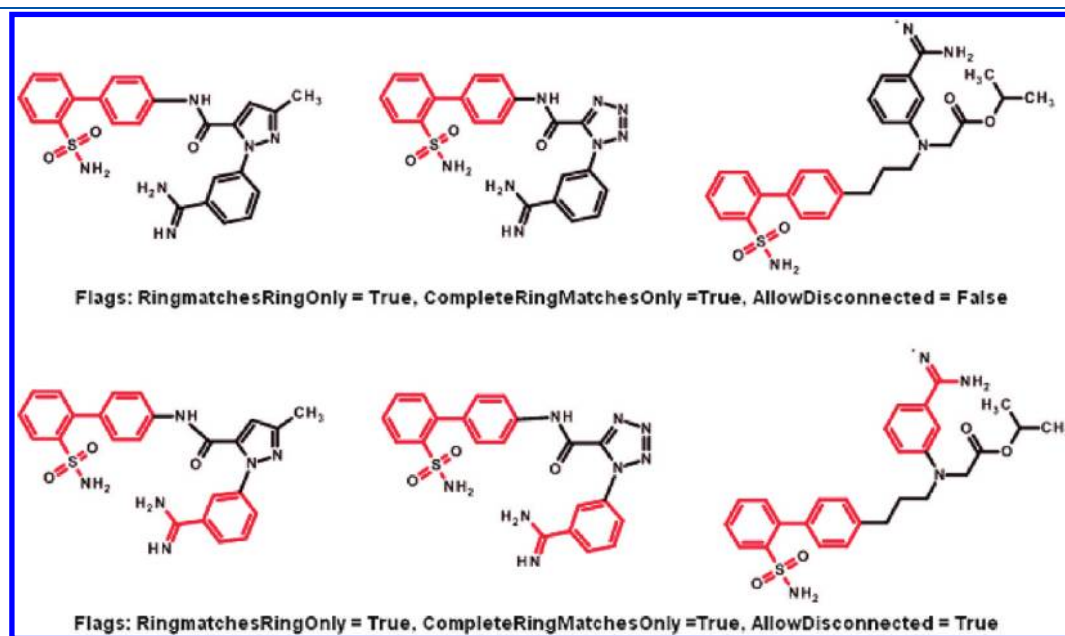


Figure 13. MCS calculation on a set of 51 factor Xa inhibitors. Settings: ring bonds can only match ring bonds, and partial ring matches are not allowed. The upper row of structures shows the connected case, and the lower row shows the disconnected case. The MCS is colored red.

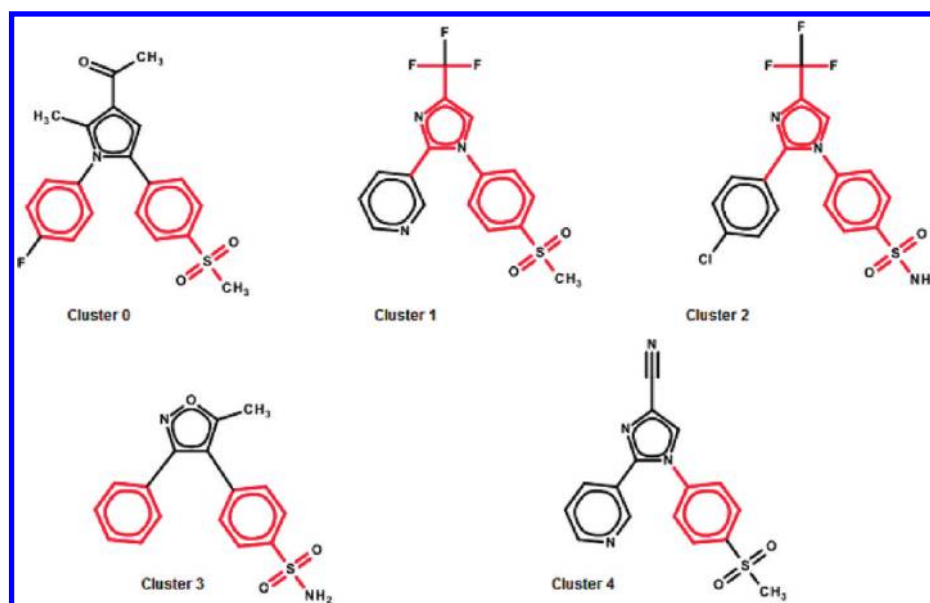


Figure 15. Clustering MCS calculation on a COX2 data set of 322 compounds. Settings: minimum MCS fragment size = 4, minimum MCS size = 4, minimum cluster size = 10, ring bonds can only match ring bonds, partial ring matches are not allowed, and disconnected MCS is allowed.

analog) and plotted MCS run time as a function of the set size. The resulting plot, not shown here, is perfectly linear, i.e., doubling the number of molecules doubles the run time.

Finding the MCS in Sets with Thousands of Molecules: Clustering. The exploratory chemist is often presented with a large number of potential actives from which to choose. This is often the case when a high-throughput screen is run. The potential actives have to be analyzed, and one (or a few series) have to be selected as leads for the chemists to work on. This is a tedious process that involves, at a minimum, the following steps:

- Applying various screens to eliminate compounds that are predicted to have undesirable properties: clogP, toxicity (hERG inhibition, cyp inhibition, AMES), bioavailability, etc.
- Grouping the compounds into clusters in a way that facilitates browsing.
- Examining the clusters for synthetic feasibility, patentability etc.

Because the above process includes manual examination of structures, the clustering step (b) is crucial. Without it, manual examination would be nearly impossible. Typically, chemists use conventional clustering (Downs and Barnard²²) or other scaffold-based methods (Nilakantan et al.²³) to group the compounds. MultiMCS, when run in clustering mode, can group a diverse set of compounds into clusters on the basis of MCS. As mentioned earlier, Stahl et al.,¹⁴ and the ChemAxon¹³ group have also developed MCS-based clustering. We have run MultiMCS on a number of test cases and ensured that the program handles a very wide variety of structural types, including large compounds with about 100 non-hydrogen atoms, large fused ring structures, ring-of-rings etc. We discuss a few examples here.

A. Medium-Sized Data Set of Related Compounds. We used a public domain collection of 322 COX2 compounds from Sutherland et al.²⁴ for this example. These were grouped into five clusters and one singleton as shown in Figure 15. It can be seen that all the compounds in the clusters contain a benzene sulfonamide or a benzene methyl-sulfone, a second benzene ring, and a connecting bridge. Variations in these features cause

Table 1. Results from a Clustering MCS Run on a COX2 Data Set from Sutherland et al.^{24a}

cluster	number of compounds	size of MCS
0	148	16
1	16	20
2	25	20
3	115	16
4	17	7
singletons	1	0

^a The settings for this run were: minimum MCS fragment size = 4, minimum MCS size = 4, minimum cluster size = 10. Ring bonds can only match ring bonds, partial ring matches are disallowed, and disconnected MCS is allowed. Time taken for this run is 2172 ms.

the compounds to be grouped into the five clusters mentioned. This example shows how MultiMCS can be used by the chemist to quickly distill the structural variation contained in 322 compounds into a single page of easily digestible information. The run took just over 2 s. The detailed statistics are shown in Table 1.

B. Large Diverse Collection. For this example, we used a subset of 1034 compounds from the NCI collection. These are compounds tested against a p388 leukemia cell line. MultiMCS grouped these into 9 clusters and 849 singletons. The detailed statistics are shown in Table 2, and a representative member of each cluster is shown in Figure 16.

It is clear from Figure 16 that the clustering procedure finds quite diverse clusters, from 10 compounds to 73 compounds, and with different MCS sizes from 10 bonds to 40 bonds. We would like to point out that our method is very efficient. For example, in cluster 0, there are 12 compounds with a large shared substructure of 39 bonds, and the MCS is in two disconnected fragments. Despite the complexity, the entire job ran in just under a minute. The clusters also illustrate the structural variety of the compounds present in this large data set, from large suramin-like compounds to much simpler compounds.

Note that there are a large number of singletons, and questions to ask are whether these can be reduced via parameter modification and whether clusters hidden in these singletons can somehow be recovered. To this end, we display the number of clusters found and number of singletons found as a function of the two main parameters, MinClusterSize and MinClusterMCSSize, in Figure 17.

As shown in Figure 17, the number of clusters broadly reduces and the number of singletons broadly increases as the MinClusterSize parameter increases. The number of singletons increases

Table 2. Results from a Clustering MCS Run on a Subset of the NCI Compounds Screened against a p388 Leukemia Cell Line^a

cluster	number of compounds	size of MCS
0	12	39
1	73	10
2	17	14
3	11	30
4	12	16
5	19	14
6	10	14
7	17	40
8	14	11
singletons	849	0

^aThe settings for this run were: minimum MCS fragment size = 10, minimum MCS size = 10, minimum cluster size = 10. Ring bonds can only match ring bonds, partial ring matches are disallowed, and disconnected MCS is allowed. Time taken for this run is 57094 ms.

as the MinClusterMCSSize parameter increases; this suggests that one can reduce the number of singletons by reducing these two parameters. However, in this data set, the number of clusters found does not seem to show simple monotonic behavior for the

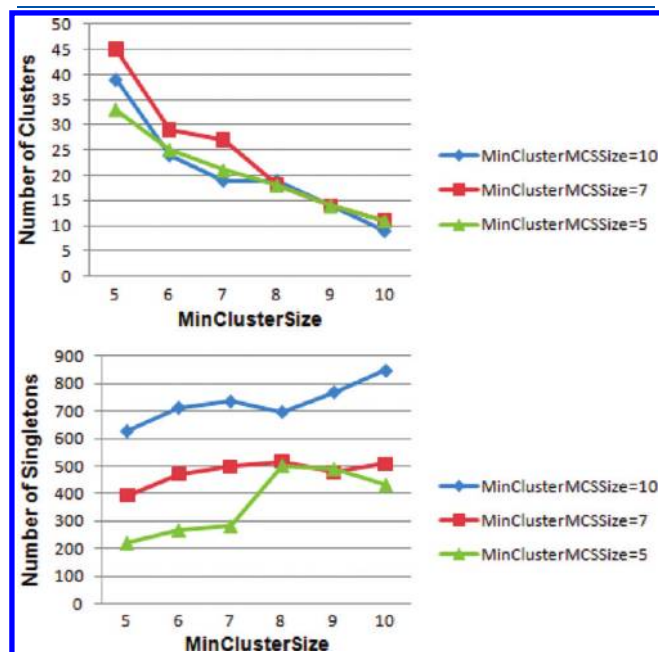


Figure 17. Number of clusters and number of singletons obtained for the 1034 NCI compounds screened against a p388 leukemia cell line, as a function of the MinClusterSize parameter, for three different values of the MinClusterMCSSize parameter.

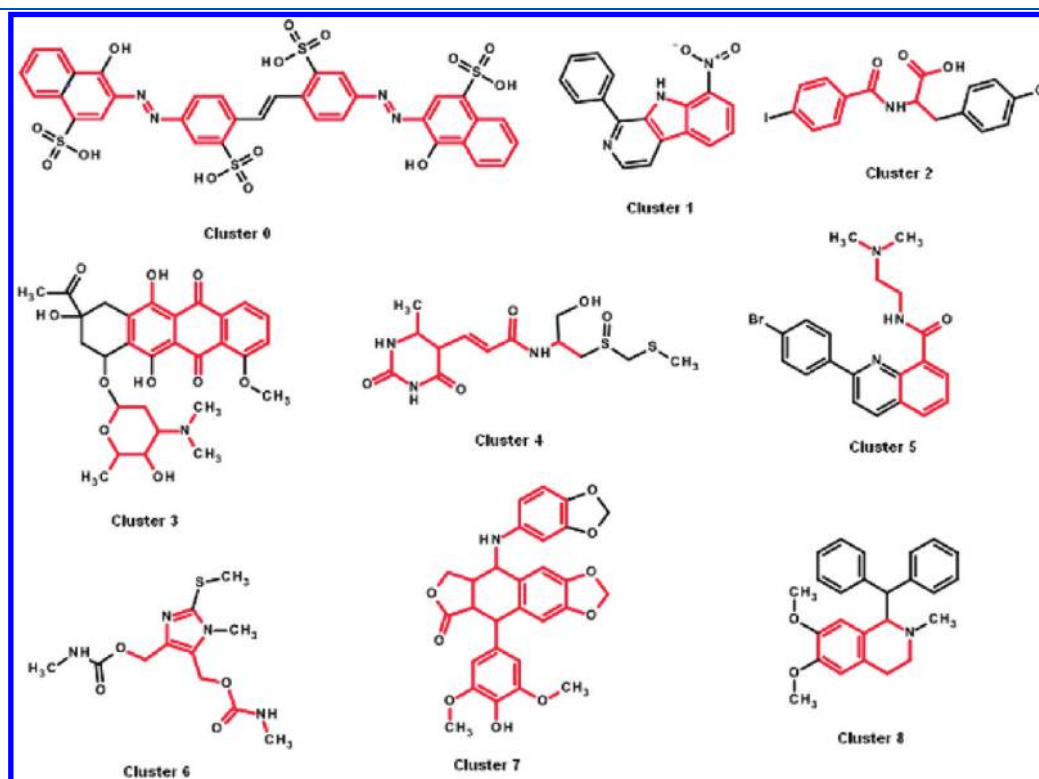


Figure 16. Clustering MCS calculation on a set of 1034 NCI compounds. Settings: minimum MCS fragment size = 10, minimum MCS size = 10, minimum cluster size = 10, ring bonds can only match ring bonds, partial ring matches are not allowed, and disconnected MCS is allowed.

Table 3. Results from a Clustering MCS Run on a Data Set from Jorissen et al.^{25a}

cluster	number of compounds	size of MCS
0	9	18
1	8	17
2	8	20
3	22	15
4	7	22
5	6	19
6	18	18
7	7	15
8	5	25
9	5	23
10	11	19
11	15	18
singletons	2021	0

^aThe settings for this run were: minimum MCS fragment size = 10, minimum MCS size = 15, minimum cluster size = 5. Ring bonds can only match ring bonds, partial ring matches are disallowed, and disconnected MCS is allowed. Time taken for this run is 92516 ms.

chosen values of 5, 7, and 10. Of course, an ideal clustering algorithm will discover the ideal values of these parameters from the data itself and also dynamically modulate these parameters for different sections of the data, i.e., rerun the clustering procedure on the singletons with lower values of these two parameters. We leave this extension of our algorithm as work for the future.

C. Large Set Containing Actives against Different Biological Targets Plus Decoy Structures. For this example, we used a set of compounds from Jorissen et al.²⁵ This set consists of CDK2 inhibitors, COX2 inhibitors, FXa inhibitors, PDE5 inhibitors, A1A antagonists, and several decoy structures: a total of 2142 structures. In large sets such as these, as mentioned previously, the resulting clusters often depend heavily on the input parameters. If we use a small value for minimum MCS size, we can end up with several large clusters with small, uninteresting MCSs, but few singletons. On the other hand, if we use a large value for minimum MCS size, we end up with several smaller clusters, but with larger more interesting MCSs, and many more singletons. To illustrate this, we ran MultiMCS with two sets of values as shown in Table 3.

Case 1: Minimum MCS Size = 15, Minimum Initial Cluster Size = 5. With the above parameters, MultiMCS produced 12 clusters and 2021 singletons. The details are shown in Table 3, and one representative structure from each cluster is shown in Figure 18. It is interesting that these clusters cover a wide variety of structural types with MCS varying in size from the minimum value of 15 bonds to 25 bonds. Moreover, this large data set was processed in about 1.5 min.

Case 2: Minimum MCS Size = 4, Minimum Initial Cluster Size = 10. With these parameters, MultiMCS produced just 5 clusters, plus 508 singletons. The MCS size varied from 6 to 12 bonds. Thus, while there were fewer singletons, no cluster had a large MCS. The details are in Table 4, and a representative compound from each cluster is shown in Figure 19.

D. Small Set of Related Structures with a Small MCS. In the earlier clustering examples, we had several hundred to thousands of compounds. In smaller data sets of around 100 compounds,

MultiMCS is typically run without clustering. However, in some of these data sets, there may be distinct congeneric series; when there are many such series in the data set, the resulting MCS may be small to nil. In such cases, it may be advantageous to run MultiMCS in the clustering mode. We illustrate this with a set of 96 collagenase inhibitors from Scozzafava and Supuran.²⁶ MultiMCS grouped these neatly into three clusters with no singletons. The detailed statistics are shown in Table 5, and a representative member of each cluster is shown in Figure 20. It can be seen that the three clusters each have large MCSs: from 16 to 25 bonds in size. A regular MCS run (in nonclustering mode) produces an MCS of just 13 bonds.

E. Mix of Complex Structures. This example illustrates how MultiMCS can handle complex structures in the clustering mode. Calicheamicin and esperamicin are two related natural products with large and complex ring structures. We used a set of 36 analogs of these compounds from PubChem²⁷ for this study. MultiMCS found three clusters, each with a large MCS. There were no singletons. The detailed statistics are shown in Table 6, and a representative compound from each cluster is shown in Figure 21. The runtime was about 43 s.

Comparison with Accelrys Pipeline Pilot. We took a data set of 808 compounds from Lounkine et al.¹⁷ and clustered these into 85 clusters comprising three compounds each using MultiMCS. We then computed the connected MCS within each cluster using MultiMCS (using the element-based atom-typing scheme, i.e., the atom type of an atom is its element type and allowing ring bonds to only match ring bonds and not allowing rings to match partially). As mentioned previously, this was done on a Windows XP system with 2GB RAM and an Intel Core 2 Duo P8400 processor running at 2.26 GHz (note that we used a single threaded implementation). We performed the same computation using Accelrys Pipeline Pilot. This was done on a Windows Server 2003 system with 16GB RAM and an Intel Xeon processor running at 3.16 GHz. We used default parameters (namely, MinimumSubgraphSize = 8, PartialRings = RejectPartialAromaticRings, and AtomAbstraction = TypeAndCharge); the only exceptions were RequiredProportion, which we set to 1 because we wanted to obtain an MCS for all compounds in the set rather than cluster by MCS and MaximumSecondsPerMolecule, which we set to 1500 to prevent possible failures of MCS detection because of time out. We compare the time taken and the MCS sizes obtained using the two methods below. See the Supporting Information for the actual instance descriptions and raw data on timings and MCS sizes.

The time taken by MultiMCS over all 85 instances was just 24 s (this timing holds even when we allow ring bonds to match chain bonds and rings to match partially). The time taken by Accelrys Pipeline Pilot was 503 s. In addition, of the 503 s taken by Accelrys Pipeline Pilot, 355 s were accounted for by just 2 of the 85 instances, with MCS sizes 34 and 40, respectively; this seems to indicate that the Accelrys Pipeline Pilot algorithm is susceptible to steep time growth on some molecules. Clearly, MultiMCS seems to offer a substantial speed advantage.

Considering MCS size next; MultiMCS obtained a larger MCS in 19 of the 85 cases, while Accelrys Pipeline Pilot obtained a larger MCS in 4 of the 85 cases. In the latter four cases, the MCS found by Accelrys Pipeline Pilot had ring bonds matching chain bonds and rings matching partially. If we run MultiMCS allowing ring bonds to match chain bonds and allowing rings to match partially, we obtain results identical to Accelrys Pipeline Pilot in

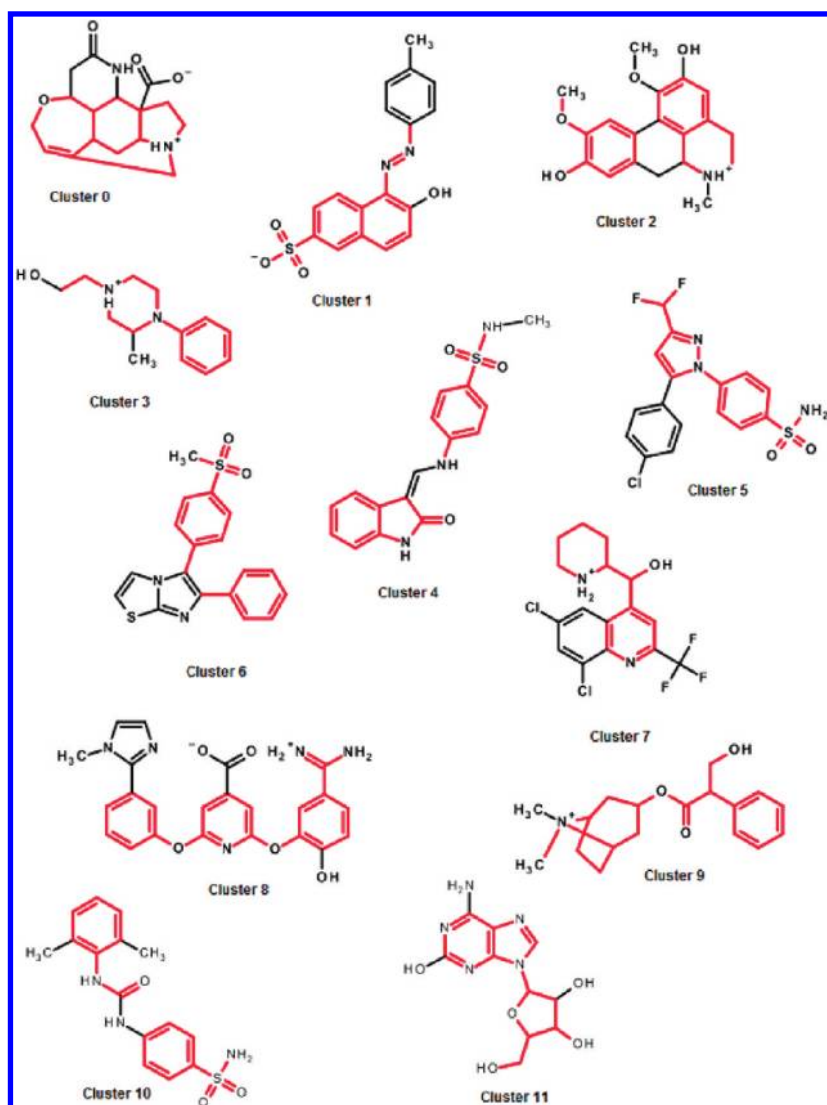


Figure 18. Clustering MCS calculation on a set of 2142 compounds. Settings: minimum MCS fragment size = 10, minimum MCS size = 15, minimum cluster size = 5, ring bonds can only match ring bonds, partial ring matches are not allowed, and disconnected MCS is allowed.

Table 4. Results from a Clustering MCS Run on a Data Set from Jorissen et al.^{25a}

cluster	number of compounds	size of MCS
0	419	7
1	1081	6
2	60	10
3	19	9
4	55	6
singletons	508	0

^a The settings for this run were: minimum MCS fragment size = 4, minimum MCS size = 4, minimum cluster size = 10. Ring bonds can only match ring bonds, partial ring matches are disallowed, and disconnected MCS is allowed. Time taken for this run is 17171 ms.

three of the four cases, and a larger MCS in one case (cluster 66). The analysis of the 19 cases where MultiMCS found a larger MCS was however confounded by a few factors:

1. In 6 of the 19 cases where MultiMCS found a larger MCS (cases 9, 11, 20, 21, 52, 58), the differences were due to

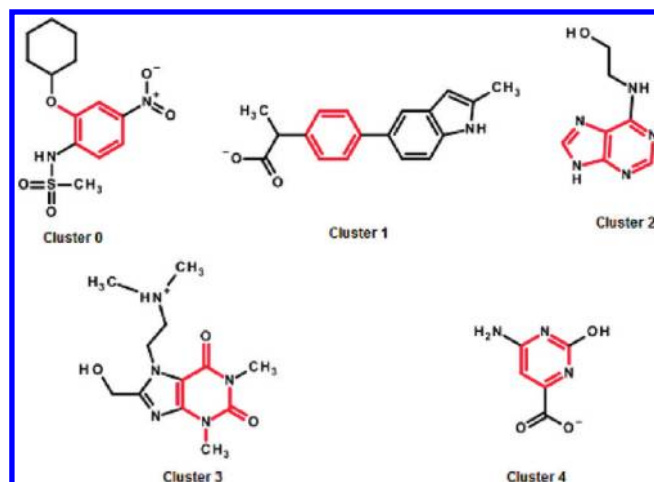


Figure 19. Clustering MCS calculation on a set of 2142 compounds. Settings: minimum MCS fragment size = 4, minimum MCS size = 4, minimum cluster size = 10, ring bonds can only match ring bonds, partial ring matches are not allowed, and disconnected MCS is allowed.

Table 5. Results from a Clustering MCS Run on a Set of Collagenase Inhibitors from Scozzafava and Supuran^{26a}

cluster	number of compounds	size of MCS
0	23	25
1	45	16
2	28	19

^aThe settings for this run were: minimum MCS fragment size = 4, minimum MCS size = 10, minimum cluster size = 10. Ring bonds can only match ring bonds, partial ring matches are disallowed, and disconnected MCS is allowed. Time taken for this run is 546 ms.

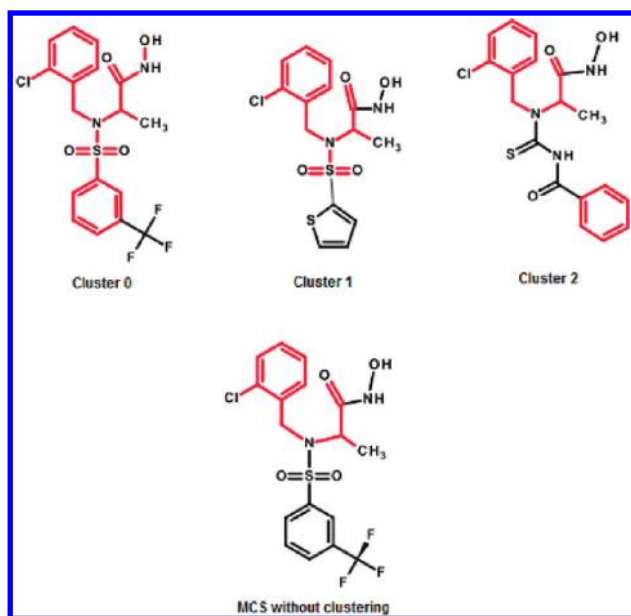


Figure 20. Top row: Clustering MCS calculation on a set of 96 compounds. Settings: minimum MCS fragment size = 4, minimum MCS size = 10, minimum cluster size = 10, ring bonds can only match ring bonds, partial ring matches are not allowed, and disconnected MCS is allowed. Bottom: The overall MCS of all 96 compounds. Note that not all molecules in cluster 1 contain the hydroxylamine fragment; hence, the relevant parts are missing from the overall MCS and the MCS for cluster 1.

Table 6. Results from a Clustering MCS Run on a Set of calicheamicin and esperamicin Analogs from PubChem²⁷

cluster	number of compounds	size of MCS
0	20	93
1	10	54
2	6	88

^aThe settings for this run were: minimum MCS fragment size = 4, minimum MCS size = 4, minimum cluster size = 5. Ring bonds can only match ring bonds, partial ring matches are disallowed, and disconnected MCS is allowed. Time taken for this run is 43234 ms.

inclusion of charge states in atom typing by Accelrys Pipeline Pilot.

- In another 7 of the 19 cases (cases 26, 39, 50, 61, 62, 63, 82), MultiMCS included a bond of type *a* (as defined in Figure 22) in the MCS that Accelrys Pipeline Pilot did not. Additionally, case 50 is interesting in that MultiMCS included two *a* type bonds in the MCS, while Accelrys

Pipeline Pilot included only one. It is not clear whether Accelrys Pipeline Pilot excludes such bonds as a deliberate rule or whether the above cases represent incorrect results. Another odd case is 57, where a bond incident on a saturated ring is included in the MCS by Accelrys Pipeline Pilot. Note that MultiMCS's definition of rings matching partially is bond-based and not atom-based, and therefore it does include such bonds in the MCS even when rings are not allowed to match partially. If needed, MultiMCS can easily be modified to exclude such bonds with negligible change to running time.

- In 2 of the 19 cases where MultiMCS found a larger MCS (cases 36 and 47), it included a bond in the MCS that Accelrys Pipeline Pilot did not; this bond is of type *b* as defined in Figure 22. Again, it is not clear whether Accelrys Pipeline Pilot excludes such bonds as a deliberate rule or whether the above cases represent incorrect results; for instance, in case 26, the MCS obtained by Accelrys Pipeline Pilot does indeed contain a *b* type bond, as shown in Figure 22.
- This leaves four cases where MultiMCS found a larger MCS (cases 57, 72, 73, 74); in all these cases, the MCS found by Accelrys Pipeline Pilot seems clearly incorrect. Three bonds are missed in case 57, and a C–Cl bond is missed in each of cases 72 and 74. In case 73, two fused rings are part of the MCS found by MultiMCS, but the sole bond shared by the two rings is (correctly) excluded. Accelrys Pipeline Pilot completely misses these rings, leading to an MCS that is smaller by 11 bonds. One could argue that an MCS comprising such structures (i.e., fused rings without the shared bond) is not very relevant. If so, MultiMCS can easily be modified to disallow such ring matches. In such a case, MultiMCS would still include one member of the fused ring pair in the MCS, thus yielding an MCS that is 7 bonds larger than that found by Accelrys Pipeline Pilot.

In conclusion, MultiMCS seems to run faster and find larger MCSs in several cases as compared to Accelrys Pipeline Pilot.

Comparison with Library MCS. Finally, we give an example to compare our algorithm with Library MCS. We focus on the overall MCS of the given collection of compounds rather than the clustering aspect. Library MCS, as described previously, couples MCS detection with clustering and produces a dendrogram. If the collection of compounds all share a common substructure, that substructure will be the MCS, and it will be at the root of the dendrogram. Because Library MCS does not do a strict MCS calculation, there are instances where it fails to find the correct MCS of the entire collection. Figure 23 shows one such case. The correct MCS of the three compounds shown in the figure is obviously the imidazole. This is correctly detected by MultiMCS. Library MCS on the other hand, clusters the two larger compounds, finds their MCS, and then compares it to the third compound, and determines that there is no MCS that spans all three compounds. The parameters used for Library MCS were the following: normal mode, minimal MCS size 3, only atom and bond type matching parameters (i.e., charge and hybridization state are ignored). For MultiMCS, we used default parameters. Note that our claim is simply the following: a hierarchical clustering procedure that relies on computing pairwise MCS may not be able to determine the common MCS of multiple compounds, and therefore newer ideas are needed for this task.

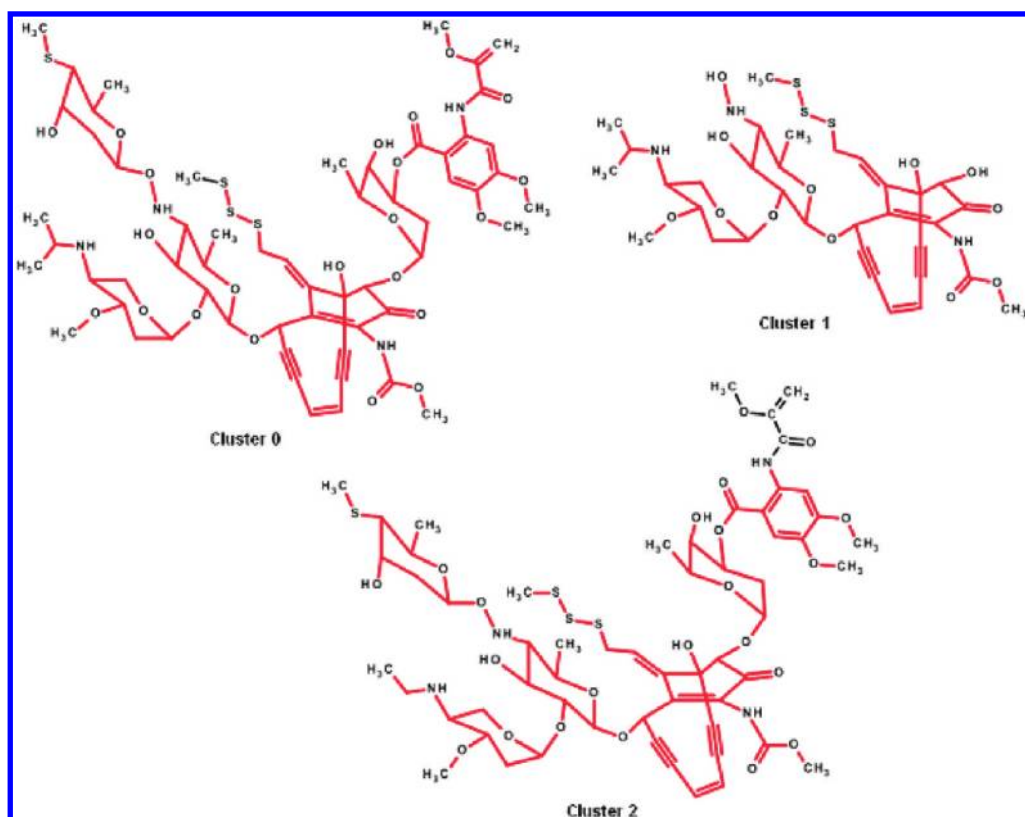


Figure 21. Clustering MCS calculation on a set of 36 analogs of calicheamycin and esperamycin. Settings: minimum MCS fragment size = 4, minimum MCS size = 4, minimum cluster size = 5, ring bonds can only match ring bonds, partial ring matches are not allowed, and disconnected MCS is allowed.

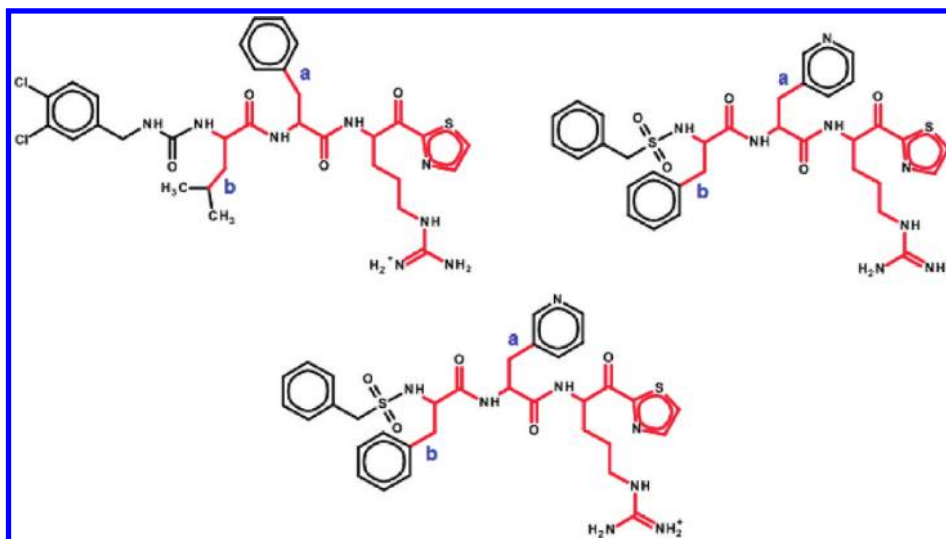


Figure 22. MCS found by MultiMCS for Case 26; the MCS found by Accelrys Pipeline Pilot does not have bond *a*, but does have bond *b*. Note that *a* differs from *b* in that one end point of *a* is a ring atom in all molecules, while one end point of *b* is a ring atom in some but not all the molecules. In both cases, the ring containing this ring atom is itself not in the MCS.

■ DEFAULT ATOM-TYPING SCHEME

We have implemented an atom-typing scheme following the work of Hattori et al.,²⁸ as shown in Figure 24. We have restricted this atom-typing scheme to carbon atoms alone; all other atoms are typed on the basis of their respective element types, i.e., all nitrogen atoms get the nitrogen atom type, and all oxygen atoms get the oxygen atom type, and so on. This atom typing serves as the default, and users retain the ability to use their own custom



Figure 23. Three compounds for the Library MCS comparison.

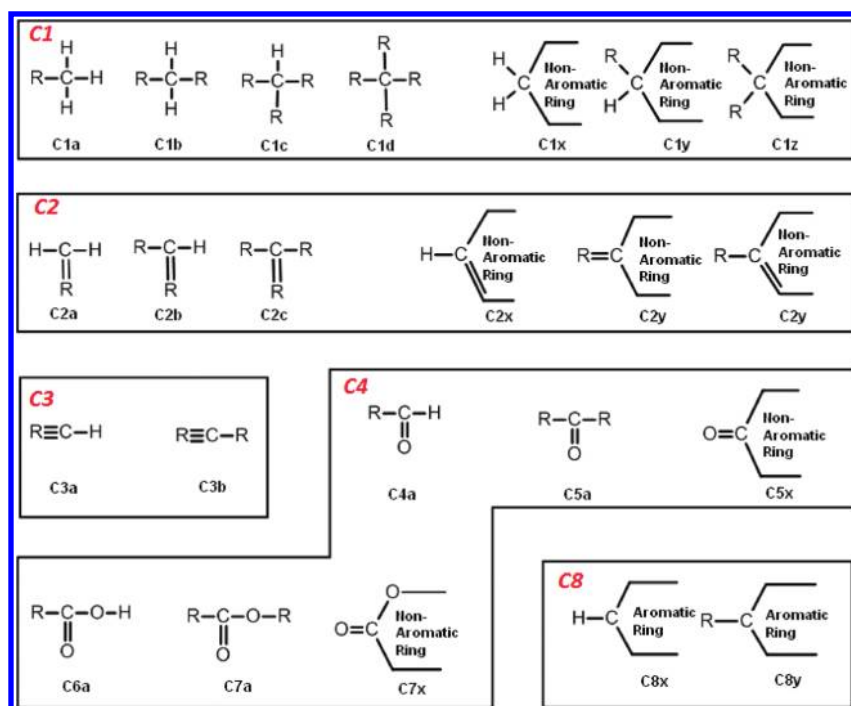


Figure 24. Default atom-typing scheme.

atom types or use just the *elemental* atom type that gives the same atom type to all carbon atoms as well. The reason for subtyping carbon atoms is to improve the speed of MCS computation quite substantially, while also enabling one to filter out irrelevant matches. Subtyping other atoms has a lesser impact on algorithm speed. The carbon subtypes we use are indicated in red.

CONCLUSION

We have designed and implemented a fast and flexible algorithm, MultiMCS, for computing the MCS of multiple molecules. MultiMCS guarantees correctness of MCS for a large class of molecules. The speed of MultiMCS is obtained via the application of a novel divide-and-conquer strategy. MultiMCS can be run on small sets with tens or hundreds of compounds or on larger sets with thousands of compounds, in which case it is coupled with a clustering step. Both connected and disconnected MCS can be discovered by the program. We have tested MultiMCS on many test data sets and have shown that it can handle a very wide variety of cases, including large numbers of complex structures.

ASSOCIATED CONTENT

S Supporting Information. Invocation flags and parameters are described in *mcsInvocation.pdf*. Details of comparison with Accelrys Pipeline Pilot appear in *ppComparison.xlsx*. Input structure data files for all the results reported in this paper appear in *structures.zip*. This file also contains a *Structure-README.xlsx* file that describes each of the structure files. This information is available free of charge via the Internet at <http://pubs.acs.org>

AUTHOR INFORMATION

Corresponding Author

*E-mail: ramesh@strands.com.

ACKNOWLEDGMENT

We thank the anonymous referees for their comments.

REFERENCES

- (1) Armitage, J. E.; Lynch, M. F. Automatic detection of structural similarities among chemical compounds. *J. Chem. Soc. C* **1967**, 521–528.
- (2) Armitage, J. E.; Crowe, J. E.; Evans, P. N.; Lynch, M. F.; McGuirk, J. A. Documentation of chemical reactions by computer analysis of structural changes. *J. Chem. Doc.* **1967**, 7, 209–215.
- (3) Levi, G. A note on the derivation of maximal common subgraphs of two directed or undirected graphs. *Calcolo* **1973**, 9 (4), 341–352.
- (4) Cone, M. M.; Venkataraghavan, R.; McLafferty, F. W. Molecular structure comparison program for the identification of Maximal Common Substructures. *J. Am. Chem. Soc.* **1977**, 99 (23), 7668–7671.
- (5) Raymond, J. W.; Gardiner, E. J.; Willett, P. RASCAL: Calculation of graph similarity using a maximum common edge subgraph algorithm. *Comput. J.* **2002**, 45 (6), 631–644.
- (6) Raymond, J. W.; Gardiner, E. J.; Willett, P. Heuristics for similarity searching of chemical graphs using a maximum common edge subgraph algorithm. *J. Chem. Inf. Comput. Sci.* **2002**, 42 (2), 305–316.
- (7) Pardalos, P.; Xue, J. The maximum clique problem. *J. Global Optim.* **1992**, 4, 301–328.
- (8) Wood, D. An algorithm for finding a maximum clique in a graph. *Oper. Res. Lett.* **1997**, 21, 211–217.
- (9) Bron, C.; Kerbosch, J. Finding all cliques of an undirected graph. *Commun. ACM* **1973**, 16 (9), 575–577.
- (10) Varkony, T. H.; Shiloach, Y.; Smith, D. H. Computer-assisted examination of chemical compounds for structural similarities. *J. Chem. Inf. Comput. Sci.* **1979**, 19 (2), 104–111.
- (11) *Pipeline Pilot*, version 7.5; Accelrys: San Diego, CA, 2008.

- (12) Takahashi, Y.; Satoh, Y.; Sasaki, S-I. Recognition of largest common structural fragment among a variety of chemical structures. *Anal. Sci.* **1987**, *3*, 23–28.
- (13) *Library MCS*, version 0.7; Chemaxon: Budapest, Hungary, 2008.
- (14) Stahl, M.; Mauser, H.; Tsui, M.; Taylor, N. R. A robust clustering method for chemical structures. *J. Med. Chem.* **2005**, *48*, 4358–4366.
- (15) Whitney, H. Congruent graphs and the connectivity of graphs. *Am. J. Math.* **1932**, *54*, 150–168.
- (16) Garey, M.; Johnson, D. *Computers and Intractability: A Guide to the Theory of NP-Completeness*; W.H. Freeman: New York, 1979; ISBN 0-7167-1045-5.
- (17) Lounkine, E.; Wawer, M.; Wassermann, A. M.; Bajorath, J. SARANEA: A freely available program to mine structure–activity and structure–selectivity relationship information in compound data sets. *J. Chem. Inf. Model.* **2010**, *50*, 68–78.
- (18) Gedeck, P.; Rohde, B.; Bartels, C. QSAR: How good is it in practice? Comparison of descriptor sets on an unbiased cross section of corporate data sets. *J. Chem. Inf. Model.* **2006**, *46* (5), 1924–1936.
- (19) Borgelt, C.; Meinl, T.; Berthold, M. R. *MoSS: A Program for Molecular Substructure Mining*, Workshop on Open Source Data Mining Software (OSDM), 2005, pp 6–15.
- (20) Agrawal, R.; Imieliński, T.; Swami, A. *Mining Association Rules between Sets of Items in Large Databases*, Proceedings of the ACM SIGMOD International Conference on Management of Data, Washington, DC, May 26–28, 1993, pp 207–216.
- (21) Fontaine, F.; Pastor, M.; Zamora, I.; Sanz, F. Anchor-GRIND: Filling the gap between standard 3D QSAR and the GRid-INdependent Descriptors. *J. Med. Chem.* **2005**, *48* (7), 2687–2694.
- (22) Downs, G. M.; Barnard, J. M. Clustering Methods and Their Uses in Computational Chemistry. In *Reviews in Computational Chemistry*; Wiley-VCH: New York, 2002, pp 1–40.
- (23) Nilakantan, R.; Bauman, N.; Haraki, K. S. Database diversity assessment: New ideas, concepts, and tools. *J. Comput.-Aided Mol. Des.* **1997**, *11* (5), 447–452.
- (24) Sutherland, J. J.; O'Brien, L. A.; Weaver, D. F. Spline-Fitting with a genetic algorithm: A method for developing classification structure–activity relationships. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1906–1915.
- (25) Jorissen, R. N.; Gilson, M. K. Virtual screening of molecular databases using a Support Vector Machine. *J. Chem. Inf. Model.* **2005**, *45* (3), 549–561.
- (26) Scozzafava, A.; Supuran, C. T. Protease Inhibitors: Synthesis of potent bacterial Collagenase and matrix metalloproteinase inhibitors incorporating N-4-Nitrobenzylsulfonylglycine hydroxamate moieties. *J. Med. Chem.* **2000**, *43* (9), 1858–1865.
- (27) Wang, Y.; Xiao, J.; Suzek, T. O.; Zhang, J.; Wang, J.; Bryant, S. H. PubChem: A public information system for analyzing bioactivities of small molecules. *Nucleic Acids Res.* **2009**, *37*, 623–633.
- (28) Hattori, M.; Okuno, Y.; Goto, S.; Kanehisa, M. Heuristics for chemical compound matching. *Genome Inf.* **2003**, *14*, 144–153.