

## Searching Fragment Spaces with Feature Trees

Uta Lessel,<sup>\*,†</sup> Bernd Wellenzohn,<sup>†</sup> Markus Lilienthal,<sup>‡</sup> and Holger Claussen<sup>‡</sup>

Department of Lead Discovery, Boehringer Ingelheim Pharma GmbH & Co. KG, Birkendorfer Strasse 65, 88397 Biberach an der Riss, Germany, and BioSolveIT, An der Ziegelei 75, 53757 St. Augustin, Germany

Received August 6, 2008

Virtual combinatorial chemistry easily produces billions of compounds, for which conventional virtual screening cannot be performed even with the fastest methods available. An efficient solution for such a scenario is the generation of Fragment Spaces, which encode huge numbers of virtual compounds by their fragments/reagents and rules of how to combine them. Similarity-based searches can be performed in such spaces without ever fully enumerating all virtual products. Here we describe the generation of a huge Fragment Space encoding about  $5 \times 10^{11}$  compounds based on established in-house synthesis protocols for combinatorial libraries, i.e., we encode practically evaluated combinatorial chemistry protocols in a machine readable form, rendering them accessible to *in silico* search methods. We show how such searches in this Fragment Space can be integrated as a first step in an overall workflow. It reduces the extremely huge number of virtual products by several orders of magnitude so that the resulting list of molecules becomes more manageable for further more elaborated and time-consuming analysis steps. Results of a case study are presented and discussed, which lead to some general conclusions for an efficient expansion of the chemical space to be screened in pharmaceutical companies.

### INTRODUCTION

At the beginning of the search for leads for different targets at Boehringer Ingelheim, a High Throughput Screening (HTS) campaign is usually started in which the screening pool with about 1 million compounds is tested. Compounds from external vendors are acquired and added to the screening collection in order to increase the chances of detecting valuable hits. But even if all externally available compounds were to be added, still only a negligibly small part of the chemical universe could be screened.

Estimates of the size of the chemical universe range from  $10^{13}$  up to  $10^{180}$  virtual compounds.<sup>1–7</sup> However, independent of the exact number, the question is whether it is possible to exploit this huge space. Problems related to this question are how to generate and store such a huge number of compounds and how to screen them. One general solution is to use Fragment Spaces. A Fragment Space is a set of small compounds (fragments) with one or more attachment points (linkers) and a set of rules describing how these fragments can be combined, i.e., which linker types are compatible with each other. Already a rather small set of fragments can span a huge set of virtual compounds due to the “combinatorial explosion”.<sup>3</sup>

One critical issue surrounding Fragment Spaces is the chemical feasibility of the compounds generated by the recombination of fragments. Because of this, in some approaches a set of retrosynthetic rules is defined to cut compounds into fragments.<sup>8,9</sup> This way, in principle, the products of combined fragments should be synthetically accessible. This type of Fragment Space is already used for

similarity searches and *de novo* design.<sup>10,11</sup> However, all retrosynthetic rules described in the literature so far consider only a very limited number of generally applicable reactions and are unable to reflect more specific reactions and reaction conditions which ultimately determine the synthetical accessibility of a particular compound. Especially the chemical feasibility of the designed products often suffers from at least one of the following drawbacks: (1) the actual availability of the reagents is not taken into account, (2) combinations of particular fragments may not be feasible though they are based on known reactions, and last but not least (3) in some cases the combination of several fragments via different reactions in the same product may exclude each other.

Here we present another way to generate Fragment Spaces that circumvents these drawbacks. We base the generation of Fragment Spaces on already existing in-house combinatorial libraries, which are validated and well established at Boehringer Ingelheim.

A combinatorial library usually consists of a core and “R-groups”. These are fragments, which may be combined in a certain way that can be expressed in quite simple rules. We encode these fragments and connection rules in a machine readable form so that they become accessible to *in silico* similarity search methods. The result of a similarity search in this Fragment Space is not only a list of similar molecules but also a link to the underlying synthesis protocols via the names of the fragments. In other words a query selects similar, complete molecules (not only scaffolds) from the combinatorial search space, which can be combined on the basis of the encoded fragments and connection rules that represent already existing synthesis protocols. These protocols can be reidentified because we can map the fragment names of the selected molecules back to particular protocols. This way, the degree of chemical feasibility is controlled by

\* Corresponding author phone: +49-7351-543062; e-mail: Uta.Lessel@boehringer-ingelheim.com.

<sup>†</sup> Boehringer Ingelheim Pharma GmbH & Co. KG.

<sup>‡</sup> BioSolveIT.

the quality of the synthesis protocols and the reagent lists used for the generation of the Fragment Space. If desired, the actual availability of reagents can be ensured in advance by using lists with currently available reagents for the generation of the Fragment Space.

There are already some existing approaches to set up Fragment Spaces that avoid splitting known molecules by retrosynthetic rules. The ChemSpace technology<sup>12,13</sup> and its successor AllChem,<sup>14</sup> both developed at Tripos, use the Topomer search methodologies<sup>15</sup> to navigate through known chemistry. The latter approach starts with about 7000 commercially available building blocks and applies a set of about 100 encoded universal reaction mechanisms to generate fragments with open valences, which are called Synthons. This step basically identifies and marks reactive functional groups within the building blocks that can react to form larger molecules. The result of an AllChem search is a list of similar molecules with often multiple synthetic routes, which consist of several reaction steps each. The runtime for searches in the full-scale AllChem database takes several hours. Nikitin et al.<sup>2</sup> set up a large virtual diversity space of  $10^{13}$  compounds based on 400 combinatorial libraries described in the literature, which they extended by adding larger collections of chemical reagents from vendor catalogs. Their structure-based *de novo* design program generates virtual candidate ligands within the time frame of several months. In comparison, the ligand-based similarity searching approach presented here requires only minutes on a single CPU. Böhm et al.<sup>16</sup> describe the generation of a huge virtual library based on a large proportion of Pfizer's combinatorial chemistry protocols. A total of 358 combinatorial libraries have been converted into a single concise Feature Trees Fragment Space<sup>3</sup> comprising a total of  $3 \times 10^{13}$  virtual products. The validation and application of FTrees<sup>17,18</sup> searches in this space was reported, and the authors conclude that this method is viable for lead finding across diverse scaffolds.

At Boehringer Ingelheim an improved version of Böhm's approach based on established combinatorial libraries has been used to generate a large virtual library, which is called BI CLAIM (Boehringer Ingelheim Comprehensive Library of Accessible Innovative Molecules). By including almost all Boehringer Ingelheim combinatorial libraries with appropriate reagent lists the whole in-house knowledge on combinatorial chemistry becomes searchable. The result of a BI CLAIM search is first and foremost a list of components that are similar to a query and that can be used for further analysis like any conventional screening library. But in addition, the names of the core fragments of these hits also refer to the protocols of how to synthesize these molecules. Therefore, the compounds are accessible via internally known combinatorial synthesis protocols. This also implies that it is easy to synthesize not only a single molecule but also a whole series of analogous compounds.

One of the new features since the work of Böhm is that during the generation of the FTrees Fragment Space the molecular environment around the linkers is taken into account. Furthermore, an option has been implemented in the FTrees Fragment Space searches, which prevents FTrees from returning results with dangling, unsatisfied linkers or incomplete compounds in terms of the underlying combinatorial library.

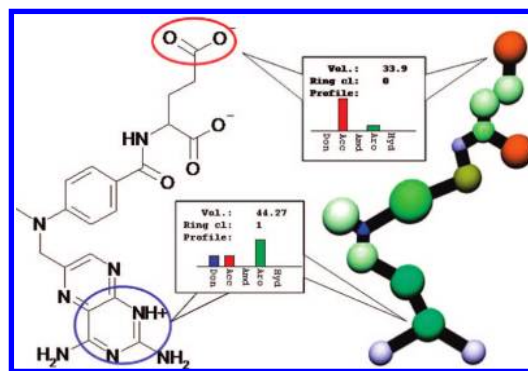


Figure 1. Representation of a molecule as Feature Tree.

In this paper the generation of the BI CLAIM Fragment Space is described. In addition we show how searches in BI CLAIM can be used as a first step in a general workflow in order to reduce the extremely huge number of virtual products by several orders of magnitude so that the resulting list becomes accessible for further more elaborated and time-consuming analysis steps. In a retrospective case study we show that our search method with FTrees Fragment Spaces is able to select the correct part of the space. Finally, we report briefly on two prospective applications of the method in recent projects.

## GENERATION OF FRAGMENT SPACES

We use the FTrees software<sup>17,18</sup> and its extension FTrees Fragment Spaces (FTrees-FS),<sup>3</sup> which is able to handle and search large virtual combinatorial libraries without explicitly enumerating all possible virtual product structures, to perform similarity searches in the generated Fragment Space.

Feature Trees are descriptors that represent the molecule as a reduced graph. The graph encodes functional groups as well as rings to single nodes as shown in Figure 1.

The physicochemical properties of the substructure represented by a node are stored in a property profile for that node. The overall topology is preserved in the Feature Tree, i.e. nodes representing fragments that are connected in the molecule are also connected in the Feature Tree. FTrees is able to select that topology-preserving mapping belonging to the best matching overlay of two molecules both represented by their corresponding Feature Trees. For each match of nodes in either tree the Tanimoto distance between the corresponding property profiles is calculated. The overall similarity of the compounds is simply the normalized sum of these local similarities.

The Feature Tree descriptor is a rather fuzzy description of the molecule, ignoring its 3D structure as well as chirality. It preserves only the overall topology of the Feature Tree nodes but within the nodes FTrees compares only property profiles without taking into account the exact positions of certain functionalities. If one thinks, for example, of substituents of a ring system, the nodes of ortho-, meta-, and para-substituted groups are all treated in the first instance in an identical way as direct neighbors of the ring node, and very small substituents like methyl groups or halogens are even subsumed to the ring node. However, within the ring node there are also optional properties for the shortest and longest path through the underlying substructure so that two rings, which are only ortho- or para-substituted, actually

differ slightly in their profiles. In other words, FTrees focuses basically only on the presence or absence of certain functionalities in roughly the right position, and thus it is able to detect remote similarities and has a potential for scaffold hopping.<sup>19</sup> Additionally this abstraction significantly reduces the runtime. Since all rings are collapsed to nodes (large macrocycles are excluded from processing by default as they would collapse to one big node), the molecule graphs can be reduced to topological trees, which can be compared by fast algorithms. Therefore, FTrees calculates the optimum similarity value between two compounds within a millisecond. However, this is still not fast enough to search through an enumerated Fragment Space with FTrees.

FTrees-FS uses dynamic programming techniques, which are described in detail in the original publication,<sup>3</sup> in order to recursively detect a sizable set of the highest similar fragments and assemble multiple fragments to virtually grow a set of Feature Trees from the Fragment Space. When searching the Fragment Space the query molecule is also converted into a Feature Tree, and the whole comparison is based on the Feature Trees descriptor. The method is deterministic. The efficiency of searching combinatorial libraries encoded in a corresponding Fragment Space versus their enumerated product space can be easily explained by the different numbers of molecules that have to be compared during a similarity search. Assuming three sets of reagents, e.g. 100 each, the enumeration would result in a library with  $100^3$  products. Thus 1 million compounds would have to be processed, whereas only 300 fragments have to be compared in the Fragment Space. This reduces the search effort by 4 orders of magnitude and allows for searching a large virtual Fragment Space comprising a total of more than  $10^{13}$  virtual products.

BI CLAIM (Boehringer Ingelheim Comprehensive Library of Accessible Innovative Molecules) contains almost all of the combinatorial libraries established at Boehringer Ingelheim with appropriate reagent lists. Thus it spans the whole space of compounds which are potentially accessible with BI's in-house knowledge on combinatorial chemistry. The BI CLAIM virtual compound space still grows because new libraries are added continuously. Currently, BI CLAIM contains about 1600 scaffolds and 30 000 reagents. Considering the compatibility of the scaffolds and the reagents, the compound space altogether encodes about  $5 \cdot 10^{11}$  virtual compounds.

To make a transition from real synthetic combinatorial chemistry to an *in silico* representation as a Fragment Space, a couple of manual preprocessing steps are required. Then CoLibri<sup>20</sup> is used to generate a Fragment Space ready for searching with FTrees-FS. The manual preprocessing steps as well as more details on how we enabled CoLibri to take into account the molecular environment around the linkers during the preprocessing and how we prevented FTrees-FS from returning partial (i.e., unsaturated) results are described in the Materials and Methods section of this paper. After a preprocessing stage, the creation of the BI CLAIM Fragment Space with CoLibri took about 8 h on a single Intel Xeon CPU. This procedure has to be done once at the beginning and has to be repeated thereafter only after an update of the list of libraries and/or the lists of reagents. Only about 270 MB disk space is needed to store the whole Fragment Space including all information needed for the BI CLAIM searches.

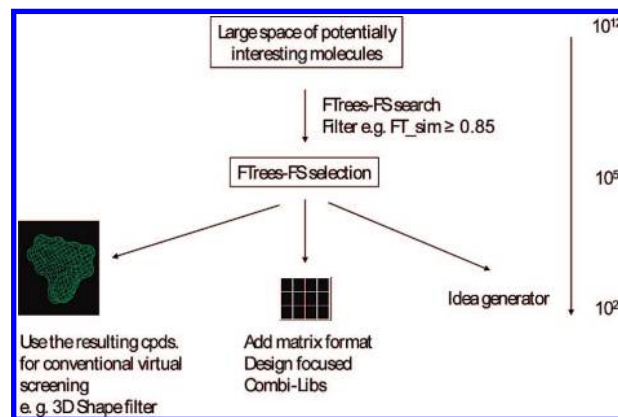


Figure 2. Basic application concept of BI CLAIM searches.

The time needed for a search based on one query in the BI CLAIM Fragment Space depends heavily on the size (number of Feature Tree nodes) of the query. Usually the searches are finished after a few minutes. Less than one minute CPU time is needed for small queries.

### FRAGMENT SPACE SEARCHES

The basic application concept for BI CLAIM (see Figure 2) starts by selecting molecules which are similar to a known hit from the huge space with nearly  $10^{12}$  virtual compounds using FTrees-FS. In most cases we limit the results to compounds with an FTrees similarity greater than 0.85 based on literature data, e.g. ref 21. This FTrees-FS selection is reduction by almost 7 orders of magnitude compared to the full BI CLAIM space. The resulting list of complete molecules can be used for conventional virtual screening like any conventional screening library. Furthermore, we get the scaffolds and thus the reaction protocols which can be used for the synthesis of focused combinatorial libraries. Last but not least, in some cases the results simply serve as idea generators.

In a typical workflow at Boehringer Ingelheim an FTrees-FS search is used in order to select a small subset of BI CLAIM to be analyzed in more detail. Pharmacophore filters or simple 3D shape filters are applied on the FTrees-FS selection followed by a visual inspection of the compounds passing the filters. At this point interesting scaffolds are selected. The prioritization of scaffolds is influenced not only by the 3D alignment and the matching of pharmacophores but also by project specific knowledge and—if available—biological data from existing compounds of the corresponding library. Interesting scaffolds are of course those which are highly populated in the FTrees-FS selection because this shows that not only a single representative of a library is similar to the query. However, it is also worth looking at lower populated scaffolds since these may lead to complete new ideas and options for decoration. For all preselected scaffolds 3D alignments with the query structures are generated in order to gain more insights into the conformational arrangements. The results of these analyses are discussed with the project chemists and the combinatorial chemists to decide which of the scaffolds should be followed up. For the interesting scaffolds focused combinatorial libraries are designed and synthesized or prototypes are selected and synthesized or purchased to support the underlying hypothesis. Again, for the design of focused



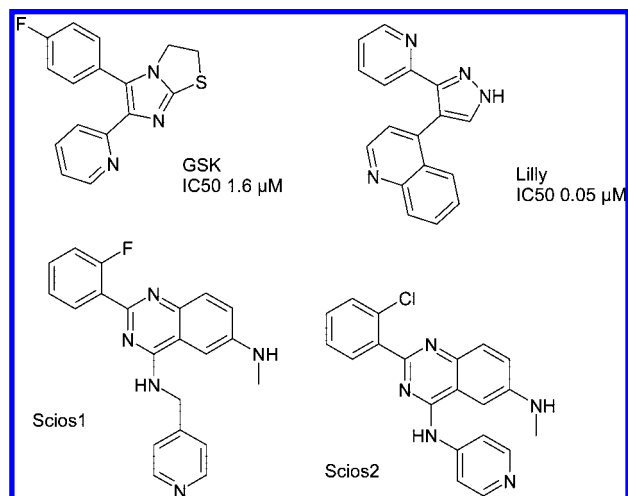


Figure 3. Query structures used for the TGF $\beta$  case study.

libraries and the selection of prototypes, the hits contained in the FTrees-FS selection can build a first basis for the reagent selection, but for the detailed design additional information such as 3D alignments and experiences of the project teams about, for example, preferred moieties are taken into account. Altogether there are strong interactions between the project teams, the combinatorial chemists, and the computational chemists performing the BI CLAIM searches.

### CASE STUDY

In principle, it is very difficult to validate Fragment Space searches because most of the compounds in large Fragment Spaces do not exist physically and consequently have not been tested against any target. Nevertheless, we wanted to show the potential of Fragment Space searches and designed a retrospective study to find compounds which have earlier been found to be TGF $\beta$  inhibitors in an HTS campaign performed at Boehringer Ingelheim. As for any other target, nearly all of the 500 billion compounds in BI CLAIM have unknown TGF $\beta$  activity. During the HTS campaign for TGF $\beta$  inhibitors 739 605 compounds were screened, and 6 343 compounds with an IC<sub>50</sub> < 100  $\mu$ M were detected. Defining these compounds as hits yields a hit rate of 0.9%. Note that if all these 739 605 compounds were contained in BI CLAIM, this would already only amount to 1.5 compounds per million of the BI CLAIM compounds. The fraction of TGF $\beta$  inhibitors with a known IC<sub>50</sub> value below 100  $\mu$ M would be about 0.013 ppm of BI CLAIM, if all 6 343 hits were contained in this set. Unfortunately, however, the number of compounds with known TGF $\beta$  activity that are also in BI CLAIM is even considerably lower than 0.01 ppm, because not all compounds tested in the HTS come from combinatorial chemistry. So the statistical chance of actually selecting such a compound that was in fact tested in the HTS campaign and found to be active is far below 1:100 000 000, i.e. you would have to select more than 100 million compounds randomly in order to find such a molecule by chance.

In spite of this statistical background, we still tried to recover these active molecules from the HTS campaign by an FTrees-FS search in the BI CLAIM Fragment Space. We used the four known TGF $\beta$  inhibitors<sup>22–25</sup> shown in Figure 3 as queries to select similar compounds from the space.

Table 1. Feature Trees Similarities of the Queries to Each Other

query 1	query 2	Feature Trees similarity
GSK	Lilly	0.7984
GSK	Scios1	0.7957
GSK	Scios2	0.7812
Lilly	Scios1	0.7761
Lilly	Scios2	0.7635
Scios1	Scios2	0.9608

Table 2. Summary of the Virtual Hits Detected

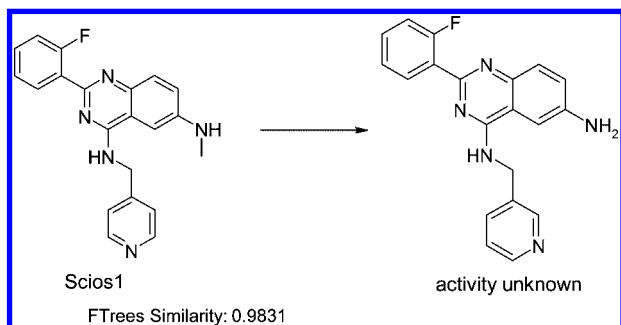
query	number of compounds in FTrees-FS selection			TGF $\beta$ inhibitors in FTrees-FS selection		
	total	not tested in HTS	tested in HTS	number	IC <sub>50</sub> ( $\mu$ M)	FTrees sim. to query
GSK	7653	7638	15	none	—	—
Lilly	7124	7109	15	1	34.4	0.9173
Scios1	23124	23106	18	1	5.8	0.9314
Scios2	19194	19136	58	none	—	—

The FTrees similarities of the queries to each other are shown in Table 1. Except for the two very similar Scios structures the similarities between the query compounds range between 0.76 and 0.80. For each query we select all molecules with an FTrees similarity above 0.85. This means we explored the environment of 3 different activity islets in the FTrees descriptor space. The results of the searches are summarized in Table 2. This table contains the total number of retrieved molecules with an FTrees similarity above 0.85 for each query and the number of compounds within this selection that are (not) tested in the HTS campaign. In addition the number of actives within the tested subset of the selection, their IC<sub>50</sub> values, and their FTrees similarities to the corresponding queries are given.

The FTrees-FS selection based on the GSK compound contains 7653 BI CLAIM compounds. Fifteen of these compounds were found to be inactive in the HTS campaign. On the other hand the FTrees-FS selection contained 7638 virtual compounds with unknown activity. The results for the other queries are similar. In the case of the Lilly compound, 15 out of 7124 compounds in the FTrees-FS selection have actually been tested before, and 1 of them showed an IC<sub>50</sub> value of 34.4  $\mu$ M in HTS. This compound has an FTrees similarity of 0.9173 to the query molecule. Scios1 provided 18 compounds (1 active), which were tested in HTS out of about 23 000 compounds in the FTrees-FS selection. Scios2 detected 58 inactive compounds together with more than 19 000 virtual compounds in the FTrees-FS selection with unknown activity. The overlap between the FTrees-FS selections of the four queries is—as expected—rather small, because the similarity threshold for each selection is higher than the FTrees similarity among the queries themselves. Altogether 47 489 virtual compounds show an FTrees similarity above 0.85 to only one of the queries and 4800 to two of the queries (mostly to Scios1 and Scios2), and two compounds are contained in the FTrees-FS selections of GSK, Scios1, and Scios2.

The active compound detected in the FTrees-FS selection of the Lilly compound is found on rank 105. This is the first of the 15 compounds with known activity in this list.

The active with an IC<sub>50</sub> of 5.8  $\mu$ M detected with the query Scios1 is located at rank 496. It is the fifth compound with known activity.



**Figure 4.** Trivial virtual hit resulting from BI CLAIM searches based on the query Scios1.

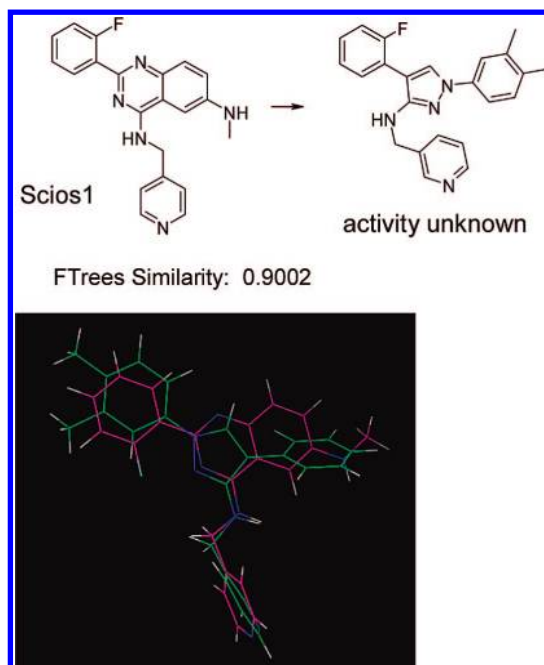
Due to the huge amount of virtual compounds with unknown activity in the FTrees-FS selection it is impossible to assess enrichment factors. This is a general problem occurring during the evaluation of Fragment Space searches. Nevertheless, altogether at least 2 known actives from the HTS campaign out of 500 billion compounds were detected. Taking into account that statistically only one known active within 100 000 000 compounds can be expected in BI CLAIM, it is remarkable that this approach is able to spot at least two of these within an overall selection of about 50 000 compounds. Note that we did not actually choose any compounds from the FTrees-FS selection and test them in a prospective manner, but these hits were found in a retrospective study. The compounds tested in HTS are a more or less arbitrary overlap between the molecules tested in the HTS campaign and the FTrees-FS selection. Of course, in a real prospective workflow further filters would be applied on the selection followed by a visual inspection in order to select compounds to be synthesized and tested.

A more detailed discussion of the scaffolds detected in the case study can further illustrate that the results are quite reasonable and useful.

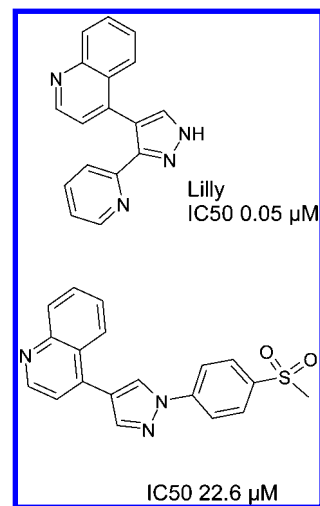
One member of the FTrees-FS selection based on the query Scios1 is the quinazoline shown in Figure 4, which has not been tested, but the high similarity (FTrees similarity: 0.9831) to the query suggests that this compound might be active. For validation purposes it is important that the scaffold of the query is detected. On the one hand, this is of course a trivial result and should be found by any program. On the other hand however, considering the vast size of the space, even such a trivial task could not be completed by most other approaches.

The same FTrees-FS selection contained also some more interesting, nontrivial virtual hits like the pyrazole shown in Figure 5a. This is one example of 2504 pyrazoles contained in the FTrees-FS selection from query Scios1. Figure 5b shows that the virtual hit and the query fit well in 3D space.

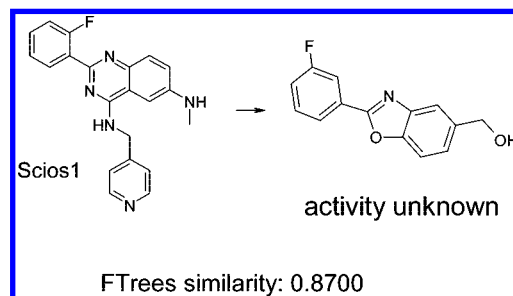
Other pyrazoles shown in Figure 6, which are known TGF $\beta$  inhibitors, also support the idea that this scaffold may be of interest. The Lilly compound belongs to the same class of pyrazoles and has an FTrees similarity of 0.8318 to the virtual hit. The second compound shown from the Boehringer Ingelheim screening pool (which is not part of BI CLAIM) has an IC<sub>50</sub> of 22.6  $\mu$ M and an FTrees similarity of 0.7823 to the virtual hit and 0.7626 to the query Scios1. While assessing this result, one has to take into account that the aim of a BI CLAIM search is not the detection of a single virtual hit but the detection of a structural class, in this case pyrazoles. In practice not only the single hit but also a whole



**Figure 5.** Virtual pyrazole hit resulting from the BI CLAIM search based on the query Scios1 a) structure and b) 3D alignment.



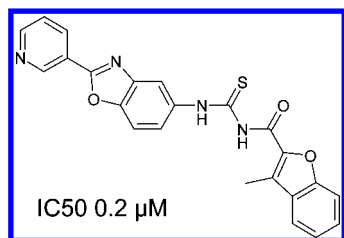
**Figure 6.** Known TGF $\beta$  inhibitors from the pyrazole class.



**Figure 7.** Virtual benzoxazole hit resulting from the BI CLAIM search based on the query Scios1.

array of pyrazoles should be synthesized and tested in the next step.

A similar example from the same BI CLAIM search is shown in Figure 7. It shows one of 33 benzoxazoles in the FTrees-FS selection. Again, a small combinatorial library or at least a series of similar compounds would be synthesized to confirm the hypothesis that benzoxazoles may be



**Figure 8.** Known TGF $\beta$  inhibitors from benzoxazole class. good starting points for the development of TGF $\beta$  inhibitors. To illustrate that this result is at least plausible Figure 8 shows a benzoxazole, which is not part of the BI CLAIM Fragment Space, but was detected in the HTS campaign. Its FTrees similarity to the benzoxazole in the FTrees FS selection shown in Figure 7 is 0.8346.

### PROSPECTIVE APPLICATION EXAMPLES

Finally, the results of two successful prospective application projects will be summarized to illustrate the potential of the BI CLAIM searches. For obvious reasons we cannot disclose too many details about these projects. Thus we summarize these only very briefly just to emphasize that we already used this set up successfully in our everyday work. In the first project an FTrees-FS search was performed based on a known GPCR ligand from the literature, which yielded some 1 000 hits. These hits were analyzed with ROCS<sup>26</sup> in a postfiltering step regarding their shape similarity to the query. In a manual selection step two compound classes were selected, and for each of them a library with a few hundred compounds was synthesized. One of these libraries provided hits with an activity of 5–10  $\mu$ M, and the second library yielded hits with IC<sub>50</sub> values of about 100 nM.

In the second project the BI CLAIM search was based on a known protease inhibitor taken from the literature. In this case, a prescreening library of about 1200 compounds from 10 different scaffolds was produced. In HTS, hits with IC<sub>50</sub> values greater than 10  $\mu$ M were detected in this library. Resyntheses yielded actives with IC<sub>50</sub> values between 5 and 10  $\mu$ M derived from 2 scaffolds. For further optimization, one of the compounds could be cocrystallized resulting in hints for a favorable decoration of the scaffold. With this knowledge, a further optimization cycle provided compounds in the 10 nM range.

### CONCLUSIONS

The examples show that FTrees searches in Fragment Spaces built on synthesis protocols for combinatorial libraries allow for an enormous expansion of the chemical space accessible for the detection of lead candidates by several orders of magnitude.

A large pool of reagents allows both diverse and focused decoration of scaffolds detected by the searches.

For the whole workflow, knowledge of combinatorial chemistry and dedicated chemistry capacity for the synthesis of virtual screening hits are needed. However, due to the validated chemistry protocols in BI CLAIM the timelines for syntheses fit to the projects' needs.

### MATERIAL AND METHODS

Besides generating a machine-readable description of reactions, it is also necessary to prepare the input structures

in a way that they can be used in FTrees-FS. Here we describe step by step how we transform real synthetic combinatorial chemistry into a Fragment Space:

First of all the synthetic reaction protocols are encoded as a virtual reaction (see Figure 9).

Functional groups involved in the reaction are replaced by unique attachment points, which we call linker atoms or simply linkers. In addition any replacement or removal of protecting groups as well as any other mechanism (e.g., solid phase support) occurring in real synthetic reactions that cannot be directly translated into a virtual reaction scheme are resolved manually (e.g., the Cbz-protection group in Figure 9). Basically all such additional groups are removed so that the fragments look like the one present in the final product. In other words, the description of a combinatorial reaction protocol is encoded such that "building blocks" remain, with which the reaction can be formally described. The compatibility between the linkers determines how the fragments may be connected (linker compatibility rules). Products of a combinatorial library are created by connecting the respective neighbor atoms of the two compatible linkers with a bond. The linker atoms themselves are subsequently eliminated, and a bond is formed.

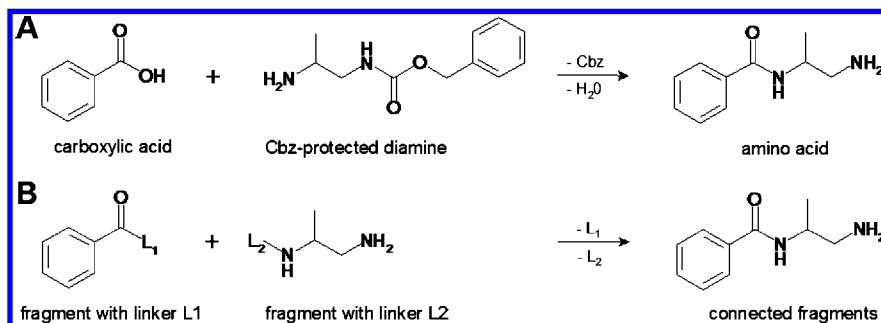
Since FTrees technology is used to search the Fragment Space the input fragments have to meet two further constraints: a) fragments may not be connected such as to close a ring, and b) linkers may not be part of a ring. These constraints are due to the Feature Trees descriptor that collapses rings to single nodes in the tree representation. Therefore, rings have to be complete at the time of generating the Feature Trees for the fragments as illustrated in Figure 10.

In comparison to the reagents themselves, the fragments are stripped off from those parts of the molecules that form the ring and replaced by the linker. The resulting ring-cores are manually sketched and stored in MDL mol format.<sup>27</sup> The reagents are coded as SMILES.<sup>28</sup>

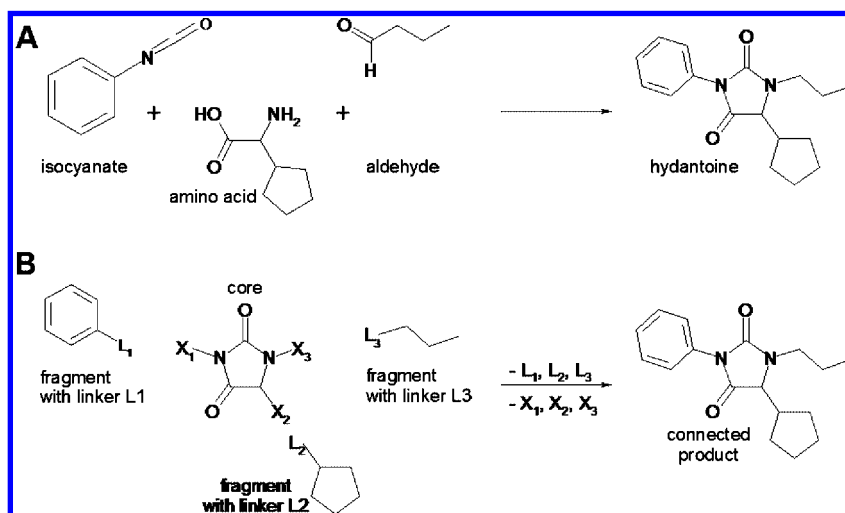
After translating each combinatorial library into such a machine readable virtual reaction scheme, we use CoLibri (Compound Library Toolkit)<sup>16,20</sup> to gather all combinatorial reaction schemes and encode them in one single Fragment Space. There are basically three tasks that are done in this step: (1) unifying the linker names, (2) identifying and removing duplicate fragments, and (3) determining the (Sybyl) atom types of the linkers. Finally, FTrees is used to generate the Feature Trees descriptors for the fragments with the linkers.

**(1). Unifying Linker Names.** In most cases the virtual reaction schemes of the different combinatorial libraries are set up independently. This is why in many cases the same linker names are used in different reactions (e.g., L1), which may cause uncertainties in the results of searches in a Fragment Space including several reaction schemes. To avoid this CoLibri renames all linkers so that they are labeled in a unique way across the entire Fragment Space. Overall BI CLAIM contains 700 unique linkers and 974 041 different raw fragments.

**(2). Removing Duplicate Fragments.** Since the reagents are used for several different scaffolds, there is a high level of redundancy of fragments. It is desirable to remove duplicate fragments in order to keep size and search times low. However, since the fragments are used in different



**Figure 9.** Chemical reaction (A) and virtual reaction scheme (B) for the generation of amides from a carboxylic acid and a Cbz-protected diamine.



**Figure 10.** Modification of a synthetic reaction scheme with a ring closure (A) to a virtual reaction (B), which can be handled by CoLibri and the FTrees Fragment Space searches.

libraries, they can have different linkers, and the fragments lists for the different libraries do not necessarily overlap completely. Thus we use CoLibri to reorganize the linkers and create nonredundant subsets of fragments that share the same compatibilities. Note that two partially overlapping fragment lists A and B from two different libraries result in three subsets: A without B ( $A \setminus B$ ), B without A ( $B \setminus A$ ), and the intersection of A and B ( $A \cap B$ ). While the original linker can be used for the first two sets, the latter overlapping subset requires a third linker type, which is compatible with both core fragments of the respective two libraries. Due to this effect the number of linkers can become quite large if a lot of identical fragments have different linker compatibilities. In the case of BI CLAIM the number of fragments is reduced to 21774 unique fragments, whereas the number of linkers increases only marginally from 700 to 766, because only very few intersections between reagents lists have to be resolved by introducing additional linker types. CoLibri maintains a look-up table, which contains all information about the combinatorial libraries, their reagents, and the corresponding linkers.

**(3). Determining Atom Types of the Linkers.** FTrees uses the Sybyl atom type for creating the property profiles of an FTrees node. However, the property profile depends not only on the properties of the atoms that are represented by the node but also the types of atoms in the neighborhood. For a link atom, the properties of the atom that will be attached when connecting fragments can be estimated. Thus, the atom types of the linkers and their neighbors in the fragments have to be known when the FTrees descriptors

are generated for the fragments. They are defined in the Fragment Space description file. So far, the atom types chosen for the neighbor atom were those that were most frequently found in the fragment list itself, and for the linker the atom types chosen were those that were found most frequently in the list of the counter fragments for the linker. However, this is a rather coarse approximation because in some cases different elements can be attached at the same attachment point and even the same element can change its Sybyl atom type depending on the resulting molecular environment. The most prominent case for this is where a nitrogen atom is either N.3 (nonplanar, trigonally coordinated N) or N.am (amide-N, see Figure 11 for an example).

For this reason in the present method we changed the way the Sybyl atom types of the linker and their neighbors are determined. For each connection that can be formed, the resulting molecular environment is analyzed. The FTrees search is based only on the FTrees descriptors of the fragments. It would be much too time-consuming if each (partial) molecule was recreated and initialized when two FTrees fragments are connected in order to get the resulting atom types. Thus, we implemented the analysis of the molecular environment as a preprocessing step in CoLibri before generating the FTrees descriptors. CoLibri distinguishes linkers representing different Sybyl atom types. The respective fragment is duplicated such that there is one instance for each linker with a different atom type. This increases the number of linkers. In the present study we finally end up with 2653 different linkers and 67 360 unique



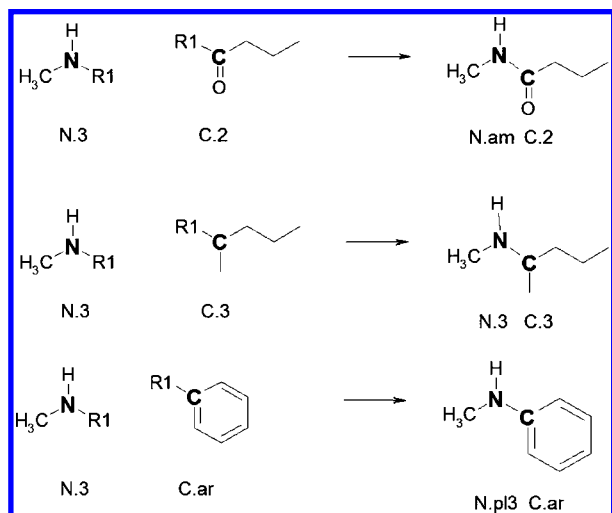


Figure 11. Assembling fragments can cause an atom type change.

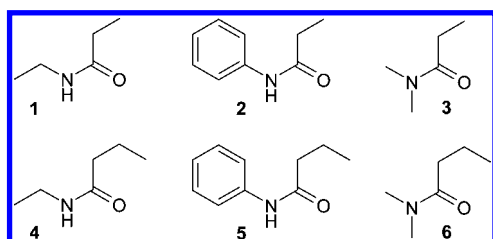


Figure 12. Amides used for analysis.

fragments overall. However, this way the representation of fragments is much refined for the later similarity searches.

In order to analyze the effect of considering atom type changes when building products from fragments, we set up the BI CLAIM Fragment Space in both ways (a) with the old method, where the potential change of atom types it not taken into account, and (b) with the new approach described above.

In both Fragment Spaces we searched for six small compounds which are part of BI CLAIM and thus are

Table 3. FTrees Similarities of 6 Amides to Themselves in the Two Fragment Spaces

compound	FTrees sim. old space (a)	FTrees sim. new space (b)
1	0.9593	0.9998
2	0.9742	0.9998
3	0.9589	0.9998
4	0.9667	0.9998
5	0.9763	0.9998
6	0.9641	0.9998

expected to be redetected with an FTrees-FS search. These compounds are generated by combining an amine with a carboxylic acid forming an amide bond (see Figure 12). We took small molecules for this analysis because the effect of changing a single atom type on the overall similarity is the greater the smaller the number of atoms. Table 3 lists the FTrees similarities of the queries to themselves in the two Fragment Spaces for each compound.

With the former approach the atom type and the profile of the linker at the amine fragment is estimated from the unconnected carboxyl environment. This estimation error leads to an FTrees similarity which is considerably below 1.0. With the new approach the atom type and the profile of the linker at the amine fragment is estimated in the connected amide environment. Thus, the predicted similarities are significantly better and almost 1.0 in all cases. The reason why we do not get a similarity of 1.0 in this case is rather technical—it lies in the slight difference between how the FTrees descriptors are generated for query molecules and for fragments and in the mapping of the fragment nodes onto the node of the whole query molecule. Again this effect is greater the smaller the molecules.

In order to see the overall effect we randomly enumerated products for diverse synthesis protocols in BI CLAIM and detected the 100 most similar compounds out of BI CLAIM in both Fragment Spaces. The overall accuracy of the similarities is increased by the new approach, which can be seen from the histograms of the Feature Trees similarities

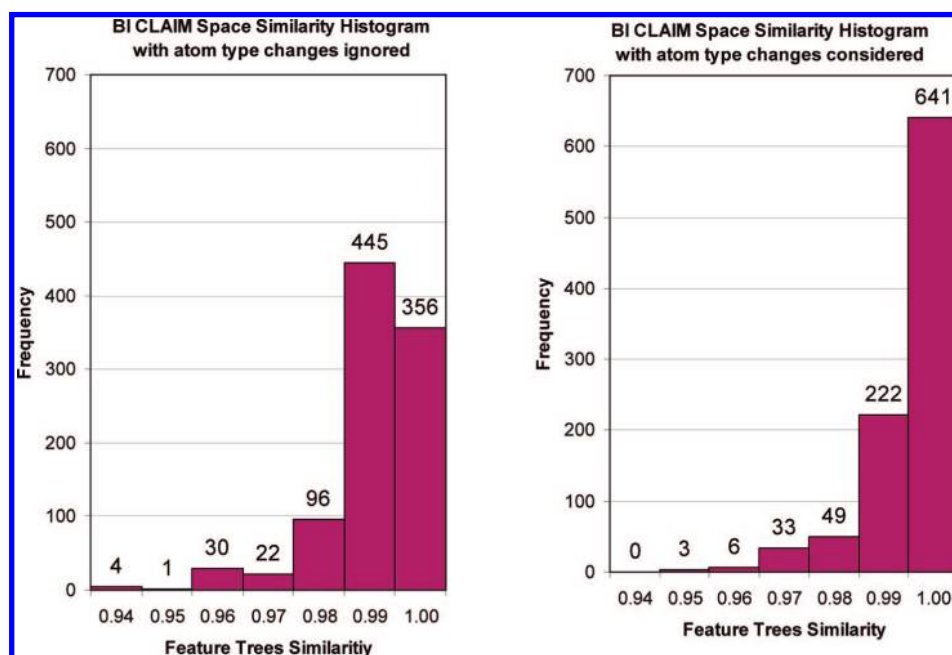
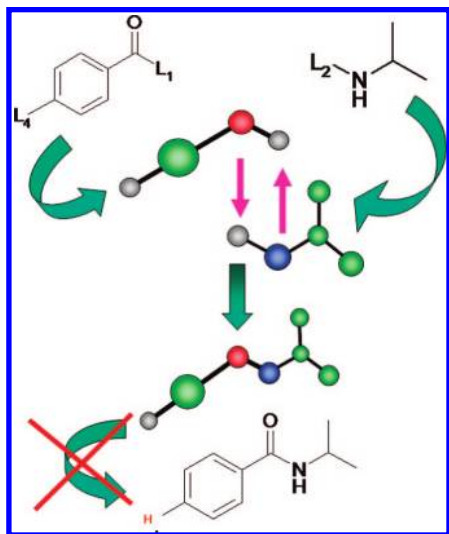


Figure 13. Feature Trees similarity histograms of redetected queries in a BI CLAIM space with the feature to consider atom type changes and in the same BI CLAIM space without this feature.





**Figure 14.** Generation of virtual structures from the Feature Trees Fragment Space. The termination of unsatisfied linkers by a hydrogen atom is prohibited.

of the queries to themselves shown in Figure 13. It includes the results from 955 queries which could be redetected among the first 100 nearest neighbors in both spaces. With the former approach (a) most results (445) have a similarity of 0.99 and only 356 results have a similarity of 1.0, whereas with the new approach the whole distribution is shifted to the right and 641 queries result in a similarity of 1.0.

However, due to the increased number of linkers and unique fragments generated by the new approach the search times rise significantly by a factor of about 5. The amount of prolongation depends heavily on the number of different linker types in the reagent lists and the size of the query molecule. If one is aiming to redetect compounds in a Fragment Space, the consideration of the atom type changes improves the results clearly. But for similarity searches, especially in procedures like the BI workflow, the small changes have only a minor influence on the overall results.

**Generating Feature Trees Descriptors.** Once a Fragment Space has been assembled by CoLibri, it is then encoded as a Feature Tree Fragment Space by FTrees-FS, where the fragments are represented by specially adapted Feature Trees. The generation of virtual structures from the Feature Trees Fragment Space is illustrated in Figure 14.

Special linker nodes define the link compatibility here in an analogous way as the linker atoms for the fragments. When the Feature Trees fragments are combined, these special linker nodes (L1 and L2 in Figure 14) are removed, and the Feature Trees fragments are connected to form a larger tree without changing the remaining FTrees nodes. FTrees-FS can either connect two fragments in this way or terminate the build up with a predefined atom, which is typically a hydrogen atom. This feature is necessary if one uses a Fragment Space that was generated by the shredding of compounds. However, in the context of encoding combinatorial chemistry this feature is distracting, because it allows an FTrees-FS search to produce an incompletely built up compound. Often, it is not even possible to synthesize such compounds as the missing fragment may be vital in the reaction protocol and cannot be left out. Therefore, we extended FTrees-FS in this study in such a way that those linkers can be marked for which the build up process may

not be terminated. Since we are working only with combinatorial libraries, we excluded a termination for all linkers in the Fragment Space. This way we made sure that complete molecules are always yielded as results, i.e. molecules which really belong to one of the combinatorial libraries in BI CLAIM. Thus, we improved the chemical feasibility of the virtual hits.

The Fragment Space definition and its Feature Trees representation are stored together ready for similarity searching with FTrees-FS.

In order to remap the results back to the original combinatorial libraries, CoLibri is used again to retrieve this information from the above-mentioned look-up table that was created during the Fragment Space generation.

## ACKNOWLEDGMENT

The authors thank Marcus Gastreich, Sally Hindle, Herbert Köppen, and Christian Lemmen for fruitful discussions.

## REFERENCES AND NOTES

- (1) Villar, H. O.; Koehler, R. T. Comments on the design of chemical libraries for screening. *Mol. Diversity* **2000**, *5*, 13–24.
- (2) Nikitin, S.; Zaitseva, N.; Demina, O.; Solovieva, V.; Mazin, E.; Mikhalev, S.; Smolov, M.; Rubinov, A.; Vlasov, P.; Lepikhin, D.; Khachko, D.; Fokin, V.; Queen, C.; Zosimov, V. A very large diversity space of synthetically accessible compounds for use with drug design programs. *J. Comput.-Aided Mol. Des.* **2005**, *19*, 47–63.
- (3) Rarey, M.; Stahl, M. Similarity searching in large combinatorial chemistry spaces. *J. Comput.-Aided Mol. Des.* **2001**, *15*, 497–520.
- (4) Ertl, P. Chemoinformatics analysis of organic substituents: identification of the most common substituents, calculation of substituent properties, and automatic identification of drug-like bioisosteric groups. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 374–380.
- (5) Bohacek, R. S.; McMartin, C.; Guida, W. C. The art and practice of structure-based drug design: a molecular modeling perspective. *Med. Res. Rev.* **1996**, *16*, 3–50.
- (6) Walters, P. W.; Stahl, M. T.; Murcko, M. A. Virtual screening—an overview. *Drug Discovery Today* **1998**, *3*, 160–178.
- (7) Gorse, A.-D. Diversity in medicinal chemistry space. *Curr. Top. Med. Chem.* **2006**, *6*, 3–18.
- (8) 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.
- (9) Schürer, S. C.; Tyagi, P.; Muskall, S. M. Prospective exploration of synthetically feasible, medicinally relevant chemical space. *J. Chem. Inf. Model.* **2005**, *45*, 239–248.
- (10) Hartenfeller, M.; Proschak, E.; Schüller, A.; Schneider, G. Concept of Combinatorial De Novo Design of Drug-like Molecules by Particle Swarm Optimization. *Chem. Biol. Drug Des.* **2008**, *72*, 16–26.
- (11) Degen, J.; Rarey, M. FLEXNovo: structure-based searching in large Fragment Spaces. *ChemMedChem* **2006**, *1*, 854–868.
- (12) Cramer, R. D.; Patterson, D. E.; Clark, R. D.; Soltanshahi, F.; Lawless, M. S. Virtual compound libraries: a new approach to decision making in molecular discovery research. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 1010–1023.
- (13) Andrews, K. M.; Cramer, R. D. Towards general methods of targeted library design: Topomer shape similarity searching with diverse structures as queries. *J. Med. Chem.* **2000**, *43*, 1723–1740.
- (14) Cramer, R. D.; Soltanshahi, F.; Jilek, R.; Campbell, B. AllChem: generating and searching 10<sup>20</sup> synthetically accessible structures. *J. Comput.-Aided Mol. Des.* **2007**, *21*, 341–350.
- (15) Jilik, R. J.; Cramer, R. D. Topomers: a validated protocol for their self-consistent generation. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 1221–1227.
- (16) Boehm, M.; Wu, T.-Y.; Claussen, H.; Lemmen, C. Similarity Searching and Scaffold Hopping in Synthetically Accessible Combinatorial Chemistry Spaces. *J. Med. Chem.* **2008**, *51*, 2468–2480.
- (17) Rarey, M.; Dixon, J. S. Feature trees: a new molecular similarity measure based on tree matching. *J. Comput.-Aided Mol. Des.* **1998**, *12*, 471–490.
- (18) *FTrees, version 2.0.2*; BioSolveIT GmbH: Sankt Augustin, Germany, 2008.

- (19) Good, A. C.; Hermsmeier, M. A.; Hindle, S. A. Measuring CAMD technique performance: A virtual screening case study in the design of validation experiments. *J. Comput.-Aided Mol. Des.* **2004**, *18*, 529–536.
- (20) *CoLibri, version 1.1.4*; BioSolveIT GmbH: Sankt Augustin, Germany, 2008.
- (21) Hessler, G.; Zimmermann, M.; Matter, H.; Evers, A.; Naumann, T.; Lengauer, T.; Rarey, M. Multiple-ligand-based virtual screening: methods and applications of the MTree approach. *J. Med. Chem.* **2005**, *48*, 6575–6584.
- (22) Callahan, J. F.; Burgess, J. L.; Fornwald, J. A.; Gaster, L. M.; Harling, J. D.; Harrington, F. P.; Heer, J.; Kwon, C.; Lehr, R.; Mathur, A.; Olson, B. A.; Weinstock, J.; Laping, N. J. Identification of Novel Inhibitors of the Transforming Growth Factor  $\beta$ 1 (TGF- $\beta$ 1) Type 1 Receptor (ALK5). *J. Med. Chem.* **2002**, *45*, 999–1001.
- (23) Sawyer, J. S.; Anderson, B. D.; Beight, D. W.; Campbell, R. M.; Jones, M. L.; Herro, D. K.; Lampe, J. W.; McCowan, J. R.; McMillen, W. T.; Mort, N.; Parsons, S.; Smith, E. C. R.; Vieth, M.; Weir, L. C.; Yan, L.; Zhang, F.; Yingling, J. M. Synthesis and Activity of New Aryl- and Heteroaryl-Substituted Pyrazole Inhibitors of the Transforming Growth Factor- $\beta$  Type I Receptor Kinase Domain. *J. Med. Chem.* **2003**, *46*, 3953–3956.
- (24) Singh, J.; Chuaqui, C. E.; Boriack-Sjodin, P. A.; Lee, W. C.; Pontz, T.; Corbley, M. J.; Cheung, H. K.; Arduini, R. M.; Mead, J. N.; Newman, M. N.; Papadatos, J. L.; Bowes, S.; Josiah, S.; Ling, L. E. Successful Shape-Based Virtual Screening: The Discovery of a Potent Inhibitor of the Type I TGF $\beta$  Receptor Kinase (T $\beta$ RI). *Bioorg. Med. Chem. Lett.* **2003**, *13*, 4355–4359.
- (25) Structures taken from Scios-Patents WO2004/010929 A2, WO 2004/047818A2, WO 2004/048930 A2, WO 2004/087056 A2, WO 2005/032481 A2.
- (26) *ROCS, version 2.3.1*; OpenEye Scientific Software, Inc: Santa Fe, NM, 2007.
- (27) Symyx Technologies, Inc. CTFile Formats, November 2007. <http://www.mdl.com/downloads/public/ctfile/ctfile.pdf> (accessed Dec 17, 2008).
- (28) Weininger, D.; Weininger, A.; Weininger, J. L. SMILES. 2. Algorithm for Generation of Unique SMILES Notation. *J. Chem. Inf. Comput. Sci.* **1989**, *29*, 97–101.

CI800272A