

*promoting access to White Rose research papers*



**Universities of Leeds, Sheffield and York**  
**<http://eprints.whiterose.ac.uk/>**

---

This is an author produced version of a paper published in **Quantitative Structure-Activity Relationships**.

White Rose Research Online URL for this paper:  
<http://eprints.whiterose.ac.uk/3571/>

---

**Published paper**

Bohl, M., Dunbar, J., Gifford, E.G., Heritage, T., Wild, D.J., Willett, R. and Wilton, D.J. (2003) *Scaffold Searching: Automated Identification Of Similar Ring Systems For The Design Of Combinatorial Libraries*, Quantitative Structure-Activity Relationships, Volume 21 (6), 590-597.

---

# Scaffold Searching: Automated Identification Of Similar Ring Systems For The Design Of Combinatorial Libraries

**Martin Bohl<sup>a</sup>, James Dunbar<sup>b</sup>, Eric M. Gifford<sup>b</sup>, Trevor Heritage<sup>c</sup>,  
David J. Wild<sup>b</sup>, Peter Willett<sup>d 1</sup> and David J. Wilton<sup>d</sup>**

Tripos GmbH, Martin-Kollar-Strasse 13, Munich D-81829, Germany<sup>a</sup>

Pfizer Global Research and Development, 2800 Plymouth Road,

Ann Arbor, MI 48105 USA<sup>b</sup>

Tripos Inc., 1699 South Hanley Road, St Louis, MO 63144, USA<sup>c</sup>

Krebs Institute for Biomolecular Research and Department of Information Studies,  
University of Sheffield, Western Bank, Sheffield S10 2TN, UK<sup>d</sup>

**Keywords:** Combinatorial libraries, Geometric searching, Ring systems, Scaffold searching, Shape similarity

**Received on** .....

---

<sup>1</sup> To whom all correspondence should be addressed. Email: p.willett@sheffield.ac.uk

## **Abstract**

Rigid ring systems can be used to position receptor-binding functional groups in 3D space and they thus play an increasingly important role in the design of combinatorial libraries. This paper discusses the use of shape-similarity methods to identify ring systems that are structurally similar to, and aligned with, a user-defined target ring system. These systems can be used as alternative scaffolds for the construction of a combinatorial library.

## Introduction

The important role played by ring systems in drug discovery has meant that much effort has been devoted over very many years to the development of automated methods for their identification, representation and searching (see, *e.g.*, [1-9]). Recent developments in combinatorial chemistry mean that it is now possible to synthesise large libraries of compounds, consisting of a central ring system to which are attached a range of different substituents [10]. An example of such a template is shown in the upper part of Figure 1, where the central ring system acts as a scaffold to position the substituents so that they can make favourable interactions with residues in a protein's binding site. Library definitions such as this are increasingly common in the literature, and the question then arises as to how one might be able to design libraries that are analogous to one that has been published. Specifically there is a need to design libraries in which the functionality can still be positioned at the required positions in 3D space but in which a different central ring scaffold is employed: we use the term *scaffold searching* to refer to the identification of such matching scaffolds. This problem was first addressed by Schneider *et al.* [11], who used similarity measures based on 2D autocorrelation vectors to find alternative topological patterns, but without focusing specifically upon ring systems. Here, we report an approach to scaffold searching that takes full account of the 3D natures of scaffolds using FBSS, a program we have developed previously for similarity searching in chemical databases [12].

## Materials And Methods

**Field-based similarity searching** The last few years have seen increasing interest in measures of inter-molecular structural similarity that are based on steric, electrostatic and hydrophobic field descriptors, an approach first suggested by Carbo *et al.* [13]. Given a molecular property  $P$  that can be calculated at any point around a molecule, a field may be created around that molecule by integrating  $P$  with respect to volume. The similarity between a pair of molecules is then determined by aligning the two

molecules so as to maximise the overlap of the corresponding fields. The similarity is normally calculated using the Carbo index, which is defined to be

$$\frac{\int P_A P_B dv}{(\int P_A^2 dv)^{1/2} (\int P_B^2 dv)^{1/2}}$$

Here,  $P_A$  and  $P_B$  are the properties of the two molecules that are being compared and the integrations are over 3D space, this normally being approximated by summing over all the components of 3D grids that surround the two molecules that are being compared. The precise form of the summation depends on the particular property that is being considered. For example, if  $P_r$  denotes the electron density at a point  $r$ , then the density is calculated from the sum of the contributions from each of the atoms in the molecule, *i.e.*,

$$P_r = \sum_{i=1}^n E_i(r - R_i)$$

where  $E_i(d)$  is the electron density contribution of atom  $i$  ( $1 \leq i \leq n$ ) at distance  $d$  from the nucleus, and where  $r - R_i$  is the Euclidean distance between  $r$  and the position ( $R_i$ ) of the  $i$ -th atom. The resulting  $P_r$  values can then be inserted back into the Carbo index for the calculation of the similarity, but this is very time-consuming unless a coarse grid spacing is used. Good and Richards have shown that Gaussians can be used to fit the curve of electron density against distance from the atomic nuclei, and that these Gaussians can be inserted into a version of the Carbo equation. The similarity is then calculated analytically rather than numerically, giving a very rapid way of calculating the shape similarity between pairs of molecules [14].

There have been several recent reports of systems that use similarity measures based on molecular fields or molecular shape (e.g., [15-18]). FBSS employs a genetic algorithm (hereafter GA) to align two molecules' fields so as to maximise the value of the Carbo index [12]. In brief, each chromosome in this GA encodes the rotations and translations that are to be applied to a database structure to align it with the target structure in a similarity search, and the GA's fitness function is the value of the Gaussian similarity coefficient resulting from that particular encoded alignment.. The program has been used previously for 3D similarity searching and for pre-processing datasets for 3D QSAR analyses [19, 20]. Here we use the program to identify ring systems that are similar in shape to a user-defined target scaffold,  $T$ , and that can be

substituted in the same approximate geometric arrangement as the points of attachment in *T*.

**Scaffold searching** When FBSS is to be used for scaffold searching, the Carbo coefficient of shape similarity (based on electron density as suggested by Good and Richards [14] and as described above) is calculated for the similarity between *T* and each of the rings in a database of ring systems. The ring systems are then ranked in descending order of the calculated shape similarities, together with the corresponding alignment. Each such alignment is then checked to see if the points of attachment in *T* (as denoted by R1, R2 etc. in Figure 1) correspond to potential points of attachment in the ring system from the database, where a potential point of attachment is a ring atom that could, given suitable chemistry, have functionality attached to it. A distance threshold is used to determine whether a substitutable ring atom is an acceptable match for a point of attachment in *T*. For brevity, we refer to this check subsequently as the *attachment search*. The output from the attachment search is hence those rings that could act as alternatives to the ring scaffold, ranked in order of decreasing shape similarity.

The effectiveness of an FBSS search depends on the parameters that are specified for the GA, these including the selection pressure, the number of generations and the population size. As used here, sensible alignments require *ca.* 5 CPU seconds on an R10000 Silicon Graphics machine for their identification, meaning that a search of a large database can be quite protracted. We have hence studied a range of techniques that, taken together, can significantly reduce the number of rings that need to be considered in the shape search. These techniques are discussed below.

**Filters for scaffold searching** The first, and most obvious, filter is to screen out those ring systems that cannot possibly support the pattern of substituents specified in the target scaffold, *T*. This can be effected by means of a 3D search in which the geometric pattern is derived from the substituents in *T*, and we have used tools in the SYBYL and UNITY systems [21] for this purpose. A query pattern is generated by removing the attached groups in *T*, replacing these with hydrogen atoms, and recalculating the molecular geometry (for which we use the PM3 forcefield in the SYBYL implementation of MOPAC). All atoms with the exception of the attachment

points and the associated hydrogens are then deleted from  $T$ , and the attachment-point atoms changed to the SYBYL atom type ANY. Distance constraints between each pair of attachment points are defined, and a UNITY 3D search for the resulting query pattern is then carried out using a tolerance (in our experiments) of  $\pm 0.5\text{\AA}$ . The procedure is illustrated in Figure 1, and results in a hit-list containing all of those ring systems with a matching arrangement of attachment points: the hit-list from the geometric search is then submitted for the shape search.

It must be emphasised that the geometric search does not remove the need for the final attachment search, which is specific to the alignment output by the GA. That said, the geometric search is able to filter out many ring systems that cannot fit the target scaffold, as we demonstrate below in Results. However, there may still be a large number of ring systems that need to undergo the shape-based attachment search, and we have hence investigated two further filters that can be employed to reduce the computational requirements of this latter search.

The first filter is extremely simple and involves calculating the molecular volume for each of the ring systems when the ring database is first set up. Then, when a target scaffold is to be searched, the ring systems are ranked in decreasing order of the magnitude of the difference between their molecular volume and that of the target scaffold. The idea here is that rings with very different volumes are unlikely to have a high degree of shape similarity, and can thus be eliminated from further consideration.

The second, and more precise, filter is obtained by taking account of the patterns of inter-atomic distances in the target scaffold and in each of the ring systems in the database that is to be searched. This filter uses a method for distance-based 3D similarity searching called *atom-mapping* [22] that has recently been applied to scaffold searching by Wild and Gifford in their program SAM [23]. Assume that inter-atomic distance matrices are available for the target scaffold,  $T$ , and for a database ring system,  $R$ . Then, as implemented for scaffold searching, a Tanimoto similarity is calculated between each heavy atom in  $T$  and each heavy atom in  $R$ , using the expression

$$\frac{C}{NT + NR - C},$$

where  $C$  is the number of inter-atomic distances in common (using a tolerance of  $\pm 0.5\text{\AA}$ ) between a chosen atom in  $T$  and a chosen atom in  $R$ , and  $NT$  and  $NR$  are the numbers of inter-atomic distances in  $T$  and  $R$ , respectively, involving the chosen atoms. Atoms from  $T$  are then paired with atoms in  $R$  in order of decreasing similarity (thus providing an approximate alignment) and the overall similarity between the two ring systems is the mean of the similarities when averaged over the Tanimoto similarity coefficients for the pairs of matched atoms. SAM is very fast in operation when used for scaffold searching [23]; however, it does not involve any specific measure of the steric overlap of the two rings that are being compared, and it can also yield very confusing alignments in many cases. This is not a problem if the atom-mapping similarities can be shown to correlate strongly with FBSS shape similarities, so that it is used as a filter prior to the full attachment search: the extent to which this occurs in practice is discussed below.

## Results And Discussion

It will be realised from the previous section that our program for scaffold searching contains several different components. Specifically, a search is carried out as shown in Figure 2, and in this section we discuss the results of several scaffold searches that seek to determine the effectiveness of the various steps in the Figure. Our experiments have used a database containing 9040 ring systems extracted from the Chemical Abstracts Service Registry System and containing only CHONS. This database was searched using the three target scaffolds shown in Figure 3, where X denotes the position of an attached group.

The UNITY 3D search retrieved totals of 5132, 4137 and 666 ring systems for targets 1, 2 and 3 respectively, when a tolerance of  $\pm 0.5\text{\AA}$  was allowed for each distance match. Thus, even just a two-substituent scaffold can result in the elimination of over 40% of the database at little computational cost, and in some cases (such as target-3) the geometric search may be all that is required prior to the shape search. In other cases, the molecular volume and/or atom-mapping filters may be required.



The degree of correlation between the molecular volume differences and the FBSS shape similarities is shown in Figure 4. This is for the ring systems remaining after the initial geometric search using target-1; entirely analogous plots are obtained with the other two example scaffolds. It is clear that there is a reasonable correlation between the two sets of values. A similar conclusion may be drawn from Figure 5, which shows the extent of the correlation between the sets of atom-mapping and shape similarity values; this is again for the ring systems remaining after the geometric search for target-1.

The scatter plots in Figures 4 and 5 suggest that both of the filters provide an effective way of post-processing the output from a geometric search: this conclusion is further demonstrated by the figures in Tables 1(a) and 1(b) (for target-1 and target-2 respectively). Each entry in the first column of one of these tables gives the number,  $N$ , of top-ranked ring systems from either the volume-difference filter or the atom-mapping filter, and each entry in the main body of the table gives the numbers (volume-difference first and then atom-mapping second) of those  $N$  top-ranked ring systems that also appeared in the top- $M$  positions in the ranking resulting from the shape search. For example, if we consider the top 2000 structures from the molecular volume search for target-1, then 690 of these ring systems occurred in the top 1000 positions of the ranking based on the FBSS shape similarities (and 1295 and 1709 of these ring systems when the top-2000 or top-3000 positions are considered). It will be clear that many of the top-ranked ring systems from the filter searches will also appear towards the top of the shape searches, and that it is hence reasonable to submit only the upper portion of the filter-search ranking to the time-consuming shape search. It will also be seen that the numbers-in-common are consistently greater for the atom-mapping similarities than for the molecular volume differences, implying that if just one filter is to be applied then the atom-mapping search is the method of choice. That said, it may still be useful to include the volume-difference filter to ensure the elimination of ring systems that are very much larger or very much smaller than the target scaffold.

Once the geometry, filter and shape searches have been carried out, the final stage is the attachment search. This is done by a SYBYL Programming Language (SPL) script that takes the alignments output from the shape search and then checks the top-

ranked ring systems to ensure that they have points of attachment in the same locations as in the target scaffold (we again use a distance tolerance of  $\pm 0.5\text{\AA}$  for a match). The fraction of the database satisfying this search criterion is not generally large. For example, when we carried out an attachment search on the top 500 ring systems from the shape search for target-1, only 31 of the systems matched the arrangement of the substituents in target-1; the corresponding figures for searches of the top 500 ring systems for target-2 and target-3 were 137 and 16, respectively. Target-2 (with three attachment points) gave more hits than target-1 (with just two points) owing to the large number of similarly-shaped 6,6 rings in the search file; more generally, the greater the number of the attachment points, the fewer the number of rings retrieved.

It is worth noting that both this attachment search and the geometric search (as shown in Figure 1) consider only the positions of the points of attachment and disregard the positions of the substituent atoms to which they are attached. The latter information can, of course, be included in a scaffold search, but the resulting increase in precision is often accompanied by a significant reduction in the size of the final output. A user-invoked SPL script is available that addresses this problem,. Specifically, a large dummy atom is substituted at each point of attachment (in both the target scaffold and a matched database ring system) for each superposition identified in the geometric search (there are often several such possible fits). These enhanced rings are then input to the shape search in the normal way. The top-ranked hits will be molecules in which the central scaffolds have a high degree of shape similarity (as previously) and in which the attached dummy atoms are also closely aligned.

We believe that the principal use of scaffold searching is an “ideas generator”, suggesting novel ring systems to a synthetic chemist that might be worth considering in a library design programme. It is our experience that the top-ranked ring systems, typically with FBSS similarities  $\geq 0.9$ , generally provide good alignments of fairly obvious alternatives to the target scaffold, *T*, with more interesting potential scaffolds appearing with FBSS similarities in the range 0.8-0.9; ring systems with still lower similarities are normally, but not consistently, of lesser interest. If only a few ring systems are output from the attachment search, as with target-1 and target-3, then it is

relatively easy to scan through the search output on a graphics terminal. When many ring systems satisfy the search constraints, as with target-2, then some form of hit-list post-processing may be required (Step 6 in Figure 2). Approaches that could be considered include clustering the output using any rapidly-computed similarity measure (such as the atom-mapping similarities or 2D fingerprint similarities) or grouping them using high-level ring descriptors such as those suggested by Bedrosian *et al.* [3], Nilakantan *et al.* [5] or Lipkus [9]; alternatively, a more precise ranking of the search output could be obtained by calculating the volume overlap for each alignment and then ranking the database ring systems in decreasing goodness-of-fit.

One hit-list post-processing approach we have adopted is to make use of the fact that FBSS can calculate not just shape similarities but also electrostatic and hydrophobic similarities or any combination of these three types of field. We found that alignments based on electrostatic or hydrophobic similarities, rather than shape similarities, led to very few matching scaffolds in the final attachment search. However, these other types of similarity can be used to rank the output from the attachment search, so as to collocate rings systems that might be expected to exhibit similar chemistries. Specifically, a search is carried out as detailed in Figure 2 and the hit-list from the attachment search identified; for each ring system in this list, the electrostatic and hydrophobic similarities are calculated given the alignment from the shape search; and finally, the hit-list is ranked in descending order of the sums of the similarities for the three types of field. We have used this approach to identify the best matches for the three target scaffolds shown in Figure 3. Specifically, each part of Figure 6 shows one of these target scaffolds in the top right, with the attachment points marked by purple-coloured dummy atoms, and with the other portions of each figure showing the top three hits based on the sums-of-similarities.

## Conclusions

In this paper we have discussed a range of tools, based on both substructure matching and similarity matching, that can be used for scaffold searching, i.e., for identifying ring systems that are similar to a user-defined target ring system, such as the central scaffold in a combinatorial library definition. Other possible applications of this work

that might be considered include starting templates for structure-based design and scaffold replacement using a previously-established 3D QSAR model.

Although the methods we have discussed seem to provide effective and efficient tools for this task, there are many variations in the precise way in which the overall search is implemented, for example in the ordering of the components of Figure 2 or by the inclusion of conformational flexibility in FBSS's GA [24]. Indeed, there are many other ways in which this sort of functionality could be provided in a library design programme. For example, CAVEAT [25] identifies pairs of ring substituents that are in a specific geometric orientation to each other, and there are many other shape similarity and alignment procedures that could be used [15, 16, 26, 27]. We thus do not claim that the approaches described here are necessarily the best that are currently available for 3D scaffold searching; however, we do believe that they provide a cost-effective way of providing an increasingly important type of search functionality.

**Acknowledgements** We thank Pfizer Global Research and Development and Tripes Inc. for funding. The Krebs Institute for Biomolecular Research is a designated Biomolecular Sciences Centre of the Biotechnology and Biological Sciences Research Council.

## REFERENCES

1. Granito, C.E., Roberts, S. and Gibson, G.W., The conversion of Wiswesser Line Notations to Ring Codes. I. The conversion of ring systems. *J Chem. Docum.* 12, 190-196 (1972).
2. Zamora, A., An algorithm for finding the smallest set of smallest rings. *J. Chem. Inf. Comput. Sci.* 16, 40-43 (1976).
3. Bedrosian, S.D. and Milne, M.B., Graphical representation for automated retrieval of a class of fused six-rings. *J. Chem. Inf. Comput. Sci.* 17, 47-49 (1977).
4. Downs, G.M., Gillet, V.J., Holliday, J.D. and Lynch, M.F., A review of ring perception algorithms for chemical graphs, *J. Chem. Inf. Comput. Sc.*, 29, 172-187 (1989).
5. Nilakantan, R., Bauman, N., Haraki, K.S., Venkataraghavan, R., A ring-based chemical structure query system: use of a novel ring complexity heuristic, *J. Chem. Inf. Comput. Sci.* 30, 1990, 65-68 (1990).
6. Domokos, L., Beilstein Ring Search System. 1. General design. *J. Chem. Inf. Comput. Sci.* 33, 663-667 (1993).
7. Balducci, R. and Pearlman, R.S., Efficient exact solution of the ring perception problem. *J. Chem. Inf. Comput. Sci.* 34, 822-831 (1994).
8. Bemis G.W. and Murcko, M.A., The properties of known drugs. 1. Molecular frameworks., *J. Med. Chem.* 39, 2887-2893 (1996).
9. Lipkus, A.H., Exploring chemical rings in a simple topological-descriptor space. *J. Chem. Inf. Comput. Sci.*, 41, 430-438 (2001).
10. Katritzky, A.R., Kiely, J.S., Hebert, N. and Chassaing, C., Definition of templates within combinatorial libraries. *J. Combin. Chem.* 2, 2-5 (2000).

11. Schneider, G., Neidhart, W., Giller, T. and Schmid, G., "Scaffold-hopping" by topological pharmacophore search: a contribution to virtual screening. *Angew. Chem. Int.* 39, 2894-2896 (1999).
12. Wild, D.J. and Willett, P., Similarity searching in files of three-dimensional chemical structures: alignment of molecular electrostatic potentials with a genetic algorithm. *J. Chem. Inf. Comput. Sci.* 36, 159-167 (1996).
13. Carbó, R., Leyda, L. and Arnau, M., How similar is a molecule to another? An electron density measure of similarity between two molecular structures. *Int. J. Quant. Chem.* 17, 1185-1189 (1980).
14. Good, A.C. and Richards, W.G., Rapid evaluation of shape similarity using Gaussian functions. *J. Chem. Inf. Comput. Sci.* 33, 112-116 (1993).
15. Hahn, M., Three-dimensional shape-based searching of conformationally flexible compounds. *J. Chem. Inf. Comput. Sci.* 37, 1997, 80-86 (1997).
16. Dean, P.M. and Perkins, T.D.J., Calculation of three-dimensional similarity. In: Martin, Y.C. & Willett, P. (editors). *Designing Bioactive Molecules: Three-Dimensional Techniques and Applications*. American Chemical Society, Washington 1998, pp 199-218.
17. Mestres, J., Rohrer, D.C. and Maggiora, G.M., A molecular-field-based similarity study of non-nucleoside HIV-1 reverse transcriptase inhibitors. *J. Comput.-Aid. Mol. Design* 13, 79-93 (1999).
18. Cramer, R.D., Patterson, D.E., Clark, R.D., Soltanshani, F. and Lawless, M.S., Virtual compound libraries: a new approach to decision making in molecular discovery research. *J. Chem. Inf. Comput. Sci.* 38, 1010-1023 (1998).
19. Gillet, V.J., Schuffenhauer, A. and Willett, P., Similarity searching in files of 3D chemical structures: analysis of the BIOSTER database using 2D fingerprints and molecular field descriptors. *J. Chem. Inf. Comput. Sci.* 40, 295-307 (2000).
20. Jewell, N.E., Turner, D.B, Willett, P. and Sexton, G.J., Automatic generation of alignments for 3D QSAR analyses. *J. Mol. Graph. Model.*, in the press.
21. SYBYL and UNITY are available from Tripos Inc. at URL <http://www.tripos.com>
22. Pepperrell, C.A., Taylor, R. and Willett, P., Implementation and use of an atom-mapping procedure for similarity searching in databases of 3-D chemical structures. *Tetrahed. Comput. Methodol.* 3, 575-593 (1990).
23. Wild, D.J. and Gifford, E.M. "Hunting for scaffolds." Presented at the Daylight MUG 2000 conference, February 22-25 2000 in Sanata Fe, NM. At URL <http://www.daylight.com/meetings/mug2000/Wild/Mug2000.html>
24. Thorner, D.A., Wild, D.J., Willett, P. and Wright, P.M., Similarity searching in files of three-dimensional chemical structures: flexible field-based searching of molecular electrostatic potentials. *J. Chem. Inf. Comput. Sci.*, 36, 900-908 (1996).
25. Lauri, G., Bartlett, P.A., CAVEAT: a program to facilitate the design of organic molecules. *J. Comput.-Aid. Mol. Design* 8, 51-66 (1994).
26. Van Geerestein, V., Perry, N., Grootenhuis, P., Haasnoot, C. 3D database searching on the basis of ligand shape using the SPERM prototype method. *Tetrahed. Comput. Methodol.* 3, 595-613 (1990).
27. Lemmen, C. and Lengauer, T., Computational methods for the structural alignment of molecules. *J. Comput.-Aid. Mol. Design* 14, 215-232 (2000)

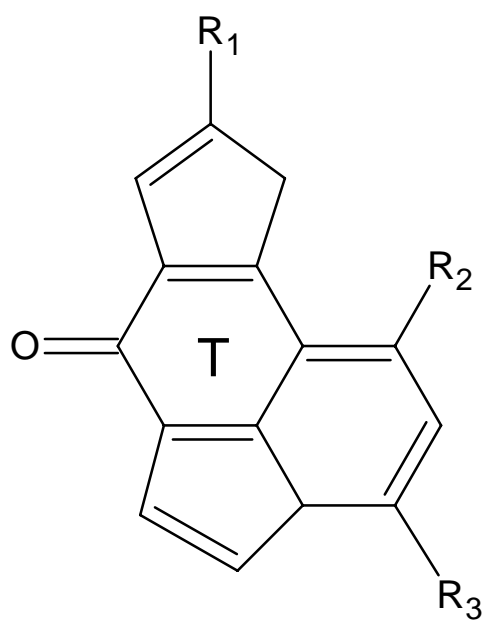
Top- $N$ scaffolds	Number of scaffolds in common with the top- $M$ in the shape search rankings		
	$M=1000$	$M=2000$	$M=3000$
1000	324 687	626 976	792 1000
2000	690 945	1295 1747	1709 1996
3000	902 999	1771 1992	2575 2844

1(a)

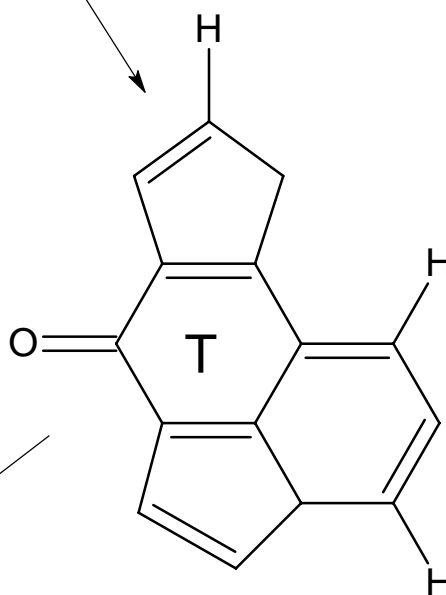
Top- $N$ scaffolds	Number of scaffolds in common with the top- $M$ in the shape search rankings		
	$M=1000$	$M=2000$	$M=3000$
1000	286 509	409 789	507 966
2000	674 813	1060 1456	1341 1890
3000	846 968	1663 1893	2232 2662

1(b)

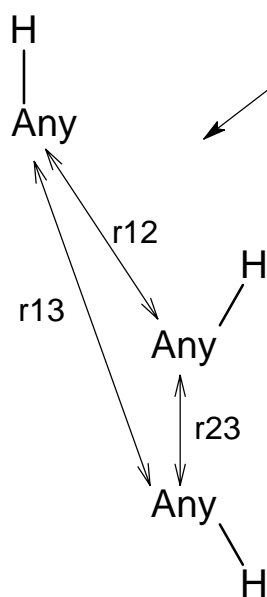
**Table 1.** Effectiveness of the molecular volume and atom-mapping filters when applied to the outputs from the initial geometric search for (a) target-1 and (b) target-2. Each entry in the first column gives the number,  $N$ , of top-ranked ring systems from either the volume-difference filter or the atom-mapping filter, and each entry in the main body of the table gives the numbers (volume-difference first and then atom-mapping second) of those  $N$  top-ranked ring systems that also appeared in the top- $M$  positions in the ranking resulting from the shape search.



Replace substituents with hydrogen atoms and re-optimize the 3D structure



Delete all atoms except those representing the attachment points. Change ring atoms to atom type Any



Set distance constraints r<sub>12</sub>, r<sub>13</sub> and r<sub>23</sub> (typically with a tolerance of  $\pm 0.5 \text{ \AA}$ )

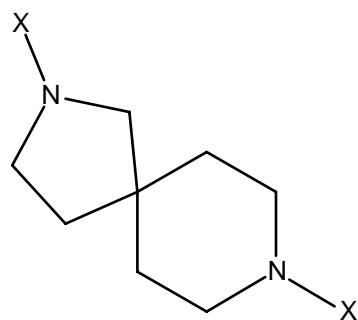
**Figure 1.** Generating a 3D query pattern for use in a scaffold search.

1. A query template is input, this consisting of a central ring scaffold,  $T$ , and the substituent positions at which functionality can be attached.
2. The template  $T$  is processed as shown in Figure 1, so that it can form the basis for a UNITY 3D search
3. If there is a large hit-list from Step 2 then a volume-difference search and/or an atom-mapping search are/is carried out to find the ring systems that are most similar to  $T$ .
4.  $T$  is used as the target for a shape search of a (possibly filtered) database of ring systems, and the resulting similarities ranked in descending order.
5. The attachment search is carried out on the top-ranked ring systems from Step 4.
6. Carry out any final post-processing steps

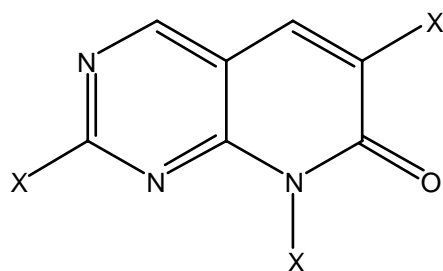
**Figure 2.** Principal components of a system for scaffold searching



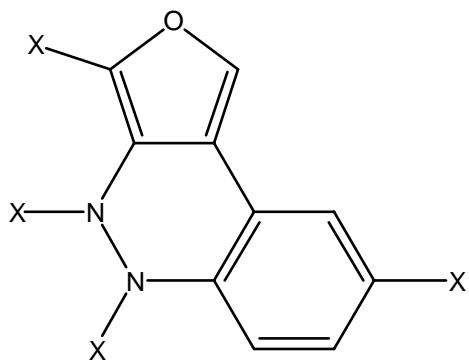
**Target 1**



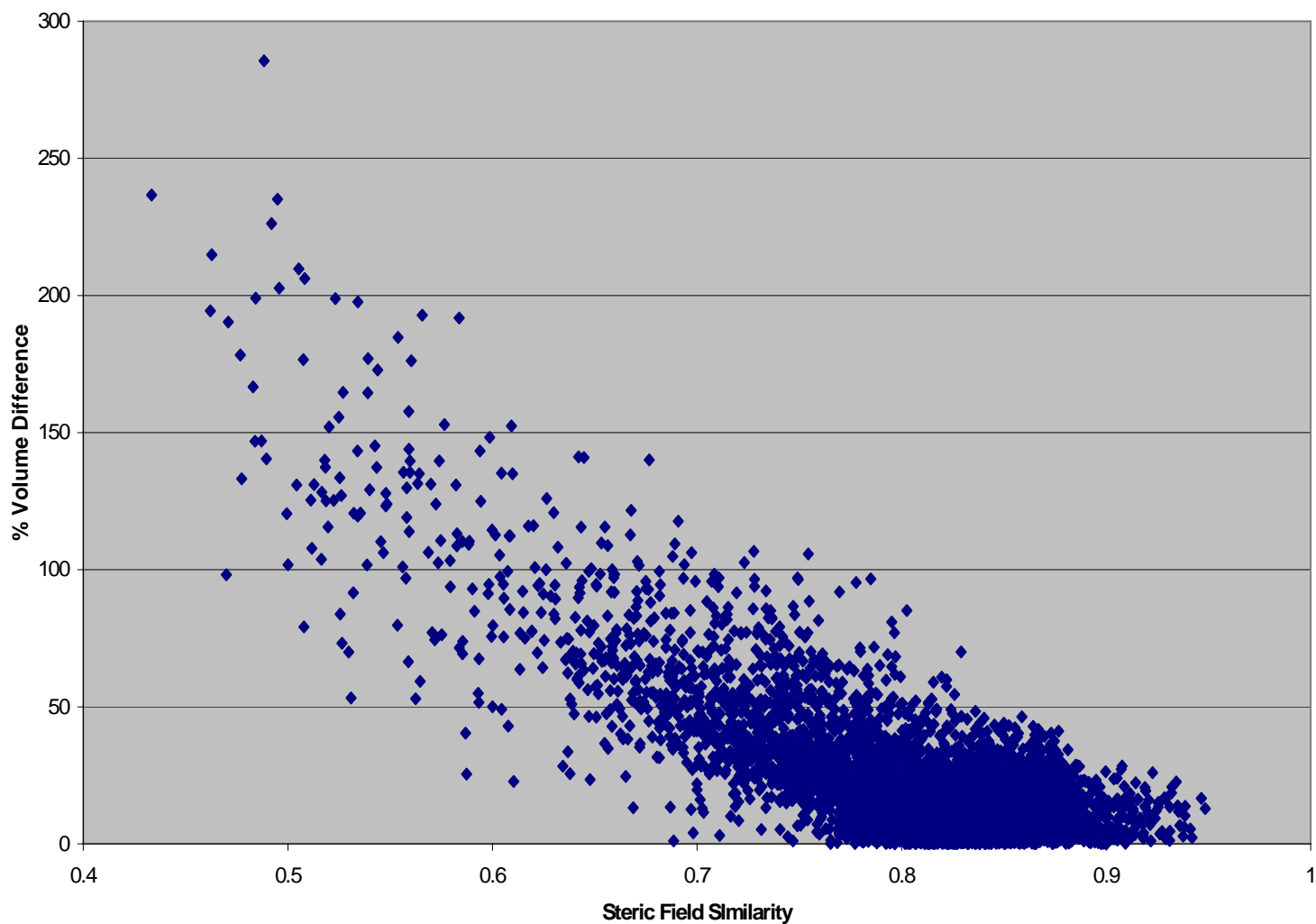
**Target 2**



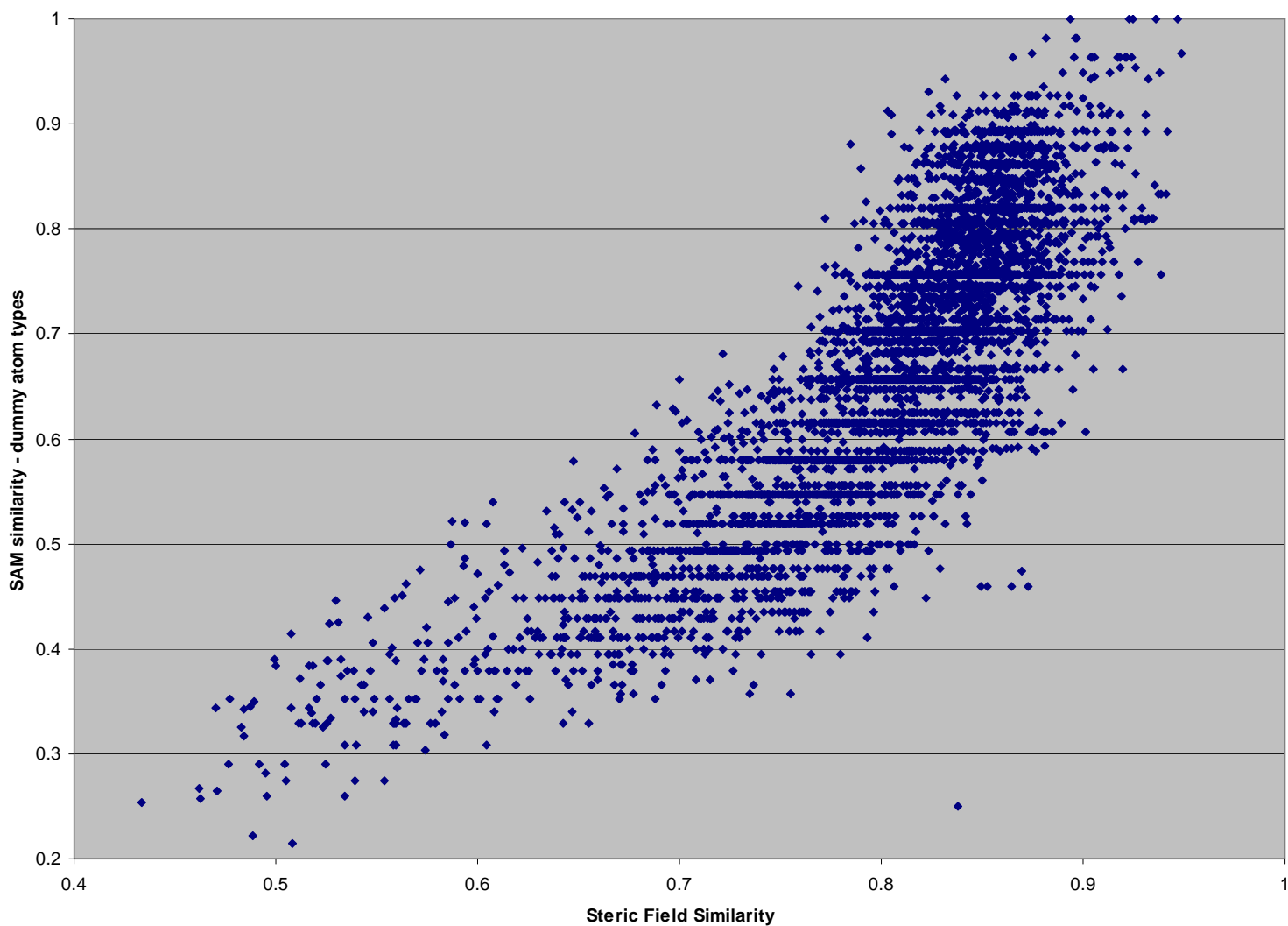
**Target 3**



**Figure 3.** The three target scaffolds used in our experiments, with the symbol 'X' denoting a point of attachment.



**Figure 4.** Scatter plot describing the effectiveness of filtering based on molecular volumes. The shape similarity was calculated for each of the ring systems passing the geometric search in a scaffold search for target-1, as was the percentage difference in the molecular volumes between target-1 and each of the selected ring systems.



**Figure 5.** Scatter plot describing the effectiveness of filtering based on atom-mapping. The shape similarity was calculated for each of the ring systems passing the geometric search in a scaffold search for target-1, as was the atom-mapping similarity for each of the selected ring systems.

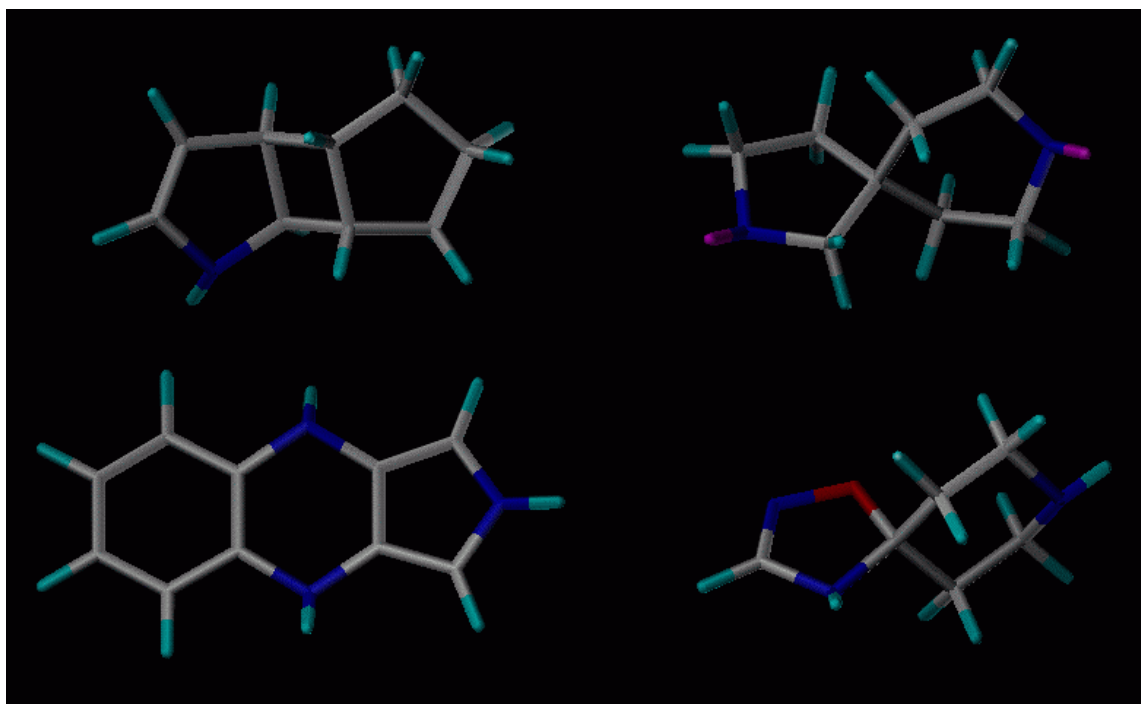


Figure 6a

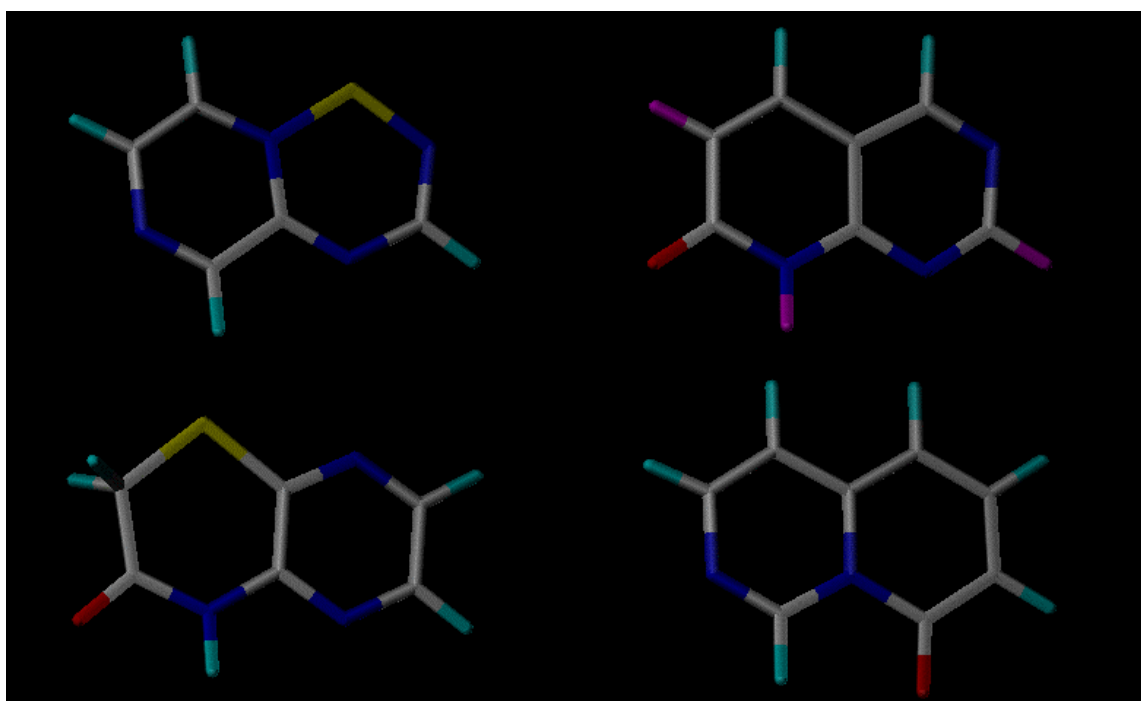
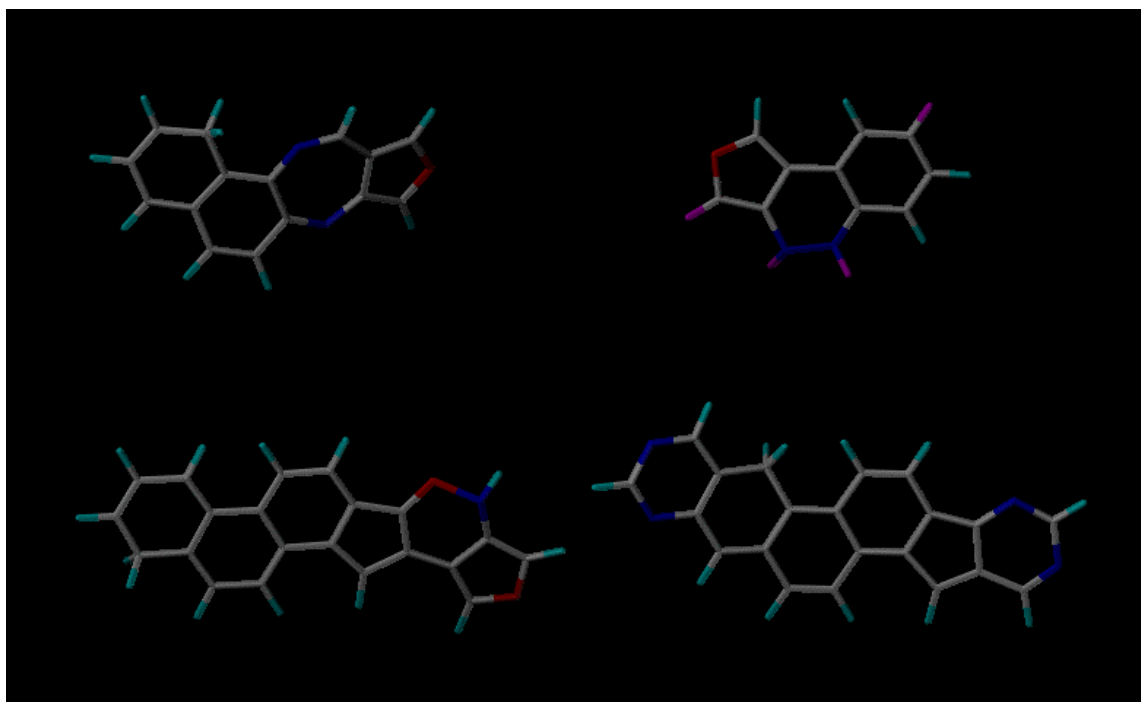


Figure 6b



**Figure 6c**

**Figure 6.** Examples of search output for (a) target-1 (b) target-2 and (c) target-3. In each case, the target scaffold is positioned in the top right of the figure (with the attachment points marked by purple-coloured dummy atoms) and with the other portions of the figure showing the top three hits based on the sums of the shape, electrostatic and hydrophobic similarities.