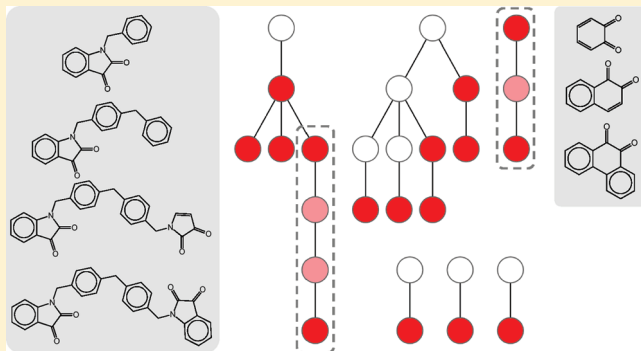


Combining Horizontal and Vertical Substructure Relationships in Scaffold Hierarchies for Activity Prediction

Ye Hu and Jürgen Bajorath*

Department of Life Science Informatics, B-IT, LIMES Program Unit Chemical Biology and Medicinal Chemistry, Rheinische Friedrich-Wilhelms-Universität, Dahlmannstrasse 2, D-53113 Bonn, Germany

ABSTRACT: For a systematic exploration of structural relationships between molecular scaffolds, ~24,000 unique scaffolds were extracted from 458 different target sets. Substructure relationships between these scaffolds were systematically determined. The scaffold tree data structure was utilized to study structural relationships between original scaffolds and derivative scaffolds obtained by rule-based decomposition. Leaf-to-root substructure relationships that resulted from rule-based decomposition were compared to leaf-to-leaf relationships between original scaffolds most of which were not part of the scaffold tree hierarchy. Decomposed scaffolds not contained in active target set compounds were prioritized on the basis of hierarchical scaffold patterns and additional substructure relationships. For high-priority virtual scaffolds, activity predictions were carried out, and these scaffolds were often found in external test compounds having the predicted activity. Taken together, our results suggest that leaf-to-root substructure relationships in scaffold trees should best be complemented with additional substructure relationships to determine high-priority virtual scaffolds for activity prediction.



INTRODUCTION

The analysis of molecular scaffolds or frameworks of bioactive compounds is highly relevant for a number of tasks in medicinal chemistry including, for example, the identification of potentially privileged substructures^{1,2} or of target class-selective chemotypes.³ Furthermore, the search for different scaffolds yielding compounds with similar activity, often referred to as scaffold hopping, is another important goal.^{4,5}

A scaffold is often obtained from an active compound by removing all substituents from ring systems and from linker segments between rings, following the Bemis and Murcko definition.⁶ Alternatively, scaffolds might also be defined on the basis of retrosynthetic criteria⁷ or other chemical rules. Large-scale scaffold analyses have been carried out, for example, to assess the structural diversity of synthetic compounds⁸ and screening libraries,⁹ survey heteroaromatic scaffolds in bioactive compounds,^{10,11} or study the distribution of scaffolds in compounds at different pharmaceutical development stages.¹²

Since conventional scaffolds⁶ contain all rings and linkers between rings, the addition of any ring to a compound (e.g., a phenyl substituent) always constitutes a new scaffold, given the underlying hierarchical scaffold definition (although the compound might in such cases better be considered an analog). This is often considered a potential caveat in scaffold analysis.¹³ Accordingly, scaffold classification schemes have been introduced that do not predominantly focus on core structures, but chemical transformations,¹³ similar to the matched molecular pair concept,¹⁴ or that organize ring systems after removal of linkers.¹⁵

Another scaffold organization scheme was introduced that iteratively removes rings from initially derived Bemis and Murcko-like scaffolds, starting at peripheral and moving to more central positions until only a single ring remains.¹⁶ Here rings are not only removed that are connected by linkers but also from condensed ring systems by dividing them into individual (parental) rings. A set of generally applicable chemical rules is applied to prioritize rings for iterative removal. For scaffolds from any source, these procedures generate a hierarchy where initially derived scaffolds ("leaves") are systematically reduced until an individual "root" ring remains. For sets of active compounds, the resulting pathways of this "leaf-to-root" hierarchy are displayed as so-called Scaffold Trees¹⁶ (STs) that currently probably represent the most general data structure to hierarchically organize scaffold populations. ST "leaf" scaffolds differ from Bemis and Murcko scaffolds only in that double bonded atoms (e.g., carbonyl oxygens) attached to rings or linkers are retained as part of the scaffold. Given the rule-based decomposition of ring systems, the ST hierarchy typically contains scaffolds that are not contained in the original set of active compounds, so-called virtual scaffolds.¹⁶ Thus, STs can be utilized to predict biological activities of such scaffolds.^{16,17} For activity prediction, STs of different compound sets can also be merged by mapping shared scaffolds and combining the pathways they are involved in.¹⁸

Received: November 15, 2010

Published: January 27, 2011

We have been interested in scaffold hierarchies to systematically analyze substructure relationships between scaffolds and their relevance for biological activity. This analysis was inspired by a previous finding that 71% of bioactive scaffolds were involved in defined substructure relationships (i.e., A is a substructure of B).¹⁹ We reasoned that it might be possible to reconcile and further explore these structural relationships on the basis of scaffold hierarchies. In this context, STs have been of particular interest because they capture substructure relationships along decomposition pathways (i.e., from leafs to roots), which we, for the purpose of our analysis, term “vertical” relationships. However, the substructure relationships we identified previously are, in the context of the ST hierarchy, leaf-to-leaf relationships, which we term “horizontal”. Such relationships have thus far not been explicitly considered in ST analysis. Therefore, we have systematically analyzed to what extent vertical and horizontal substructure relationships between scaffolds complement each other. For this purpose, a large-scale analysis of target set-dependent scaffold hierarchies has been carried out. Prioritized candidate scaffolds that were not contained in target set compounds have been mapped to external compound sources and their biological activity has been predicted.

MATERIALS AND METHODS

For scaffold generation, bioactive compounds were extracted from the ChEMBL²⁰ database (CDB) and BindingDB²¹ (BDB). These databases are two major publicly available repositories of active compounds from medicinal chemistry sources with defined target and activity annotations. ST scaffolds were generated using the Scaffold Tree Generator program.¹⁶ Figure 1 illustrates the rule-based decomposition of ST scaffolds and the formation of a tree branch. Resulting STs were drawn with Cytoscape.²² Hierarchically organized scaffolds were prioritized on the basis of defined scaffold patterns and substructure relationships and mapped to scaffolds of active compounds from the Molecular Drug Data Report²³ (MDDR) and approved drugs from DrugBank.²⁴ The scaffold analysis reported herein was carried out with in-house generated Molecular Operating Environment²⁵ (MOE) Scientific Vector Language (SVL), Perl, and Pipeline Pilot²⁶ scripts.

Upon publication the scaffold hierarchies generated for our analysis can be freely obtained via the following URL: <http://www.lifescienceinformatics.uni-bonn.de> (please, see the “Downloads” section).

RESULTS AND DISCUSSION

Compound and Scaffold Statistics. From the pool of CDB and BDB compounds, we extracted compound activity classes (with specific target annotations) under the condition that each activity class (target set) had to contain at least 10 compounds with at least 1 μ M potency. On the basis of these criteria, we obtained 458 target sets containing a total of 34,916 active compounds that yielded 23,879 unique ST scaffolds.

Scaffold Hierarchies. For each of our 458 target sets, an ST was generated. Figure 2 shows a representative example and illustrates how different leaf scaffolds form (or do not form) converging scaffold pathways toward a root scaffold, i.e. an individual ring. A consequence of this hierarchical scaffold decomposition scheme is that not all ST scaffolds are represented

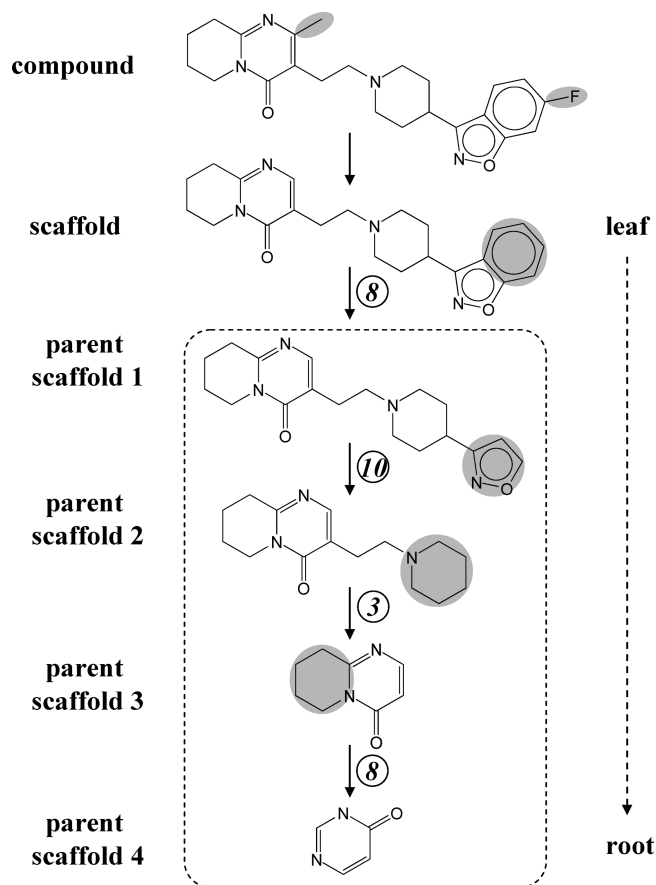


Figure 1. Scaffold hierarchy. Shown is an exemplary scaffold branch¹⁶ generated from an active compound. The leaf scaffold is extracted from the compound by removing all single-bond substituents from rings or linkers between rings. During each decomposition step, a smaller (parent) scaffold is generated. One ring is iteratively removed per step from the current scaffold according to 13 predefined chemical rules¹⁶ until only a single ring remains as the root scaffold. In this example, parent scaffold 1 was generated by removing a benzene ring on the basis of rule 8 (i.e., “remove rings with the least number of heteroatom first”). In the following steps, rule 3 (“choose a parent scaffold having the smallest number of acyclic linker bonds”) and rule 10 (“smaller rings are removed first”) were applied and, finally, rule 8 again to yield the root scaffold.

by active compounds from which leaf scaffolds are derived. This leads to the distinction of “real” ST scaffolds (R) that are contained in active compounds and “virtual” scaffolds (V) that do not occur in source compounds, as illustrated in Figure 2. Following this classification scheme, the 23,879 unique scaffolds comprising 458 target set STs yielded 13,377 real scaffolds and 10,502 virtual scaffolds.

Substructure Relationships. Scaffold trees capture hierarchical leaf-to-root substructure relationships between scaffolds along decomposition pathways but do not explicitly account for horizontal leaf-to-leaf substructure relationships. A central point of our study has been to determine to what extent such horizontal substructure relationships are implicitly captured by the ST data structure. Therefore, we first identified all pairs of leaf scaffolds that represented a defined substructure relationship. For this analysis, the most generic scaffold, the benzene ring, was not considered. As reported in Table 1, we detected 13,181 pairs that involved a total of 9712 leaf scaffolds, i.e. 73% of all original

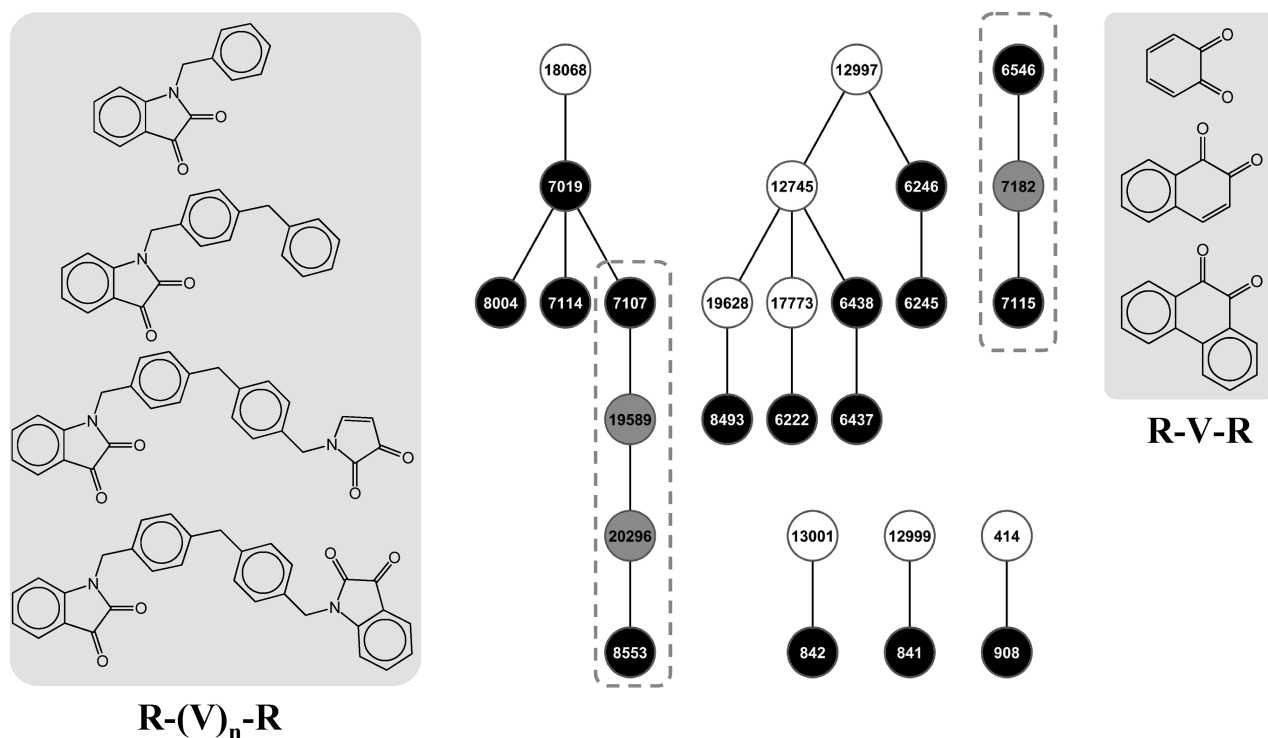


Figure 2. Real and virtual scaffolds. An exemplary scaffold tree is shown for the carboxylesterase-2 inhibitor set. Nodes represent scaffolds that are labeled with IDs and gray-scaled according to different scaffold types including “real” scaffolds (black), “virtual” scaffolds (white), and “prioritized virtual” scaffolds (gray). Edges connect scaffolds in a tree branch (decomposition pathway). In this orientation, leaf scaffolds are at the bottom and root scaffolds at the top. Two scaffold branches are highlighted using dashed rectangles. In the left branch, virtual scaffolds 19589 and 20296 are located between two “real” scaffolds 7107 and 8553, forming an R-(V)₂-R pattern. In the right branch, virtual scaffold 7182 has two neighboring real scaffolds (6546 and 7115), forming an R-V-R pattern. Scaffold sequences comprising these patterns are displayed on a light gray background. Virtual scaffolds from such patterns are considered prioritized virtual scaffolds for activity prediction.

Table 1. Substructure Relationships^a

substructure relationship	scaffold pairs	scaffolds
horizontal	13,181	9712 (73%)
vertical	4217	5205 (39%)

^a Each scaffold pair represents a substructure relationship. Horizontal relationships represent leaf-to-leaf and vertical leaf-to-root substructure relationships between scaffolds. Vertical relationships are determined by the ST hierarchy.

scaffolds. Thus, the majority of all leaf ST scaffolds were involved in pairwise substructure relationships (similar to 71% of all Bemis and Murcko scaffolds extracted from CDB compounds¹⁹). We then determined how many of these scaffold pairs also represented vertical ST substructure relationships. Only 4217 of all 13,181 pairs (32%) were detected in ST pathways. These pairs involved a total of 5205 scaffolds, i.e. 39% of all leaf scaffolds (Table 1). Thus, for all 458 target sets, the STs only contained about one-third of the substructure relationships that were present between the original scaffolds. On the basis of these findings, we then asked the question how the additional substructure relationship information might be utilized for tree analysis.

Substructure Information Content. Substructure relationships can be added to the ST structure by annotating trees with nonpathway substructure pair information, as illustrated in Figure 3. The hierarchy of the exemplary ST on the left in Figure 3 contains three pairwise substructure relationships, and

scaffolds 3 and 4 are each involved in two pairs. Within a branch, a scaffold can be a part of at most two pairs and hence these scaffolds cannot be further distinguished by pair numbers. However, on the right in Figure 3, the tree is annotated with all additional substructure relationships involving leaf scaffolds (two in this case). Now scaffold 2 is also involved in two pairs, and leaf scaffold 1 and root scaffold 4 are each involved in three pairs. Thus, taking this additional information into account, scaffolds can be further differentiated by the number of substructure pairs they participate in and scaffolds involved in most pairs can be prioritized on the basis of substructure information content. We next evaluated how added substructure information might affect activity predictions. For this purpose, all possible pairs between virtual and real scaffolds were systematically analyzed.

Pairs of Virtual and Real Scaffolds. One of the most interesting aspects of the ST data structure is the opportunity to predict the activity of virtual scaffolds.^{16,17} Prime candidates for activity prediction are virtual scaffolds that are proximal to real scaffolds in the tree because of their structural relatedness,¹⁷ which represents a rather intuitive approach, leading to a number of successful predictions.^{17,18}

In order to systematically explore relevant scaffold pairings, we isolated all V-R scaffold pairs (i.e., pairs formed by a virtual and a real scaffold) from all target set STs. As reported in Table 2, 53,220 V-R pairs were found in 442 target sets. When we limited the magnitude of structural differences within a pair to a maximum of two rings, the number of V-R pairs was reduced to

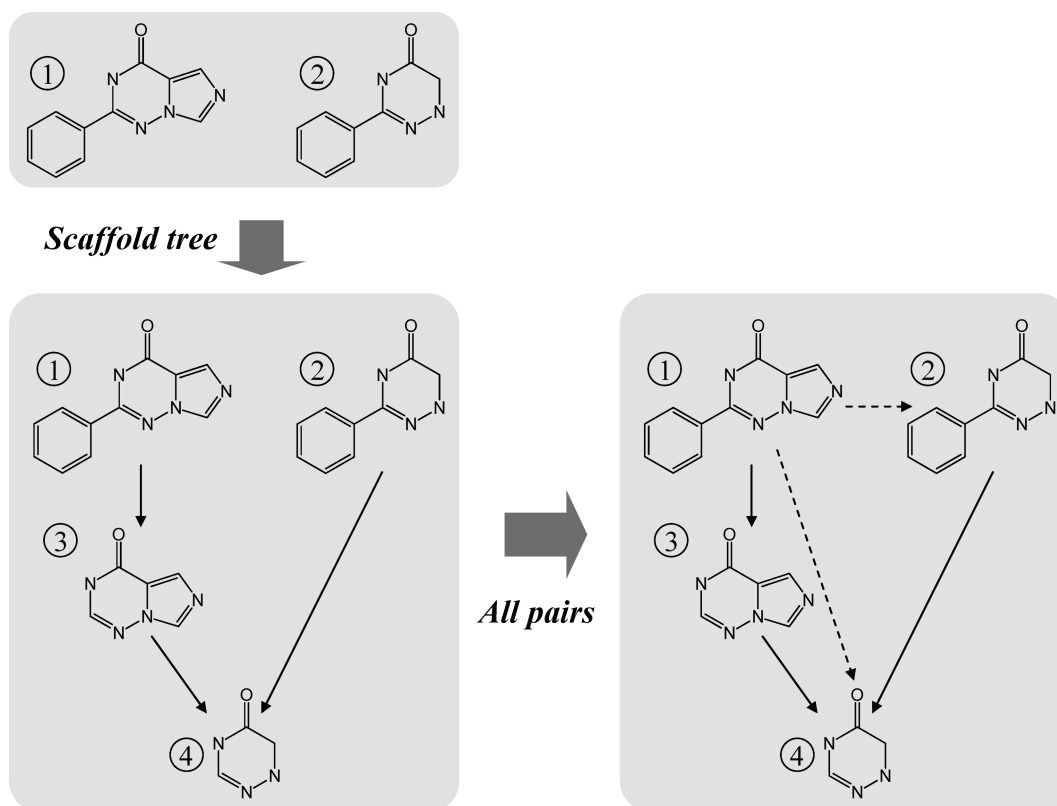


Figure 3. Scaffold pairs representing substructure relationships. For two scaffolds 1 and 2, the parent scaffolds 3 and 4 are generated by iteratively removing one ring from each child. The resulting scaffold tree contains three scaffold pairs that form substructure relationships, i.e. 1–3, 3–4, and 2–4 (solid arrows). However, by examining additional substructure relationships for leaf scaffolds, two additional pairs, i.e. 1–2 and 1–4 (dashed arrows), are identified. Thus, scaffolds 1 and 4 are now each involved in three substructure relationships, and this additional information can be used to further prioritize scaffolds.

Table 2. V-R Scaffold Pairs^a

V-R pair category	V-R pairs	virtual scaffolds	real scaffolds	target sets
V - R = n rings	53,220	9668	12,541	442
V - R ≤ 2 rings	25,664	8750	11,799	442
V - R = 1 ring	10,886	6584	8933	440

^aReported are the total number of V-R scaffold pairs within ST hierarchies ($|V - R| = n$ rings, i.e. no structural constraint), the number of V-R pairs where the virtual and real scaffold differed by at most two ring systems, and the number of pairs where V and R differed by only one ring. In addition, for each V-R pair category, the number of virtual and real scaffolds and the number of target sets from which the pairs originated are provided.

approximately half. Moreover, when structural deviations in pairs were limited to one ring, 10,886 V-R pairs were obtained from 440 target sets that involved 6584 virtual and 8933 real scaffolds (Table 2). These V-R pairs provided a pool of scaffold pairs for pattern definition and activity prediction, as discussed in the following.

Prioritized Scaffold Patterns. We next defined scaffold patterns formed by V-R pairs that were of increasing attractiveness for activity prediction. Virtual scaffolds were considered attractive candidates if they were “framed” by real scaffolds. In Figure 2, such $R-(V)_n-R$ patterns are highlighted. For example, two virtual scaffolds appear in a pathway between two real scaffolds (i.e., $n = 2$). However, the most attractive pattern is formed when a virtual scaffold has two real scaffolds as neighbors

(i.e., $n = 1$), yielding an R-V-R pattern. In this case, a virtual scaffold is involved in substructure relationships to a known active child and parent scaffold (and each of these real scaffolds differs in one ring from the virtual scaffold). For the purpose of our analysis, we regarded virtual scaffolds involved in $R-(V)_n-R$ and R-V-R patterns as “prioritized virtual” scaffolds.

As illustrated in Figure 3, nonterminal scaffolds in tree branches are always involved in at least two substructure relationships in the ST hierarchy, and these structural relationships correspond to two scaffold pairs (for example, pairs 1–3 and 3–4 in Figure 3). However, scaffolds might also participate in additional substructure relationships/pairs with leaf (or other) scaffolds that are not a part of the hierarchy. These pairwise relationships can be systematically detected and added to the tree structure, as also illustrated in Figure 3. Hence, by further extending the ST hierarchy scaffolds can also be evaluated by taking additional substructure information into account. For prioritized virtual scaffolds in R-V-R patterns, nonhierarchy substructure pairings (i.e., relationships to leaf scaffolds and other real scaffolds) add further activity-relevant structural information. Moreover, for many virtual scaffolds that are not part of R-V-R patterns, other V-R pairs might also be found that provide additional substructure information. Thus, the likelihood of activity would be expected to further increase for virtual scaffolds involved in additional substructure relationships to known active scaffolds.

Scaffold Mapping. Different categories of virtual scaffolds were mapped to scaffolds extracted from MDDR compounds and approved drugs. We reasoned that prioritized virtual scaffolds

Table 3. Scaffold Mapping^a

scaffold type	no. of scaffolds		
	total	MDDR	drugs
real	13,377	2658 (20%)	226 (1.7%)
virtual	10,502	1005 (10%)	59 (0.6%)
prioritized virtual R - (V) _n - R	997	174 (17%)	30 (3%)
R - V - R	544	120 (22%)	23 (4.2%)
R - V - R or ≥ 2 V-R pairs	1678	449 (27%)	75 (4.5%)

^a Five sets of CDB/BDB scaffolds extracted from ST hierarchies including real and virtual scaffolds and prioritized virtual scaffolds in different patterns were mapped to scaffolds of MDDR compounds and of approved drugs from DrugBank. The pattern designated “R-V-R or ≥ 2 V-R pairs” combines all prioritized virtual scaffolds of R-V-R patterns with other virtual scaffolds that were involved in at least two additional substructure relationships (i.e. nonhierarchy pairs). The total number of scaffolds comprising each set is given, and the number of these scaffolds that matched MDDR or drug scaffolds is reported.

should display an increasing tendency to match scaffolds from bioactive compounds. From 157,522 MDDR compounds and 1247 drugs, 71,649 and 722 scaffolds were obtained, respectively. Mapping of CDB/BDB scaffolds from ST hierarchies to the MDDR was considered a meaningful exercise because, as reported in Table 3, only 20% of all real and 10% of all virtual scaffolds were found to match MDDR scaffolds. For approved drugs, the numbers of matching scaffolds were much smaller.

We first considered virtual scaffolds from R-(V)_n-R patterns, which reduced the number of candidate scaffolds from 10,502 to 997 (~9%). However, the match rate of these scaffolds increased to 17% in the MDDR and 3% in drugs (Table 3). We then focused on virtual scaffolds from R-V-R patterns, which further reduced the number of candidate scaffolds to 544. In this case, 22% of these scaffolds matched MDDR scaffolds. Thus, compared to nonprioritized virtual scaffolds, the use of R-V-R virtual scaffolds essentially doubled the match rate. Furthermore, we also found that nearly all virtual scaffolds from R-V-R patterns were also involved in additional substructure relationships. Therefore, we complemented the set of R-V-R scaffolds with other virtual scaffolds that were involved in at least two V-R substructure pairs outside the ST hierarchy, which resulted in a total of 1,678 prioritized virtual scaffolds. This extended scaffold set produced a match rate of 27% in the MDDR (and 4.5% in approved drugs). Thus, scaffold mapping revealed that prioritized virtual scaffolds were generally more likely to match bioactive scaffolds than nonprioritized virtual scaffolds.

Activity Prediction. On the basis of these findings, we went a step further and predicted the activity of high-priority virtual scaffolds. For this purpose, we ranked the two sets of matching virtual scaffolds from R-V-R patterns on the basis of additional substructure information content, i.e. additional V-R pairs these scaffolds were involved in. Tables 4 and 5 report the rankings for prioritized virtual scaffolds matching MDDR compounds (120 scaffolds) and approved drugs (23 scaffolds), respectively. For these scaffolds, up to 58 additional substructure relationships (V-R pairs) were detected. Each prioritized V-scaffold was then predicted to have the same activity as the neighboring R-scaffolds, and this prediction was compared to the target annotations of matching MDDR compounds or drugs. Correct predictions were found at a high rate for 26 of 120 virtual scaffolds in the MDDR and for six of 23 virtual scaffolds in DrugBank, although these compound sources have different target distributions.

Table 4. Activity Prediction for Prioritized Virtual Scaffolds in the MDDR^a

ScaffID	#additional V-R pairs	correct prediction	SMILES
11893	58	<i>cathepsin B</i>	<i>c1ccc(cc1)c2ccccc2</i>
12745	53	<i>serotonin receptor 2c, renin</i>	<i>c1ccc2ccccc2(c1)</i>
10815	38		<i>c1ccc(cc1)Cc2ccccc2</i>
12771	22	<i>matrix metalloproteinases 1, 3</i>	<i>c1ccc2ncccc2(c1)</i>
10620	17	<i>matrix metalloproteinases 2, 3, 9, 13</i>	<i>c1ccc(cc1)COc2ccccc2</i>
12537	17		<i>c1ccc2CCCc2(c1)</i>
12634	17		<i>c1ccc2[nH]cnc2(c1)</i>
12772	16		<i>c1ccc2ncccc2(c1)</i>
22794	16		<i>c1ccc2occcc2(c1)</i>
11896	13		<i>c1ccc(cc1)c2ccccc2</i>
21195	13		<i>c1ccc(cc1)C2CCNCCC2</i>
15362	12		<i>O=C(CC(c1ccccc1)-c2ccccc2)N3CCCC3</i>
21850	12	<i>matrix metalloproteinases 3, 8</i>	<i>c1ccc(cc1)OCc2cnc3ccccc23</i>
11700	11		<i>c1ccc(cc1)Oc2ccccc2</i>
9220	9		<i>O=S(=O)(Nc1ccccc1)-c2ccccc2</i>
9823	9		<i>c1ccc(cc1)C2CCCCC2</i>
12828	9	<i>serotonin receptor 1d, 1b</i>	<i>c1ccc3c(c1)[nH]cc3-(C2CC[N+](C2))</i>
13145	9	<i>adenosine receptor A2A</i>	<i>c1nc2nc[nH]c2(n1)</i>
22561	9		<i>c1ccc2[n+](c1)cccc2(c1)</i>
4470	8		<i>O=C(NCCc1ccccc1)c2ccccc2</i>
6312	8		<i>O=C(c1ccccc1)c2ccccc2</i>
9366	8		<i>O=S(=O)(c1ccccc1)N2CCCC2</i>
9812	8		<i>c1ccc(cc1)C2CC2</i>
22733	8		<i>c1ccc2nc(ccc2(c1))N3CC-[N+](CC3)C1COC(C1)n3cnc2cncnc23</i>
600	7		<i>O=C(CCc1ccccc1)-NCCc2ccccc2</i>
2054	7		<i>c1ccc(cc1)CCc2ccccc2</i>
10321	7		<i>c1[nH]cc(n1)C2CC2</i>
9597	6		<i>c1ccc(cc1)Cc2ccc3ccccc3(c2)</i>
10813	6	<i>aldose reductase</i>	<i>c1ccc(cc1)Nc4ncnc3c4-(ncn3(C2CCCCO2))</i>
11486	6	<i>adenosine receptor A1</i>	<i>c1ccc(cc1)c2cc3ccccc3-([nH]2)</i>
11769	6	<i>serotonin receptor 2a</i>	<i>c1ccc(cc1)c2cc3ccccc3-([nH]2)</i>
11909	6		<i>c1ccc(cc1)c2cncnc2</i>

^a The 120 virtual scaffolds from R-V-R patterns that matched MDDR scaffolds were ranked according to the number of additional non-hierarchy V-R substructure pairs they were involved in and their activity was predicted. Thirty-two scaffolds were involved in more than five additional pairs and are listed. Prioritized virtual scaffolds with correct activity prediction are shown in italics and their activities are reported. In addition, SMILES²⁷ representations of ranked scaffolds are provided.

As shown in Tables 4 and 5, correct predictions were preferentially observed for highly ranked virtual scaffolds having high substructure information content.

In Table 4, the top-ranked prioritized virtual scaffold is the biphenyl scaffold, which occurred in the cathepsin B scaffold tree and was hence predicted to be present in compounds active against cathepsin B. Figure 4 shows the scaffold tree environment of the biphenyl scaffold and its two immediate real scaffold neighbors, representing an R-V-R pattern. Also shown is a cathepsin B inhibitor that was found to contain this prioritized scaffold. The second-ranked scaffold in Table 4 is naphthalene, which was

Table 5. Activity Prediction for Prioritized Virtual Scaffolds in DrugBank^a

ScaffID	#additional V-R pairs	correct prediction	SMILES
11893	58		<chem>c1ccc(cc1)c2ccccc2</chem>
12745	53	<i>beta-2 adrenergic receptor, cyclooxygenase 2</i>	<chem>c1ccc2ccccc2(c1)</chem>
10815	38	<i>dopamine transporter</i>	<chem>c1ccc(cc1)Cc2ccccc2</chem>
12771	22		<chem>c1ccc2ncccc2(c1)</chem>
10620	17		<chem>c1ccc(cc1)COc2ccccc2</chem>
12537	17		<chem>c1ccc2CCc2(c1)</chem>
11700	11		<chem>c1ccc(cc1)Oc2ccccc2</chem>
12828	9	<i>serotonin receptor 1d, 1b, 2a</i>	<chem>c1ccc3c(c1)[nH]cc3(C2CC[N+](CC2))</chem>
13145	9	<i>adenosine receptor A3</i>	<chem>c1ncc2nc[nH]c2(n1)</chem>
22561	9		<chem>c1ccc2[n+](c1)cccc2(c1)</chem>
4470	8		<chem>O=C(NCCc1ccccc1)c2ccccc2</chem>
6312	8		<chem>O=C(c1ccccc1)c2ccccc2</chem>
9812	8		<chem>c1ccc(cc1)C2CC2</chem>
600	7	<i>purine nucleoside phosphorylase (PNP)</i>	<chem>C1COC(C1)n3cnc2enonc23</chem>
10321	7		<chem>c1ccc(cc1)CCc2ccccc2</chem>
11909	6		<chem>c1ccc(cc1)c2ccnnc2</chem>
154	5		<chem>C1CC2CCC(C1)[N+](C2)</chem>
11248	3	<i>alpha-2a adrenergic receptor</i>	<chem>c1ccc(cc1)Nc2ccccc2</chem>
17561	3		<chem>O=C(Nc1nccs1)c2ccccc2</chem>
10075	2		<chem>c1ccc(cc1)CC2CCCC2</chem>
21088	2		<chem>c1ccc(cc1)C(CCC[N+](C)CCCC2)c3ccccc3</chem>
21181	2		<chem>c1ccc(cc1)C2CCCC2</chem>
9415	0		<chem>O=S(=O)(c1ccccc1)c2ccccc2</chem>

^a The 23 virtual scaffolds from R-V-R patterns that matched approved drugs were ranked according to the number of additional nonhierarchy V-R substructure pairs they were involved in and their activity was predicted. Prioritized virtual scaffolds with correct activity prediction are shown in italics and their activities are reported. In addition, SMILES²⁷ representations of ranked scaffolds are provided.

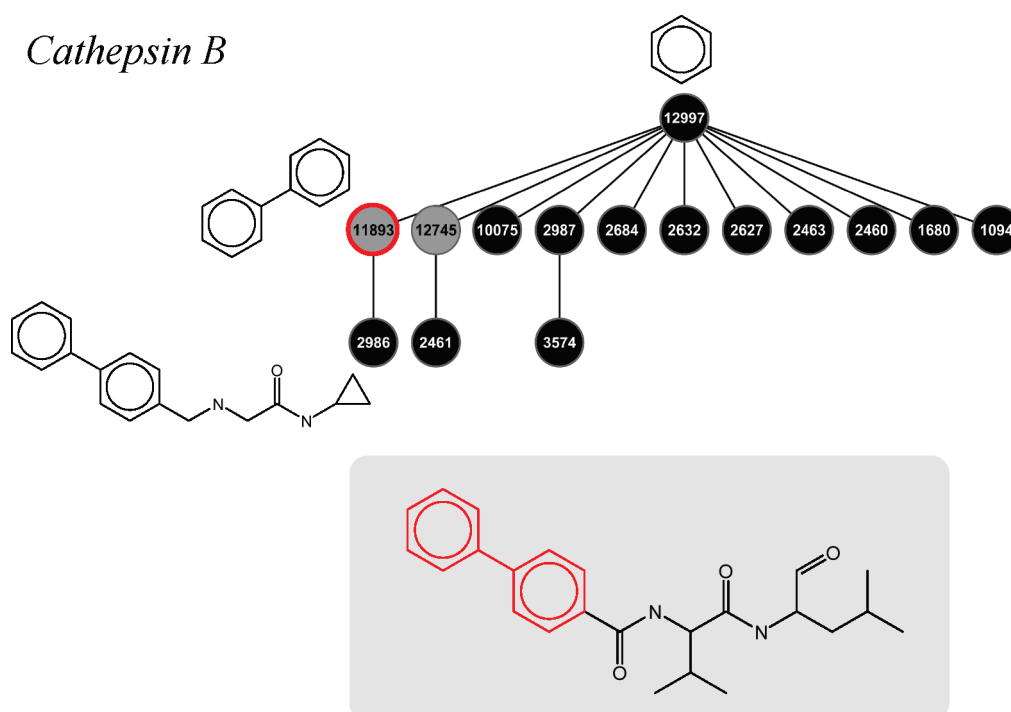
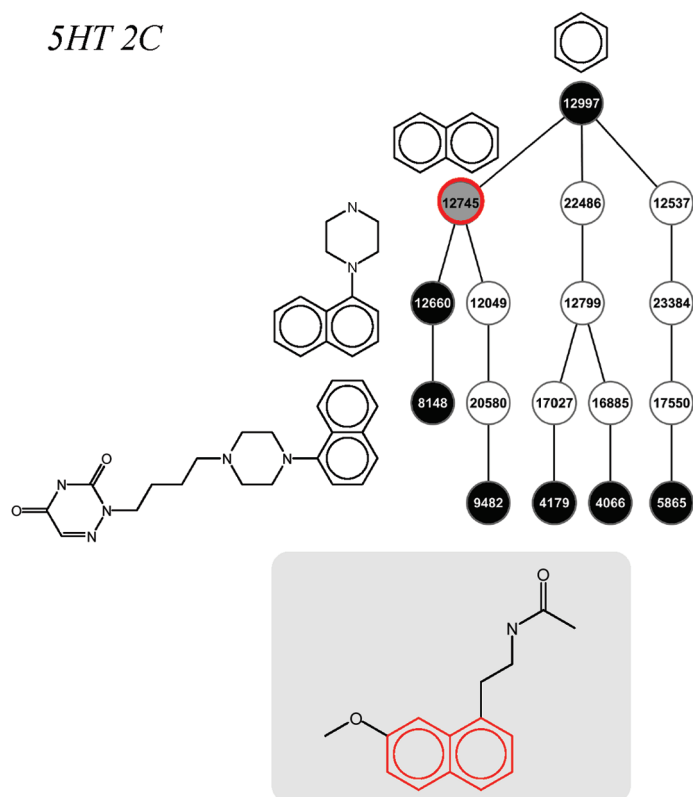


Figure 4. Activity prediction for the biphenyl scaffold (bioactive compounds). Scaffold tree branches for the cathepsin B inhibitor set contained the prioritized “virtual” biphenyl scaffold (labeled with ID 11893 and identified by a red circle). Neighboring “real” scaffolds of the same branch are shown. A representative cathepsin B inhibitor containing the biphenyl scaffold (red) is shown on a light gray background.

a

5HT 2C



b

Renin

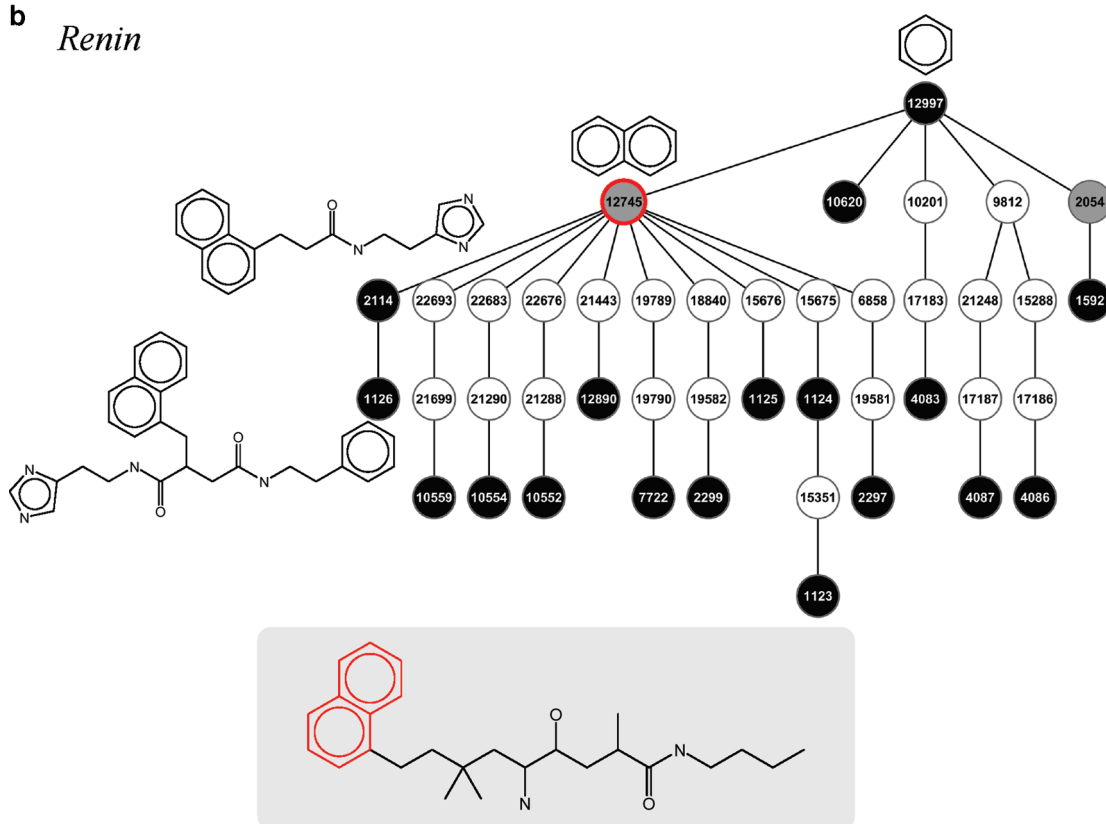


Figure 5. Activity prediction for the naphthalene scaffold (bioactive compounds). Scaffold tree branches for (a) 5HT 2C antagonists and (b) renin inhibitors contained naphthalene as a prioritized virtual scaffold. The presentation is according to Figure 4. Matching bioactive compounds containing this scaffold (red) are shown.

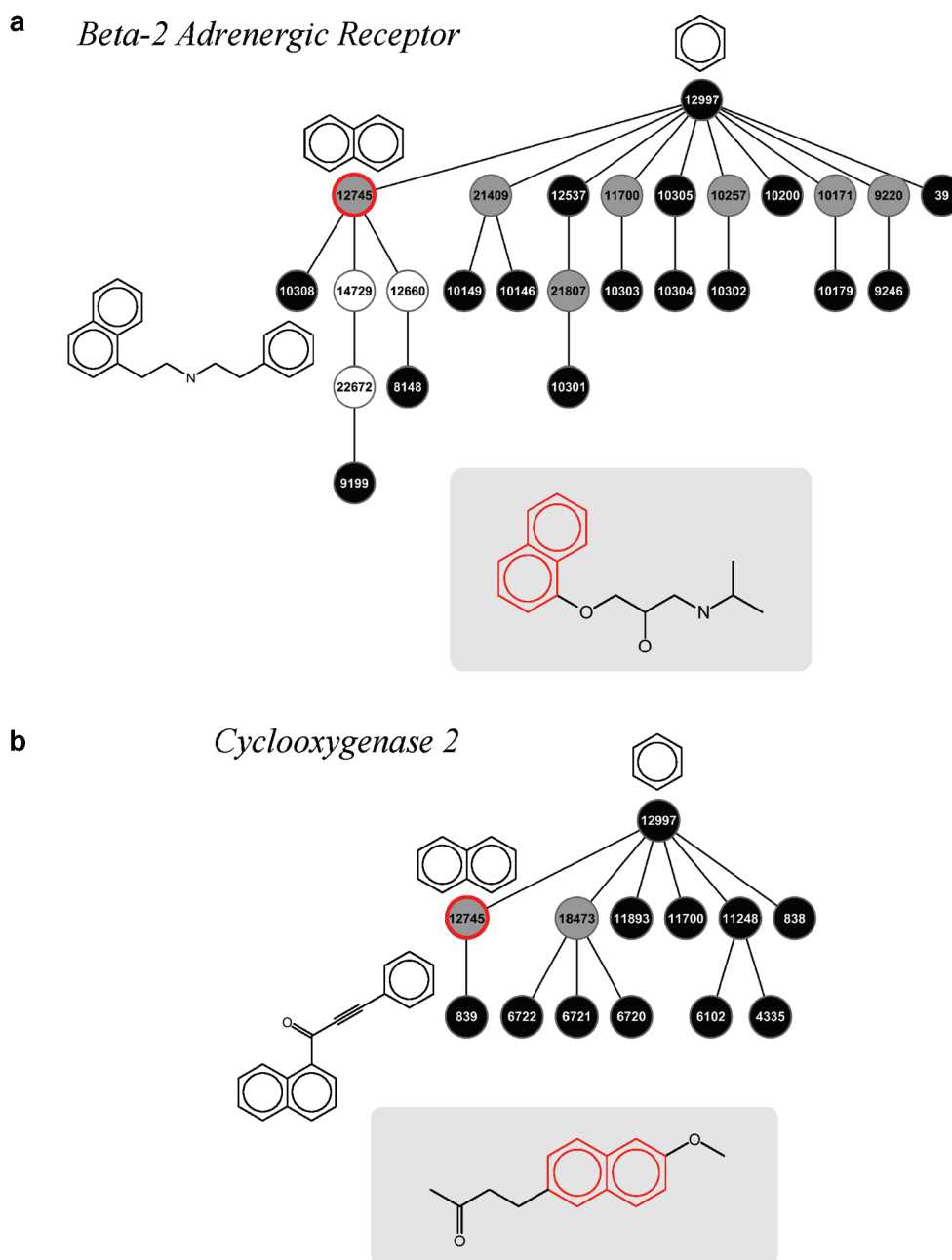


Figure 6. Activity prediction for the naphthalene scaffold (drugs). Scaffold tree branches for (a) beta-2 adrenergic receptor antagonists and (b) cyclooxygenase 2 inhibitors contained naphthalene as a prioritized virtual scaffold. The presentation is according to Figure 4. Matching drugs are shown.

prioritized in two different scaffold trees originating from 5HT 2C serotonin receptor antagonists and renin inhibitors, respectively, both of which correctly matched MDDR compounds. In Figure 5, the different scaffold tree environments of the naphthalene scaffold are shown. In the 5HT 2C scaffold tree in Figure 5a, this high-priority scaffold occurs in a peripheral branch, whereas in the renin inhibitor-derived tree in Figure 5b it is the most central scaffold from which 10 sub-branches originate. Hence, these target-set derived scaffold tree environments of naphthalene differed substantially. However, for both targets, active compounds containing the naphthalene moiety were identified. As reported in Table 5, the biphenyl scaffold did not yield a correct match, i.e. no drug was found directed against a target representing one of the trees where the biphenyl

scaffold was prioritized. However, for the naphthalene scaffold, drugs acting against two of its targets were identified, the beta-2 adrenergic receptor and cyclooxygenase 2. The corresponding scaffold tree environments and matching drugs are shown in Figure 6a and Figure 6b, respectively. Thus, these targets differed from those for which matching MDDR compounds were identified. In Table 5, the third-ranked scaffold is diphenylmethane, which is closely related to the biphenyl scaffold. In this case, no matching MDDR compound was identified. However, the diphenylmethane scaffold was found to correctly match a drug active against the dopamine transporter. In Figure 7, the corresponding scaffold tree environment of diphenylmethane is shown. Here, this prioritized virtual scaffold is also the most central scaffold and involved in seven partly

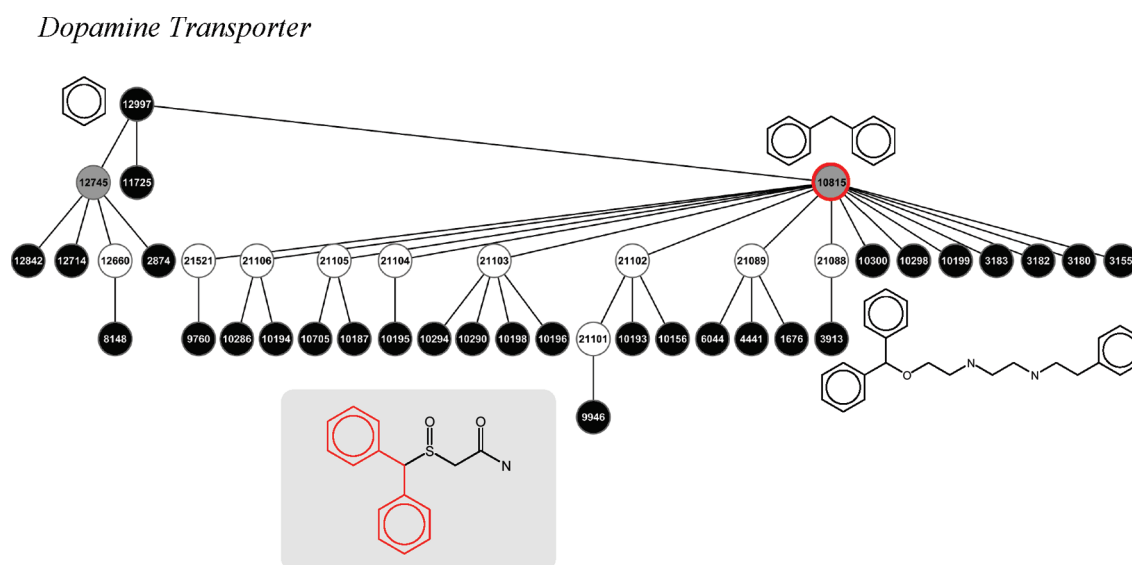


Figure 7. Activity prediction for the diphenylmethane scaffold (drugs). Scaffold tree branches for dopamine transporter inhibitors contained diphenylmethane as a prioritized virtual scaffold. The presentation is according to Figure 4. A matching drug is shown.

overlapping R-V-R patterns, making it a prime candidate for activity prediction.

CONCLUSIONS

In this study, we have systematically explored the overlap between horizontal and vertical substructure relationships in scaffold hierarchies of many different target sets. Only about a third of all leaf-to-leaf substructure relationships detected in our large-scale analysis were found to be implicitly covered by scaffold hierarchies. Hierarchical and nonhierarchical substructure relationships are complementary in nature. Thus, the additional substructure information was included in the analysis to further differentiate between scaffolds. On the basis of our findings, virtual scaffolds were successfully prioritized for scaffold mapping and activity prediction by combining scaffold pattern and substructure pair information. Given the wealth of available scaffold substructure relationships, scaffold prioritization scheme introduced herein should also be useful for practical applications of scaffold hierarchies to predict novel active compounds.

AUTHOR INFORMATION

Corresponding Author

*Phone: +49-228-2699-306. Fax: +49-228-2699-341. E-mail: bajorath@bit.uni-bonn.de.

REFERENCES

- (1) Evans, B. E.; Rittle, K. E.; Bock, M. G.; Dipardo, R. M.; Freidinger, R. M.; Whitter, W. L.; Lundell, G. F.; Veber, D. F.; Anderson, P. S. Methods for Drug Discovery: Development of Potent, Selective, Orally Effective Cholecystokinin Antagonists. *J. Med. Chem.* **1988**, *31*, 2235–2246.
- (2) Schnur, D. M.; Hermsmeier, M. A.; Tebben, A. J. Are Target-Family-Privileged Substructures Truly Privileged? *J. Med. Chem.* **2006**, *39*, 2000–2009.
- (3) Hu, Y.; Wassermann, A. M.; Lounkine, E.; Bajorath, J. Systematic Analysis of Public Domain Compound Potency Data Identifies Selective Molecular Scaffolds across Druggable Target Families. *J. Med. Chem.* **2010**, *53*, 752–758.
- (4) Brown, N.; Jacoby, E. On Scaffolds and Hopping in Medicinal Chemistry. *Mini Rev. Med. Chem.* **2006**, *6*, 1217–1229.
- (5) Zhao, H. Scaffold Selection and Scaffold Hopping in Lead Generation: A Medicinal Chemistry Perspective. *Drug Discovery Today* **2007**, *12*, 149–155.
- (6) Bemis, G. W.; Murcko, M. A. The Properties of Known Drugs. 1. Molecular Frameworks. *J. Med. Chem.* **1996**, *39*, 2887–2893.
- (7) 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.
- (8) Lipkus, A. H.; Yuan, Q.; Lucas, K. A.; Funk, S. A.; Bartelt, W. F., III; Schenck, R. J.; Tripp, A. J. Structural Diversity of Organic Chemistry. A Scaffold Analysis of the CAS Registry. *J. Org. Chem.* **2008**, *73*, 4443–4451.
- (9) Krier, M.; Bret, G.; Rognan, D. Assessing the Scaffold Diversity of Screening Libraries. *J. Chem. Inf. Model.* **2006**, *46*, 512–524.
- (10) Ertl, P.; Jelfs, S.; Mühlbacher, J.; Schuffenhauer, A.; Selzer, P. Quest for the Rings. In Silico Exploration of Ring Universe To Identify Novel Bioactive Heteroaromatic Scaffolds. *J. Med. Chem.* **2006**, *49*, 4568–4573.
- (11) Pitt, W. R.; Parry, D. M.; Perry, B. G.; Groom, C. R. Heteroaromatic Rings of the Future. *J. Med. Chem.* **2009**, *52*, 2952–2963.
- (12) Hu, Y.; Bajorath, J. Scaffold Distributions in Bioactive Molecules, Clinical Trials Compounds, and Drugs. *ChemMedChem* **2010**, *5*, 187–190.
- (13) Katritzky, A. R.; Kiely, J. S.; Hebert, N.; Chassaing, C. Definition of Templates within Combinatorial Libraries. *J. Comb. Chem.* **2000**, *2*, 2–5.
- (14) Wilkens, S. J.; Janes, J.; Su, A. I. HierS: Hierarchical Scaffold Clustering Using Topological Chemical Graphs. *J. Med. Chem.* **2005**, *48*, 3182–3193.
- (15) Kenny, P. W.; Sadowski, J. Structure Modification in Chemical Databases. In *Chemoinformatics in Drug Discovery*; Oprea, T. I., Ed.; Wiley-VCH: Weinheim, Germany, 2004; pp 271–285.
- (16) 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.
- (17) Renner, S.; van Otterlo, W. A. L.; Seoane, M. D.; Möcklinghoff, S.; Hofmann, B.; Wetzel, S.; Schuffenhauer, A.; Ertl, P.; Oprea, T. I.; Steinhilber, D.; Brunsvel, L.; Rauh, D.; Waldmann, H. Bioactivity-guided

Mapping and Navigation of Chemical Space. *Nat. Chem. Biol.* **2009**, *5*, 585–592.

(18) Wetzel, S.; Wilk, W.; Chammaa, S.; Sperl, B.; Roth, A. G.; Yektaoglu, A.; Renner, S.; Berg, T.; Arenz, A.; Giannis, A.; Oprea, T. I.; Rauh, D.; Kaiser, M.; Waldmann, H. A Scaffold-Tree-Merging Strategy for Prospective Bioactivity Annotation of γ -Pyrones. *Angew. Chem.* **2010**, *122*, 3748–3752.

(19) Hu, Y.; Bajorath, J. Structural and Potency Relationships between Scaffolds of Compounds Active against Human Targets. *ChemMedChem* **2010**, *5*, 1681–1685.

(20) ChEMBL; European Bioinformatics Institute (EBI): Cambridge, 2010. <http://www.ebi.ac.uk/chembl/> (accessed May 11, 2010).

(21) Liu, T.; Lin, Y.; Wen, X.; Jorissen, R. N.; Gilson, M. K. BindingDB: a Web-Accessible Database of Experimentally Determined Protein-Ligand Binding Affinities. *Nucleic Acids Res.* **2007**, *35*, D198–D201.

(22) Shannon, P.; Markiel, A.; Ozier, O.; Baliga, N. S.; Wang, J. T.; Ramage, D.; Amin, N.; Schwikowski, B.; Ideker, T. Cytoscape: a Software Environment for Integrated Models of Biomolecular Interaction Networks. *Genome Res.* **2003**, *13*, 2498–2504.

(23) *Molecular Drug Data Report (MDDR)*; Symyx Software: San Ramon, CA, 2008.

(24) Wishart, D. S.; Knox, C.; Guo, A. C.; Cheng, D.; Shrivastava, S.; Tzur, D.; Gautam, B.; Hassanali, M. DrugBank: A Knowledgebase for Drugs, Drug Actions and Drug Targets. *Nucleic Acids Res.* **2008**, *36*, D901–D906.

(25) *Molecular Operating Environment (MOE)*; Chemical Computing Group Inc.: Montreal, Quebec, Canada, 2009.

(26) *Pipeline Pilot*, Student ed., version 6.1; Accelrys, Inc.: San Diego, CA, 2007.

(27) Weininger, D. SMILES, a Chemical Language and Information System. 1. Introduction to Methodology and Encoding Rules. *J. Chem. Inf. Comput. Sci.* **1988**, *28*, 31–36.