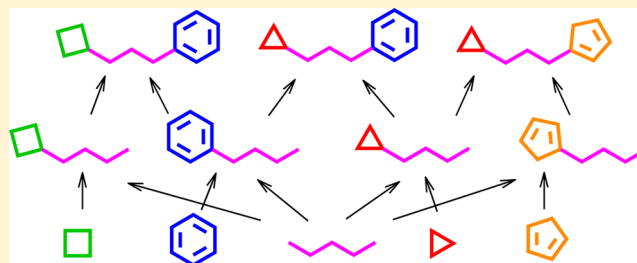


# Construction and Use of Fragment-Augmented Molecular Hasse Diagrams

Peter Lind\*

Medivir AB, Box 1086, 14122 Huddinge, Sweden

**ABSTRACT:** Collections of molecules can be organized in many different ways based on substructures that are common to two or more of the molecules. The article describes a method that builds on the ideas of partial orders and Hasse diagrams and which organizes molecules in a particularly simple and natural way using only sub- and superstructure relations. The method outputs the original molecule collection together with common substructures and a set of relations between fragments and molecules. The result is a complete deconstruction of the original structures into those fragments or building blocks that are shared between two or more molecules. Scaffolds for the R-group analyses that can be performed on the data set are automatically detected. Cyclic and linear substituents are treated in the same way. No rules are incorporated that express any form of domain expertise or judgment. The method should be useful for library profiling, data set navigation, fragment-based screening, identification of activity cliffs, and identification of library subsets that are amenable to fragment-based QSAR.



## ■ INTRODUCTION

The goal of a medicinal chemistry project is to identify one or more candidate drugs which meet a target profile specification. The development process starts with the identification of hit and lead compounds. These compounds are structurally modified and tested in iterations in order to improve compound properties such as activity, selectivity, and bioavailability while also considering synthesizability and patentability. The influence that molecular structure has on these properties is usually understood and communicated in terms of fragment concepts like functional groups, scaffolds, and substituents. Drug design supporting software will then have an advantage if it can express its results in terms of fragments, because team members representing different disciplines can then all appreciate results and discuss how to best simultaneously optimize the different target compound properties. Medicinal chemistry software which is not fragment-based may relate results to more or less obscure predictor variables which must first be calculated from the molecules. Such two-step procedures are appropriate when compounds contain untested fragments but are otherwise harder to use in compound design because the relation between structure and predicted target value is indirect.

Chemistry software which uses structural fragments for calculation and result presentation has been used to support many tasks. Examples are design and profiling of libraries and development of structure–activity relationship (SAR) models. Fragment-based software used for such purposes will typically organize fragments or molecules into classes, clusters, or hierarchical structures. These classes or structures are displayed in a table or chart format to visualize characteristics of the set, or they are used in regression models where the class

membership of a molecule can be predictive of some activity or property.

The Methods section of this paper describes software for the construction of fragment-augmented molecular Hasse diagrams. These diagrams are drawings in which molecules and fragments are organized in a simple and natural way, using only sub- and superstructure relations. The diagrams build on the ideas of partial orders and Hasse diagrams. The standard definitions of these concepts are given in a separate section.

## ■ PREVIOUS WORK

**Clustering.** The most important method that has been used to organize sets of molecules is clustering. Clustering methods in general organize sets by partitioning their elements into classes. A distance function is used to evaluate proximities or similarities between elements in the set. Elements which are close by the metric are then grouped together into classes. Distance functions based on many types of molecular properties and descriptors have been used to cluster molecules.<sup>1,2</sup>

**Fingerprint Based Clustering.** A type of distance metric which is often used for clustering of chemical structures uses molecular fingerprints.<sup>3,4</sup> Fingerprints are bit vectors where an individual bit expresses the presence or absence of a structural feature. A typical procedure for fingerprint creation will start from a graph representation of the atom connectivity table. Then all paths and cycles within predetermined size ranges are enumerated. Atoms and bonds are assigned numeric values

Received: July 29, 2013

depending on hybridization and local environments, and hash values for each path or cycle are calculated based on these bond and atom values. The hash values correspond to positions in the bit vector, and bits are set to 1 at those positions. Fingerprints created by such procedures will be based on a mapping from fragments to bits, but the mapping from fingerprints back to specific fragments is lost because of the hashing step. Other types of fingerprints will not use hashing, and their chemical interpretation will be clearer. Fingerprint-based clustering methods will group together compounds that have similar bit strings. This means that compounds in a specific cluster will differ from compounds in other clusters in that certain bits will be on more often. There will be no cluster defining set of bits where all bits are strictly required to be 1. The main advantage of using fingerprints for clustering is that calculations are very fast,<sup>5</sup> in particular if bit-level parallel hardware can be used.<sup>6</sup> Because of this, clustering based on 2D similarity is usually performed via fingerprints, and clusters can be characterized by the selection of a few representative compounds from each of them.<sup>7</sup> Fingerprint classification results can be different than those from human classifications because fingerprints are in general not sensitive to the relative position and orientation of chemical groups separated by a distance.<sup>8</sup>

Direct comparison of graph representations of chemical structures is computationally expensive and therefore seldom used for the calculation of similarity metrics.

#### MCS and Scaffold Based Classification, Scaffold Trees.

The development of methods that structures sets of chemical compounds based on functional groups and fragments started in the late 1990s. Cosgrove and Willet created the program SLASH<sup>9</sup> with the purpose of giving a measure of the probability that a particular fragment is contributing to activity. The program uses predefined rules to extract fragments from molecules in the data set. The fragments, which may overlap, are sorted into three classes: rings, functional groups, and chains. Fragments are described in a hierarchy of levels of detail so that two fragments may match on a less specific description level while they do not match on a more specific description level.

The commercial program Leadscape<sup>10</sup> organizes compounds in a library using a number of predefined “structural features,” which are functional groups and substituents. Features are arranged in a hierarchy following an idea similar to that of the SLASH program. In a similar spirit, Xu later proposed a Scaffold-based Classification Approach (SCA),<sup>11</sup> which groups compounds into the same class if they share the same scaffold. The first step in the procedure of Xu is to find scaffolds by removing side chains from compounds. Scaffolds are typically shared between several compounds. Then, in the second step, compounds are assigned class membership values based on similarity to scaffolds representing class centers. A complexity measure is used to sort scaffolds, and that measure is used as one of the position coordinates when compounds are plotted in a chart, with the other coordinate being a class membership value. The overall procedure is similar to that of Leadscape, except that scaffolds are generated in the procedure and are not predefined. Methods that use Maximum Common Substructures (MSCs) for class centers have subsequently been reported by Miller,<sup>12</sup> by Cross et al.,<sup>13</sup> and by Stahl et al.<sup>8</sup> The distinction between scaffolds and MSCs in this context is that scaffolds are molecules with side chains stripped off, while an MSC can be a partial scaffold or a scaffold with a partial side chain.

Nicolaou et al.<sup>14</sup> describes what they call a “phylogenetic-like tree growing algorithm” (PGLT). The input consists of chemical structures with associated biological data. The outcome from calculation is a hierarchy of chemical classes, where each class is represented as a node in a tree. The authors describe their algorithm as being “of a hybrid nature employing various techniques ranging from neural networks and genetic algorithms to expert rules and chemical substructure searching.”

The program HierS, developed by Wilkens et al.,<sup>15</sup> clusters compounds based on their chemical graphs and builds a hierarchical relationship between ring scaffold classes. Molecules are grouped by shared ring structures, which are generated using a scheme similar to that of Bemis and Murcko,<sup>16</sup> which is also used by Xu. A given compound can have more than one scaffold substructure, and each of those scaffolds can also belong to several classes represented by smaller sub scaffolds. A hierarchical structure is formed with complete compounds at the first level, at the second level linked ring systems without side chains, at higher levels partial scaffolds, and at the top level isolated ring systems without linkers or chains. An approach similar to that of HierS was reported by Koch et al., who studied a principle for structural classification of natural products (SCONP).<sup>17</sup>

The commercial program Classpharmer classifies compounds based on a proprietary algorithm which is not disclosed but which has been described<sup>18</sup> as based on the calculation of approximated maximum common substructures between all pairs, triplets, quadruplets, and so on, up to a user defined “homogeneity” level.

Schuffenhauer et al.<sup>19</sup> have published a method for generating scaffold trees which produces class hierarchies similar to those of HierS, but with the important difference that also ring systems can be digested into smaller fragments. A set of 13 rules determines the allowed cleavages and the order of removals.

Lounkine and Bajorath<sup>20</sup> describe a method that generates what they call core trees. A tree represents fragmentation pathways of a structure. The full structure is at the root of the tree. Algorithmically, the method uses iterations of random fragmentation. Substructures that appear in more than one compound in the same activity class are described as “activity class characteristic substructures” if they do not appear combined in compounds outside the class. The output for a compound corresponds to its deconstruction by removal of atoms and fragments. The deconstruction can go along several paths depending on order of removal, and the paths are shown as the visual output. The ends of paths correspond to tree leafs.

Clark<sup>21</sup> extends on the scaffold tree concept by adding a procedure that generates layouts of molecules in tree diagrams. The aim is to create aesthetically appealing diagrams that clarify the relations between fragments. The fragmentation is roughly as in ref 17, but the sequential breaking of fused aromatic rings is collapsed into a single step. Molecules and fragments are oriented in a way that makes it easy to see how fragments in different molecules correspond to each other. Figure 1 shows how this method organizes a set of molecules based on a common constituent fragment.

Cerruela García<sup>22</sup> proposed a visualization method applicable to data sets where all compounds share a maximum common substructure. The root node of a tree represents the most common MSC in the data set, and this is put in the center of the picture. Branches go out from the root in all directions, to nodes that represent larger fragments. Angles and distances

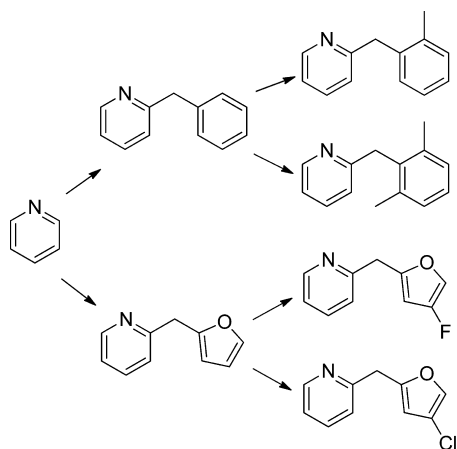


Figure 1. A fragment hierarchy.

between nodes are calculated based on structural similarity or some other notion of proximity.

**Scaffold Networks.** Varin et al.<sup>23</sup> have reported work on scaffold-based organization of sets of chemical structures which extends the concept of scaffold trees. The scaffold tree method was modified so that no rules were applied to select only one particular scaffold per molecule and hierarchy level. Instead, all branches from all scaffolds were included. The method was used to identify active scaffolds from primary screening data where it performed better than the scaffold tree method. Selection rules prioritize scaffolds with few acyclic linker bonds, spiro or bridged ring systems, and ring systems with sizes other than 3, 5, or 6. Figure 2 shows a scaffold network.

**Generation of Fragments.** It is possible to base a classification procedure on a set of predefined fragments, but typically the classification software generates fragments as the process goes along. There are two different ways to extract fragments from a molecule set: The first way is to use rules to disassemble each molecule. One molecule will then give the same fragments regardless of what other compounds are present. The other way is to compare molecules in the set to identify common substructures, and it is the combination of molecules in the set that determines what fragments will be used.

**Rule-Based Fragmentation.**<sup>11,16,21,24,25</sup> A set of fragmentation rules can be applied to each molecule in a set. Each molecule will then be fragmented independently from others. Figure 3 shows the often used protocol for rule-based fragmentation that was published by Bemis and Murcko.<sup>16</sup> Molecules are disassembled into side chains, ring systems, and linkers. Ring systems with linkers are termed frameworks. Another widely used fragmentation scheme is RECAP<sup>25</sup> (Rethrosynthetic Combinatorial Analysis Procedure). This procedure disassembles molecules into building blocks by disconnecting those bonds that can usually be formed with well-known synthetic procedures. Examples are ester, amide, ether, and olefinic bonds.

**Common Substructures.** Another way to create fragments from a set of molecules is to identify substructures common to several molecules.<sup>17,26,27</sup> There are many ways to do this, and practically useful methods will fall in between two extreme possibilities: The first extreme is to collect only those fragments that are shared between strictly all molecules; then only few and very small fragments will be found if the set is diverse. The other extreme is to identify and collect the maximum common

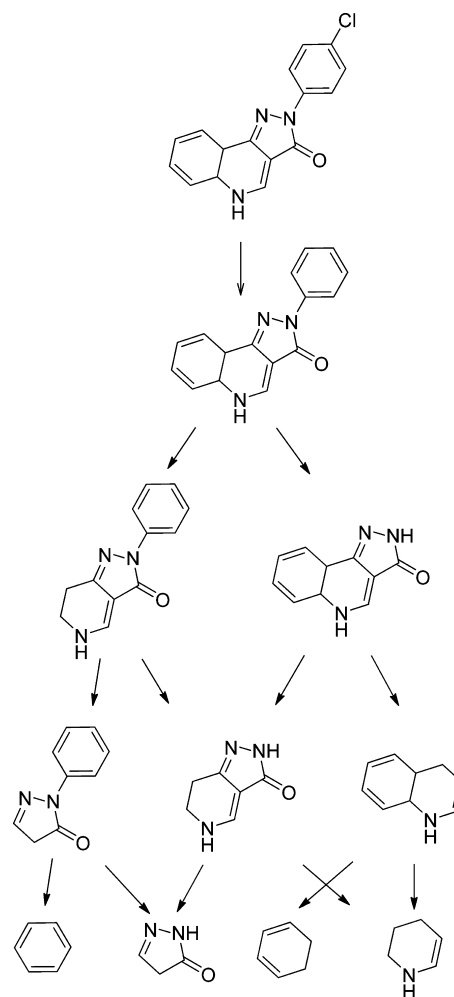


Figure 2. A scaffold network.

substructure (MSC) from each possible pair of molecules. Many common substructures will be found in a diverse set. Realization of a strict version of such a method requires that the maximum common subgraph problem, which is computationally hard,<sup>28</sup> is solved for each molecule pair. It is also possible that the output contains many times more structures than the input, which is certainly not desired when the goal is simplification or systemization. It is clear that practical systemization methods must use fragment selection schemes, which favors selection of substructures that are common to several but not all compounds.

In summary, work on technologies for the organization of sets of molecules has focused on clustering methods and on methods that organize molecules in a hierarchical fashion based on constituent fragments. The methods that construct hierarchies use various pruning rules to identify the set of substructures that are used for compound classification. The rules are more or less complicated and must be learned by anyone who wants to understand results from these methods. Most methods produce tree-like models, but the scaffold network method recognizes that a molecule can have many substructures while at the same time each substructure can occur in several molecules, and it therefore models fragment–molecule relations by using a lattice-like network of common scaffold substructure relationships. The scaffold network method generates fragments by ring pruning, and it uses

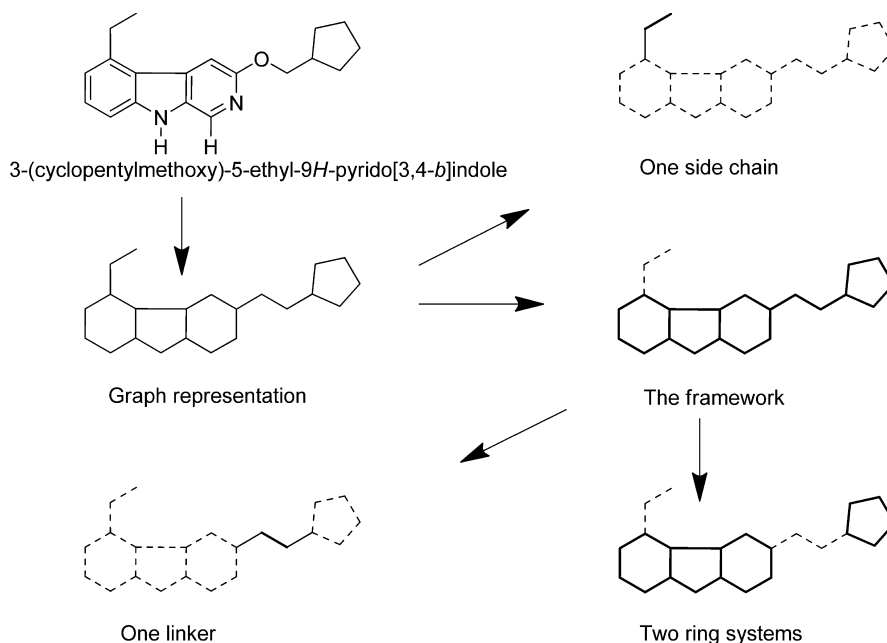


Figure 3. Bemis–Murcko fragmentation.

scaffold selection rules designed to give preference to substructures that are judged to be favored in terms of medicinal chemistry.

## PARTIAL ORDERS AND HASSE DIAGRAMS

**Binary Relation.** A binary relation<sup>29</sup> on a set  $A$  is a subset of all possible ordered pairs of elements from  $A$ .

A partial order is a binary relation " $\leq$ " over a set  $P$  which is antisymmetric, transitive, and reflexive. Antisymmetry means that if  $a \leq b$  and  $b \leq a$  then  $a = b$ ; transitivity means that if  $a \leq b$  and  $b \leq c$  then  $a \leq c$ ; reflexivity means that  $a \leq a$  for all  $a$  in  $P$ . Elements  $a$  and  $b$  are *comparable* if  $a \leq b$  or  $b \leq a$ ; otherwise they are *incomparable*.

A partially ordered set is a set with a partial order. Partial orders formally describe sets where some elements relate to others by, for example, precedence, inclusion, or size.

**Cover Relation.** If elements  $a$  and  $b$  are distinct and there is no element  $x$  for which  $a \leq x \leq b$ , then  $a$  is *covered by*  $b$ , or  $b$  *covers*  $a$ .

**Hasse Diagrams.** Hasse diagrams are pictures that represent partially ordered sets. Hasse diagrams are constructed such that if an element  $a$  precedes an element  $b$  in the partial order, then  $a$  is drawn below  $b$  on the page. A line is drawn from a preceding element  $a$  up to element  $b$  if there is no intermediate element  $x$  so that  $a \leq x \leq b$ . A Hasse diagram fully describes any partially ordered set.

Figure 4 is a Hasse diagram showing the set  $\{A, B, C, D\}$  and all its subsets ordered by subset inclusion. Figure 5 is a Hasse diagram showing the string ABCD and all its substrings ordered by substring inclusion.

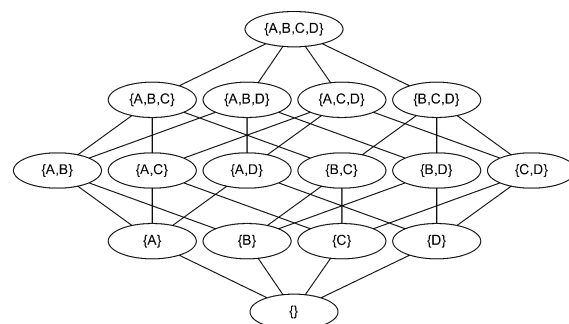


Figure 4. A Hasse diagram showing the set  $\{A, B, C, D\}$  and all its subsets ordered by subset inclusion.

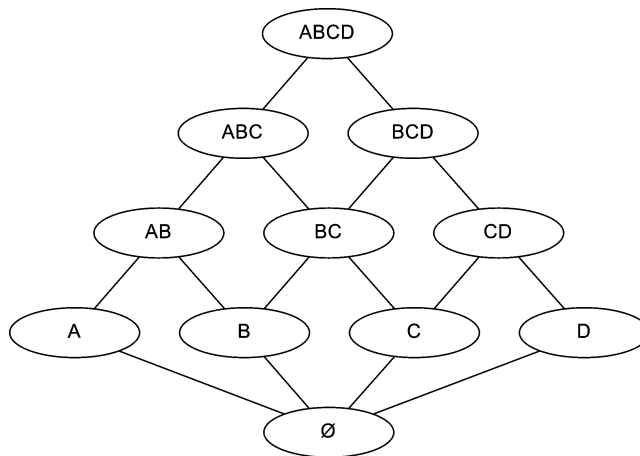


Figure 5. A Hasse diagram showing the string ABCD and all its substrings ordered by substring inclusion.

## METHODS

The present method works by building a computer representation of a molecular Hasse diagram where structures are ordered based on substructure inclusion. The diagram will contain the input molecules along with maximum common substructures (MSCs) which are computed for every new addition. The algorithm is of the online type, meaning that

molecules are added one by one, with all data structures being completely updated after every new addition.

**Software.** The program is written in C#. The source code is deposited in the public domain.<sup>30</sup> Chemistry functionality is implemented using the Indigo toolkit from GGA Software

```
Declare list NodeList;
Declare queue FragmentInsertionQueue;

Procedure Main
for each compound in the input set
create corresponding object NewNode;
AddNode(NewNode);
end for
End procedure

Procedure AddNode(NewNode)
InsertNode(newNode);
while (FragmentInsertionQueue.Count > 0)
    dequeue nodes and insert them using InsertNode();
end while
End procedure

Procedure InsertNode(NewNode)
if a node identical to NewNode already exists in NodeList then return;

make List LowerNodesCollection containing all nodes lower than NewNode;
remove from LowerNodesCollection all nodes having a higher node in that
collection;

make List HigherNodesCollection containing all nodes higher than NewNode;
remove from HigherNodesCollection all nodes having a lower node in that
collection;

remove existing edges between nodes in LowerNodesCollection and
HigherNodesCollection;

add edges from all nodes in LowerNodesCollection up to NewNode;
add edges from NewNode up to all nodes in LowerNodesCollection;

FragmentInsertionQueue = FindMaxCommonSubstructures(NewNode,NodeList);
NodeList.Add(NewNode);
End procedure
```

**Figure 6.** Pseudocode describing how the algorithm inserts new nodes.

Services.<sup>31</sup> Graphviz<sup>32</sup> software is used to render Hasse diagram drawings.

**Data Structures.** Hasse diagram element objects serve to store chemical structures and have lists for upward and downward edges. Unique structures have unique string keys. Edge objects, corresponding to lines in the diagram, have upper and a lower element objects. A Hasse diagram object has methods for insertion and deletions of elements and for calculation of common substructures. An empty root element object is always present in the diagram.

**Algorithm.** The procedure used to maintain the Hasse diagram is outlined in Figure 6. New structures are added one by one. Chemistry objects are instantiated from molfile or other file format and are given a unique key based on the chemical structure. Duplicates are not added. The correct insertion point in the diagram is found by identifying the set of elements that

will cover the new element and the set of elements that will be covered by the new element. The elements to be covered are those which are smaller than the new object and are not covered by any other such element. This set of elements will possibly contain only the root object. The elements that will become covering after insertion are those which are larger than the new object and that are not covering any other such element. This set of elements will possibly be empty. The edges that will exist between elements from these two sets are broken, and new edges are inserted upward from members of the set of elements to be covered by the new element, and upward from the new element to members of the set of elements that will cover the new element.

Maximum common substructures are calculated between the new structure and other structures. They are inserted the same way as full structures. The test version of the algorithm



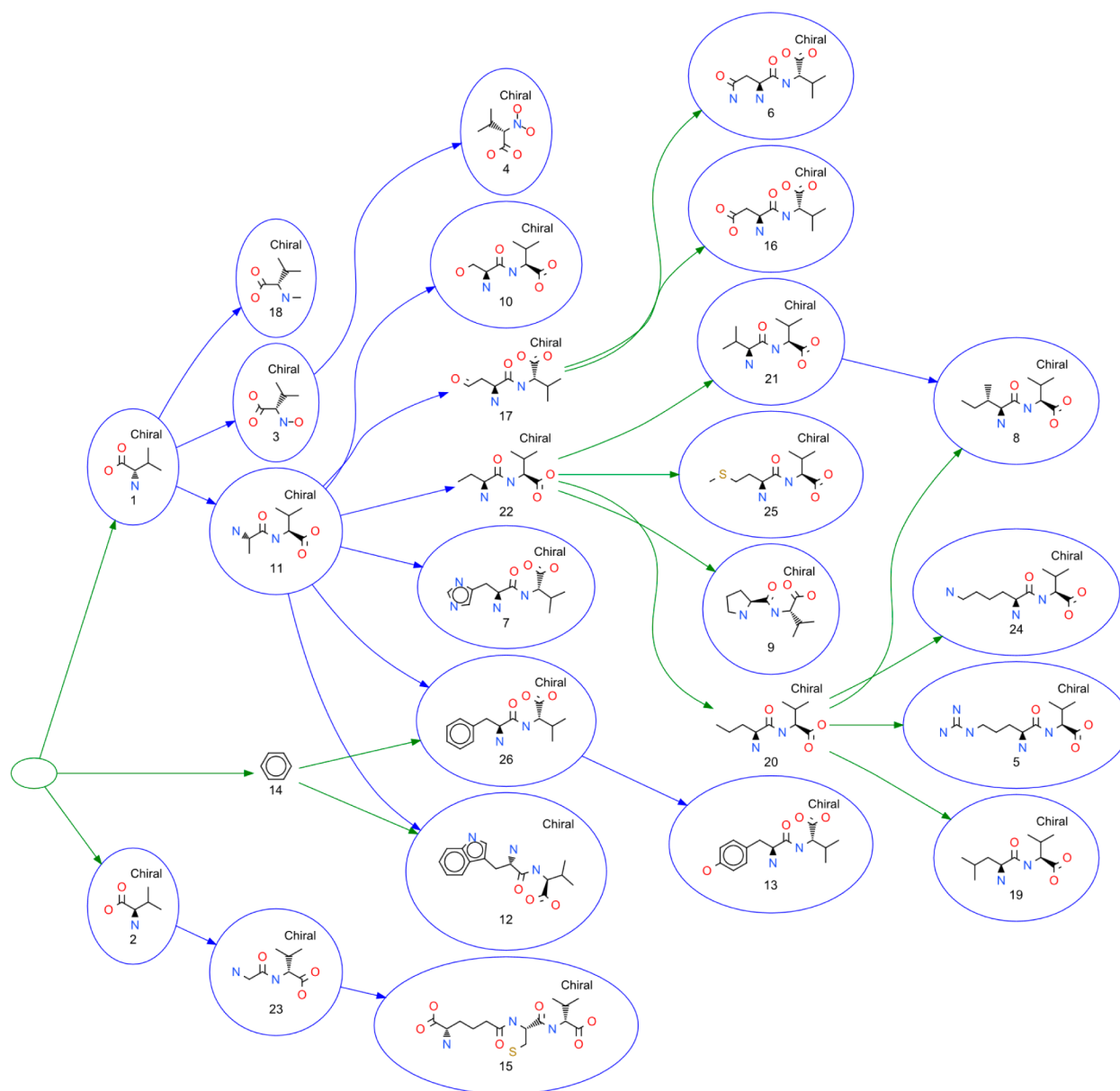


**Figure 7.** A diagram showing substructure/superstructure relations between 32 benzodiazepine inverse agonist site binders.

calculates MSCs between the new structures and all other structures. Most found MCS structures will then be duplicates of already inserted MCSs and are rejected for reinsertion. Future research may show if heuristics can be used to limit the set of structures that are candidates for productive MCS matching against the newly added structures.

Removal of elements follows a reverse procedure.

**Implementation Details That Influence Results.** The present method is based entirely on the concepts of partial ordering by substructure inclusion and maximum common substructures (MSCs). The exact rules that are used to determine if a structure is a substructure of another structure, for example charge and stereo matching principles, or hydrogen atom count rules, will influence the result. Also the precise way



**Figure 8.** A diagram showing relations between 24 valines from the CheBI database.

to calculate maximum common substructures is important, and different MCS calculation implementations can give different results. The Indigo toolkit used in the test implementation has both an exact and an approximate mode for MCS, and the choice may affect the result. The MSC scaffolds were output in the simplest possible format with hydrogen atoms at the substituent attachment points. An alternative version which would give another level of organization could output scaffolds with placeholder symbols, such as R labels, at substitution points and use a rule saying that structures with placeholders are taken to be substructures of the full compounds that they represent. Some common substructures that are generated by the Indigo toolkit were set to be rejected in the tested version of the method; these are chemically nonsensical structures such as partial rings having aromatic bonds. Otherwise no rules are incorporated that express any form of chemistry expertise or judgment.

A conventional substructure determination function is sufficient when molecules and fragments have no special

features, but extended functions are needed if one wants to work with Markush structures or molecules with query features. Diagrams with Markush structures may be useful in the context of fragment combination enumerations.

## RESULTS

Figure 7 shows the substructure–superstructure relations of 32 compounds from the data set of benzodiazepine inverse agonist site binders that was studied by Silverman and Platt.<sup>33</sup> Compounds present in the input data set are shown in circles. Maximum common substructures generated by the program are shown without circles. Green arrows go from common substructures to larger structures. Blue arrows go from real compounds to larger compounds. A setting allowing only common substructures with at least six heavy atoms was used. The diagram is rotated 90° to best fit the page, with the root to the left and larger structures to the right.

Figure 8 shows the relations between 24 valine containing compounds from the CheBI database.<sup>38</sup> This diagram shows

that the method, in contrast to most other methods, does not rely on rings or ring systems for the classification of compounds. One of the extracted common substructures is cyclic (compound 14, benzene), but three noncyclic structures (17, 22, 20) are also extracted as common substructures. The most commonly occurring fragment in the data set is L-valine (1) which in the diagram has paths to all compounds except the three which are derivatives of D-valine (2). The next most important fragment is glycyl-L-valine (11), which is a substructure of the majority of compounds.

**Comparison with Earlier Methodologies.** Some methods produce hierarchical structures representing the buildup of many larger structures from a single fragment by sequential additions.<sup>17,19,21,22</sup> This is exemplified by the fragment hierarchy method<sup>21</sup> shown in Figure 1. Other methods produce diagrams representing the ways that a compound can be disassembled into smaller fragments.<sup>18,20,23</sup> The scaffold network<sup>23</sup> procedure shown in Figure 2 is an example.

The fragment-augmented Hasse diagram methodology extracts fragments that are common to at least two larger fragments or molecules, and it keeps track of paths from smaller structures to their superstructures. This is done in a complete fashion so that all paths between a larger structure and all its substructures are included; also all paths from a substructure to all its superstructures are included. Relations from smaller to larger and from larger to smaller are treated symmetrically and equally.

**Library Profiling.** Important substructures which are present in many compounds can be identified by counting for each scaffold the number of structures that lie on upward paths. A diverse library is one which is free from dominating scaffolds of significant size.

**Library Browsing and Navigation.** Two or more compounds which share a common scaffold have a sibling relation. The present method has potential use in library browsing applications where one may want to link a compound to its sibling compounds. Compounds 19 and 33 in Figure 7 are siblings. A compound can have several scaffold substructures and will then have several sets of siblings. In the example with sequences in Figure 4, structure BC will have sibling AB relative to scaffold B and also a sibling CD relative to scaffold C.

**Identification of Compound Classes.** The method identifies all substructures (scaffolds) which are common to two or more other scaffolds or molecules. The substructures are maximal in the sense that they cannot be expanded in any way and still be substructures of the same set of compounds. A compound class can be defined as a set of compounds sharing the same scaffold. Then, compounds belonging to a class are the ones that lie on paths upward from their class defining scaffold. Conversely, the classes that a compound belongs to are represented by the scaffolds that lie on paths down from the compound.

**Screen for Scaffolds to Be Used for R-Group Analysis.** Fragment-based QSAR can be performed when different substitutions on a scaffold are independent and can be combined. Activity can then be modeled as additive contributions from substituents.<sup>34</sup> The present method detects shared fragments which can be used for QSARs of this type.

**Fragment-Based Screens.** Fragments associated with activity can be added to a set of screening compounds. In a Hasse diagram, compounds on paths up from the fragments will be superstructures of those. The same set of superstructures can be identified via fingerprinting, but the Hasse

diagram method can directly show how fragments and full structures are related to each other.

**Implicit Fragments.** Not all fragments that make up the compounds in the diagram are represented as nodes in the graph. The nodes represent full compounds or substructures which are common to two or more compounds or other substructures. Fragments which are not shared between compounds are not part of the diagram itself. In Figure 8, compound 20 is extended to compounds 5, 19, and 24 by the addition of guanidine, methyl, and methylamine fragments.

**Activity Cliffs.** The concept of activity cliffs was introduced to describe the phenomenon that similar compounds do not always have similar activities.<sup>35</sup> It is important in medicinal chemistry that the situations can be detected and described where a small structure change causes a large activity change. Activity cliffs are easily identified by first detecting siblings in the diagram and then sorting them on activity.

**Relation to Formal Concept Analysis.** The method described here has some resemblance to Formal Concept Analysis (FCA),<sup>36</sup> which is a way to organize collections of objects into lattices based on object properties. FCA uses a partial order induced by subset inclusion, while the present method uses substructure inclusion and does not in general give a proper lattice. The data structures of the two methods have different properties and are created using different algorithms. Lounkine et al.<sup>37</sup> have used FCA in combination with Bemis–Murcko fragmentation<sup>16</sup> to study fragments that appear together in active compounds.

**Time Complexity.** By far the most computation time was spent on the MCS calculation step. Every new structure is compared to all other structures in the diagram with the naïve algorithm that was used, and from this  $O(n^2)$  time complexity is expected if problem size is measured as the number of compounds. This makes the method computationally demanding. Extensive tests on several data sets will be required to produce meaningful experimental timing statistics. This is because MCS calculation time depends strongly on the size and nature of structures, with some pairs of drug-sized molecules taking many orders of magnitude longer than others to process, so influences from this effect may easily overshadow the effect of compound numbers. Calculation of the diagram in Figure 7 took 10 s on a standard PC with a 2.9 GHz processor.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: peter.lind@medivir.se.

### Notes

The authors declare no competing financial interest.

## REFERENCES

- (1) Downs, M. G.; Barnard, J. M. Clustering Methods and Their Uses in Computational Chemistry. In *Reviews in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; John Wiley & Sons, Inc.: Hoboken, NJ, 2002; Vol. 18, Chapter 1, pp 1–40.
- (2) Matter, H. Selecting Optimally Diverse Compounds from Structure Databases: A Validation Study of Two-Dimensional and Three-Dimensional Molecular Descriptors. *J. Med. Chem.* **1997**, *40*, 1219–1229.
- (3) Raymond, J. W.; Blankley, C. J.; Willett, P. Comparison of chemical clustering methods using graph- and fingerprint-based similarity measures. *J. Mol. Graphics Modell.* **2003**, *21*, 421–433.
- (4) Willett, P.; Barnard, J. M.; Downs, G. M. Chemical Similarity Searching. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 983–996.



- (5) Nasr, R.; Hirschberg, D. S.; Baldi, P. Hashing Algorithms and Data Structures for Rapid Searches of Fingerprint Vectors. *J. Chem. Inf. Model.* **2010**, *50*, 1358–1368.
- (6) Haque, I. S.; Pande, V. S.; Walters, W. P. Anatomy of High-Performance 2D Similarity Calculations. *J. Chem. Inf. Model.* **2011**, *51*, 2345–2351.
- (7) Gardiner, J. E.; Gillet, V. J.; Willet, P.; Cosgrove, D. A. Representing Clusters Using a Maximum Common Edge Substructure Algorithm Applied to Reduced Graphs and Molecular Graphs. *J. Chem. Inf. Model.* **2007**, *47*, 354–366.
- (8) Stahl, M.; Mauser, H.; Tsui, M.; Taylor, N. R. A Robust Clustering Method for Chemical Structures. *J. Med. Chem.* **2005**, *48*, 4358–4366.
- (9) Cosgrove, D. A.; Willet, P. SLASH: A program for analyzing the functional groups in molecules. *J. Mol. Graphics Modell.* **1998**, *16*, 19–32.
- (10) Roberts, G.; Myatt, G. J.; Johnson, W. P.; Cross, K. P.; Blower, B. E., Jr. LeadScope: Software for Exploring Large Sets of Screening Data. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 1302–1314.
- (11) Xu, J. A New Approach to Finding Natural Chemical Structure Classes. *J. Med. Chem.* **2002**, *45*, 5311–5320.
- (12) Miller, D. W. A chemical class-based approach to predictive model generation. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 568–578.
- (13) Cross, K. P.; Myatt, G.; Yang, C.; Fligner, M. A.; Verducci, J. S.; Blower, P. E., Jr. Finding Discriminating Structural Features by Reassembling Common Building Blocks. *J. Med. Chem.* **2003**, *46*, 4770–4775.
- (14) Nicolaou, C. A.; Tamura, S. Y.; Kelley, B. P.; Basset, S. I.; Nutt, R. P. Analysis of Large Screening Data Sets via Adaptively Grown Phylogenetic-Like Trees. *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 1069–1079.
- (15) Wilkens, S. J.; Janes, J.; Su, A. I. HierS: Hierarchical scaffold clustering using topological chemical graphs. *J. Med. Chem.* **2005**, *48*, 3182–3193.
- (16) Bemis, G. W.; Murcko, M. A. The Properties of Known Drugs. 1. Molecular Frameworks. *J. Med. Chem.* **1996**, *39*, 2887–2893.
- (17) Koch, M. A.; Schuffenhauer, A.; Scheck, M.; Wetzel, S.; Casaulta, M.; Odermatt, A.; Ertl, P.; Waldmann, H. Charting biologically relevant chemical space: A structural classification of natural products (SCONP). *Proc. Natl. Acad. Sci. U. S. A.* **2005**, *102*, 17272–17277.
- (18) Krier, M.; Bret, G.; Rognan, D. Assessing the Scaffold Diversity of Screening Libraries. *J. Chem. Inf. Model.* **2006**, *46*, 512–524.
- (19) Schuffenhauer, A.; Ertl, P.; Roggo, S.; Wetzel, S.; Koch, M. A.; Waldmann, H. The Scaffold Tree – Visualization of the Scaffold Universe by Hierarchical Scaffold Classification. *J. Chem. Inf. Model.* **2007**, *47*, 47–58.
- (20) Lounkine, E.; Bajorath, J. Core Trees and Consensus Fragment Sequences for Molecular Representation and Similarity Analysis. *J. Chem. Inf. Model.* **2008**, *48*, 1161–1166.
- (21) Clark, A. M. 2D Depiction of Fragment Hierarchies. *J. Chem. Inf. Model.* **2010**, *50*, 37–46.
- (22) Cerruela Garcia, G.; Luque Ruiz, I.; Gómez-Nieto, M. Á. Analysis and Study of Molecule Data Sets Using Snowflake Diagrams of Weighted Maximum Common Subgraph Trees. *J. Chem. Inf. Model.* **2011**, *51*, 1216–1232.
- (23) Varin, T.; Schuffenhauer, A.; Ertl, P.; Renner, S. Mining for Bioactive Scaffolds with Scaffold Networks: Improved Compound Set Enrichment from Primary Screening Data. *J. Chem. Inf. Model.* **2011**, *51*, 1528–1538.
- (24) Schuffenhauer, A.; Brown, N.; Ertl, P.; Jenkins, J. L.; Selzer, P.; Hamon, J. Clustering and Rule-Based Classifications of Chemical Structures Evaluated in the Biological Activity Space. *J. Chem. Inf. Model.* **2007**, *47*, 325–336.
- (25) Lewell, X. Q.; Judd, D. B.; Watson, S. P.; Hann, M. M. RECAP—Retrosynthetic Combinatorial Analysis Procedure: A Powerful New Technique for Identifying Privileged Molecular Fragments with Useful Applications in Combinatorial Chemistry. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 511–522.
- (26) Böcker, A. Toward an Improved Clustering of Large Data Sets Using Maximum Common Substructures and Topological Fingerprints. *J. Chem. Inf. Model.* **2008**, *48*, 2097–2107.
- (27) Hariharan, R.; Janakriaman, A.; Nilakantan, R.; Singh, B.; Varghese, S.; Landrum, G.; Schuffenhauer, A. MultiMCS: A Fast Algorithm for the Maximum Common Substructure Problem on Multiple Molecules. *J. Chem. Inf. Model.* **2011**, *51*, 788–806.
- (28) Garey, M. R.; Johnson, D. S. Appendix: A List of NP-Complete Problems. In *Computers and Intractability: A Guide to the Theory of NP-Completeness. A Series of Books in the Mathematical Sciences*; Klee, V., Ed.; W. H. Freeman and Co.: New York, 1979; p 202.
- (29) Kim, H. S.; Neggers, J. Definitions and Examples. In *Basic Posets*; World Scientific Publishing Co. Pte. Ltd.: Singapore, 1998; pp 1–6.
- (30) GitHub code repository. <https://github.com/peter-lind/hasse-manager> (accessed July 28, 2013).
- (31) GGA Software Services. <http://ggasoftware.com/opensource/indigo> (accessed July 28, 2013).
- (32) Graphviz – Graph Visualization Software. <http://www.graphviz.org/> (accessed July 28, 2013).
- (33) Silverman, B. D.; Platt, D. E. J. Comparative Molecular Moment Analysis (CoMMA): 3D-QSAR without Molecular Superposition. *J. Med. Chem.* **1996**, *39*, 2129–2140.
- (34) Free, S. M.; Wilson, J. W. A Mathematical Contribution to Structure-Activity Studies. *J. Med. Chem.* **1964**, *7*, 395–399.
- (35) Maggiora, G. M. On Outliers and Activity Cliffs - Why QSAR Often Disappoints. *J. Chem. Inf. Model.* **2006**, *46* (4), 1535–1535.
- (36) Wille, R. Formal Concept Analysis as Mathematical Theory of Concepts and Concept Hierarchies. In *Formal Concept Analysis. Foundations and Applications*; Ganter, B.; Stumme, G.; Wille, R., Eds.; Springer: Hedielsberg, Germany, 2005; pp 1–33.
- (37) Lounkine, E.; Auer, J.; Bajorath, J. Formal Concept Analysis for the Identification of Molecular Fragment Combinations Specific for Active and Highly Potent Compounds. *J. Med. Chem.* **2008**, *51* (17), 5342–5348.
- (38) Chemical Entities of Biological Interest. <http://www.ebi.ac.uk/chebi/> (accessed November 7, 2013).