

*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 **Organic & Biomolecular Chemistry**.

White Rose Research Online URL for this paper:  
<http://eprints.whiterose.ac.uk/77602>

---

#### **Published paper**

Hert, J, Willett, P, Wilton, DJ, Acklin, P, Azzaoui, K, Jacoby, E and Schuffenhauer, A (2004) *Comparison of topological descriptors for similarity-based virtual screening using multiple bioactive reference structures*. *Organic & Biomolecular Chemistry*, 2 (22). 3256 - 3266.

---

#### AUTHOR FOR CORRESPONDENCE

Prof. Peter Willett, Department of Information Studies, University of Sheffield, Western Bank, Sheffield S10 2TN, UK

#### TITLE

Comparison of topological descriptors for similarity-based virtual screening using multiple bioactive reference structures

#### AUTHORS

Jérôme Hert, Peter Willett (\*) and David J. Wilton (Krebs Institute for Biomolecular Research and Department of Information Studies, University of Sheffield, Western Bank, Sheffield S10 2TN, UK)

Pierre Acklin, Kamal Azzaoui, Edgar Jacoby and Ansgar Schuffenhauer (Novartis Institutes for BioMedical Research, Discovery Technologies, CH-4002 Basel, Switzerland).

#### ABSTRACT

This paper reports a detailed comparison of a range of different types of 2D fingerprints when used for similarity-based virtual screening with multiple reference structures. Experiments with the MDL Drug Data Report database demonstrate the effectiveness of fingerprints that encode circular substructure descriptors generated using the Morgan algorithm. These fingerprints are notably more effective than fingerprints based on a fragment dictionary, on hashing and on topological pharmacophores. The combination of these fingerprints with data fusion based on similarity scores provides both an effective and an efficient approach to virtual screening in lead-discovery programmes.

## INTRODUCTION

Virtual screening is widely used to enhance the cost-effectiveness of drug-discovery programmes, by ranking databases of chemical structures in decreasing probability of activity; this prioritisation then means that biological testing can be focussed on just those few molecules that have significant *a priori* probabilities of activity [1, 2]. There are many different ways in which a database can be prioritised; here, we focus on similarity searching methods [3, 4]. Similarity searching is one of the most widely used virtual-screening approaches, and involves matching a known active molecule, the *reference structure*, against each of the database molecules, computing a measure of structural similarity in each case, and then ranking the database molecules in order of decreasing similarity score. Structurally similar molecules are likely to have similar biological activities [5-7] and there is hence an extensive literature associated with the similarity measures that can be used to quantify the degree of resemblance between a reference structure and each of the database molecules.

Most of the studies of similarity searching that have been reported thus far have considered the use of only a single bioactive reference structure. It is, however, increasingly the case that several, structurally diverse, reference structures may be available, e.g., published competitor compounds or hits from high-throughput screening (HTS), and this has stimulated interest in the use of multiple reference structures to identify further molecules for biological screening [8, 9]. We have recently reported a detailed comparison of several different search algorithms that can be used in such circumstances, and identified two, data fusion and binary kernel discrimination (BKD), that provided a high level of effectiveness in simulated virtual screening experiments [10].

An important component of any similarity procedure is the structure representation that is used to encode the molecules that are to be searched, with 2D fragment bit-strings (or fingerprints) of various types being by far the most commonly used in current cheminformatics systems [11, 12]. A fingerprint is a binary string that encodes the presence of substructural fragments, i.e. topological patterns of atoms and bonds, in a molecule. This is clearly a very simple representation of molecular structure but one that has been used with considerable success ever since the earliest reports of similarity searching [13, 14], and also for related cheminformatics tasks such as molecular diversity analysis [15] and database clustering [16]. In two much-cited papers, Brown and Martin compared several different types of fingerprint when used for cluster-based physicochemical property prediction [17, 18]; here, we report an analogous comparison of fingerprints when they are used for similarity-based virtual screening using multiple reference structures.

## RESULTS AND DISCUSSION

Taking account of the different search algorithms, fingerprint types and normalisation schemes described in the EXPERIMENTAL section, there is a total of 30 different similarity procedures available for evaluation. Each such procedure was used with each of 11 different activity classes from the MDL Drug Data Report (MDDR) [19] database, with ten searches being carried out for the actives in each

particular activity class. A different set of ten active reference structures was used for each of the ten searches, this set remaining constant across the 30 different similarity procedures. The results of the searches are shown in Tables 1-4, the first two listing the average recall obtained from the top-1% of the rankings for each of the activity classes, and the second two listing the average recall from the top-5% of the rankings. The mean recalls, averaged over all of the 11 activity classes, are shown in Figures 1 and 2 (which also show the mean recall for data fusion and BKD averaged over all of the different fingerprint types considered here).

Inspection of these tables reveals the very marked superiority of the circular substructure descriptors; indeed, there was only a single case where one of these fingerprints did not provide the best result, *viz* the average recall at 1% using BKD for the set of cyclooxygenase inhibitors. This general effectiveness of the circular substructures (with the notable exception of FCFP\_2) is highlighted in Figures 1 and 2. Of these circular substructure fingerprints, the ECFP\_4 ones, irrespective of the normalisation method (method-A or method-B) or of the search algorithm (data fusion or BKD) are the descriptors of choice for virtual screening of the sort advocated here. The FCFP\_4 and ECFP\_2 descriptors are also very effective: the former fingerprints seem to perform relatively better with the more heterogeneous (i.e., low self-similarity) classes, such as the cyclooxygenase and protein kinase C inhibitors, while the ECFP\_2 fingerprints yield better results with the more homogeneous (i.e., high self-similarity) classes, such as the renin inhibitors.

As an alternative way of considering the figures in Tables 1-4, consider the enrichment factors [20] to which these results correspond. The enrichment factor is the number of times better (in terms of active molecules retrieved) that a particular search algorithm is than a random selection of molecules from the database. Thus, the average enrichment values for ECFP\_4B at 1% are 42.3 and 43.5 for BKD and data fusion, respectively, with the corresponding 5% values being 13.1 and 13.5, respectively, demonstrating the utility of the methods discussed here for virtual screening purposes.

Circular substructures of various sorts have been widely used for applications such as structure and substructure searching [21-23], constitutional symmetry [24], structure elucidation [25] and, most recently, probabilistic modelling of bioactivity where a full training-set is available [26, 27]. The work reported here demonstrates that this type of fragment is also very well suited to virtual screening using multiple reference structures.

When comparing the normalisation methods used for the circular substructure representations (see EXPERIMENTAL section), method A, where all the initial features are just assigned a new bit-position, always provides descriptors that are more effective than method B. However the differences are generally very small, and we would hence recommend the use of method-B for the processing of these descriptors as this method is faster and, more importantly, is reproducible over different databases. There is little to choose between data fusion and BKD over the entire class of circular substructures, although it does appear that the use of these substructures with data fusion was particularly successful for the more heterogeneous classes like the cyclooxygenase and protein kinase C inhibitors. Conversely, when these descriptors were used with BKD, they worked particularly

well for the more homogeneous activity classes, such as the renin inhibitors and the angiotensin II AT1 antagonists.

The dictionary-based descriptors, represented here by the BCI fingerprints, were ranked second overall, returning generally higher recalls than the hashed fingerprints, i.e., Unity, Daylight and Avalon. This finding is in agreement with the studies of cluster-based property prediction by Brown and Martin [17, 18] (although they used different types of dictionary and hashed fingerprints from those studied here).

Perhaps our most surprising finding is the performance of the pharmacophore descriptors, with both the CATS and Similog fingerprints yielding consistently poorer recall values. Previous studies of these descriptors, for chemogenomics and scaffold-hopping applications [9, 28, 29], have demonstrated that they can be highly effective in operation, but this was certainly not the case for the present application. We note in the EXPERIMENTAL section that both of these molecular characterisations are based on the encoding of the occurrences, rather than the incidences, of substructural fragments in a molecule, yielding an integer vector rather than a binary fingerprint. Here, however, they have been encoded in a binary form since the kernel function used in our BKD implementation requires a binary string. It is hence possible that the poor performance of the two pharmacophore fingerprints arose from the use of an inappropriate encoding mechanism. To test this, searches were carried out with the original, occurrence-based vectors; these searches used just data fusion, as this search algorithm does not necessarily require binary fingerprints for the generation of the ranked sets of scores that are fused together. Specifically, the rankings for the individual reference structures were computed using the non-binary form of the Tanimoto similarity coefficient and the Floersheim distance, as defined in EXPERIMENTAL below. The use of the integer vectors and the non-binary coefficients did not improve the recall of the data fusion searches, and we hence conclude that the use of binary representations does not explain the poor performance of these 2D pharmacophore descriptors that is observed in Tables 1-4. It is perhaps worth noting in passing that previous comparisons of 2D fingerprints with 3D pharmacophore descriptors have often shown the former to be superior [17, 18, 30], despite the claimed effectiveness of the latter methods for diversity analysis and similarity searching [31].

Thus far, we have evaluated the various approaches solely in terms of the numbers of active molecules that have been retrieved. It is, however, also of importance to consider the diversity of these sets of retrieved actives, since it is clearly preferable for the outputs also to maximise the numbers of chemotypes that are identified. We have hence analysed the outputs summarised in Tables 1-4 and Figures 1 and 2 in terms of the numbers of distinct ring systems identified in the sets of retrieved actives. We have considered two levels of ring description, as illustrated in Figure 3, and as discussed previously by Bemis and Murcko [32] and by Xu and Johnson [33, 34]; these authors refer to these levels of description as molecular frameworks or cyclic systems (Figure 3a), and frameworks or skeletal cyclic systems (Figure 3b), respectively. Figure 4 shows the percentages of the molecular frameworks in the complete set of actives that are retrieved in the top-1% of the ranking by each of the search procedures when averaged over all of the activity classes (i.e., as in Figure 1); Figure 5 gives the top-1% distribution for the frameworks and Figures 6

and 7 the corresponding top-5% distributions. It will be seen that the relative performance of the various procedures in terms of retrieving chemotypes (and hence in their suitability for scaffold-hopping applications) mirrors closely the relative performance based on numbers of actives (as shown in Figures 1 and 2).

All the experiments carried out so far were performed using a version of the MDDR database in which every molecule was characterised by its neutral structure. However, drugs are used *in vivo* and further searches were hence carried out in order to see if any improvements in recall could be obtained by using the protonated states of the MDDR molecules. The pH component of Scitegic's Pipeline Pilot software [35] was used to derive protonated molecular representations corresponding to a pH 6.8. However, very little difference was observed in the recalls obtained from the compounds in their protonated and neutral forms, with the latter normally being the more effective. There would hence appear to be little point in carrying out the additional processing required to produce the protonated representations.

The results presented here provide further evidence of the general effectiveness of the BKD and data fusion methods for virtual screening applications where multiple reference structures are available, and evidence of the general effectiveness of fingerprints based on 2D circular substructures, in particular the ECFP\_4 fingerprints. If a single choice is required, then the best overall performance would seem to result from data fusion of the similarity scores of searches based on the ECFP\_4B fingerprints. This is indicated as the combination of choice for several reasons. If we consider the choice of fingerprint first, then whilst the ECFP\_4 descriptors achieved an excellent overall level of performance, they gave particularly good results when searching for structurally heterogeneous sets of molecules, a more challenging task than for highly self-similar sets of molecules. For this descriptor, the type-B binning scheme results in a very compact, reproducible representation that is only marginally inferior to the much larger, non-reproducible type-A binning scheme. Turning now to the search algorithms: data fusion is far less demanding of computational resources than is BKD and also does not require the specification of values for the latter's tuneable parameters; and an inspection of the standard deviations in Tables 1-4 shows that these tend to be larger (corresponding to a high level of variation in search performance) for BKD than for data fusion, suggesting a greater degree of consistency for the latter algorithm.

When considering the two search algorithms, it must be emphasised that we are dealing with a combination of characteristics, as evidenced by the fact that BKD does better than data fusion for some of the fingerprint types (e.g., Unity or Daylight): however, when used in combination with ECFP\_4, the data fusion searches are to be preferred. It should also be emphasised that this preference for score-based data fusion over BKD is specific to the circumstances of these experiments, which involve just a limited number of active reference structures, as we have found that BKD is to be preferred when a proper training-set is available containing large numbers of both known actives and known inactives [36].

## EXPERIMENTAL

Our virtual screening system involves three main parts: a structural representation

that is used to encode the molecules that are being searched; a searching method that ranks a database of molecules in order of decreasing probability of activity in response to a set of active reference structures; and a quantitative measure of the effectiveness of those rankings. The focus of this paper is the first of these factors, but it is appropriate to describe briefly the last two factors before discussing the many different types of 2D fingerprint that were evaluated.

### Searching algorithms

A previous study [10] investigated a range of search algorithms that could be used when multiple reference structures were available. These experiments all involved a single type of fingerprint, specifically the Unity fingerprints produced by Tripos Inc. [37]. Two of the algorithms were found to provide a consistently high level of screening effectiveness: these algorithms were data fusion using the maximum of similarity scores and an approximate form of the BKD machine learning technique.

Data fusion (or consensus scoring) involves combining the results of different similarity searches of a chemical database. Previous studies have involved the use of a single reference molecule, but characterised by several different representations or using several different similarity coefficients (see, e.g., [38, 39]). An alternative approach, and the one used here, is to have a fixed representation and similarity coefficient, but to combine the search outputs obtained with several different reference structures. Assume that some database molecule  $i$  yields similarity scores of  $s_1, s_2..s_n$  with the  $n$  different reference structures, then we have shown that effective searches are obtained by ranking the database molecules on the basis of the maximum of these scores, i.e.,  $\max\{s_1, s_2..s_i..s_{n-1}, s_n\}$ ; such searches are more effective than those resulting from the use of ranks, rather than scores, or the use of a fusion rule based on averaging [9, 10].

The similarity scores were computed using the Tanimoto coefficient; for a molecule having a fingerprints with  $a$  bits set, of which  $c$  are also set in the fingerprint for a molecule that has  $b$  bits set, then the Tanimoto coefficient,  $T_c$ , is defined to be

$$T_c = \frac{c}{a + b - c}.$$

Some of the similarity scores necessitated the use of two non-binary similarity coefficients. Let  $x_{jA}$  denote the occurrence of the  $j$ -th fragment (1 ≤  $j$  ≤  $n$ , the length of the integer vector) in molecule  $A$  (and similarly for molecule  $B$ ). Then the similarity coefficients used were the non-binary form of the Tanimoto similarity coefficient [3]

$$S_{A,B} = \frac{\mathring{a} \sum_{j=1}^{j=n} (x_{jA} x_{jB})}{\mathring{a} \sum_{j=1}^{j=n} (x_{jA})^2 + \mathring{a} \sum_{j=1}^{j=n} (x_{jB})^2 - \mathring{a} \sum_{j=1}^{j=n} (x_{jA} x_{jB})},$$

and the Floersheim distance

$$D_{A,B} = \frac{\mathring{a} \sum_{j=1}^{j=n} \frac{|x_{jA} - x_{jB}|}{1 + \min(x_{jA}, x_{jB})}}{\mathring{a} \sum_{j=1}^{j=n} x_{jA} + \mathring{a} \sum_{j=1}^{j=n} x_{jB}},$$

which is a Novartis coefficient that has been used in-house with the Similog

descriptors.

Binary kernel discrimination (BKD) is a machine learning technique that was first applied to virtual screening by Harper *et al.* [40]. The similarity between two compounds  $i$  and  $j$ , characterised by binary fingerprints of length  $M$ , that differ in  $d_{ij}$  positions, is computed by the kernel function  $K_1(i, j)$ ,

$$K_1(i, j) = \left( |^{M-d_{ij}} (1 - |)^{d_{ij}} \right)^{k/M}$$

where  $|$  is a smoothing parameter to be determined and where  $k$  is an integer less than  $M$ . This kernel was developed for use with a training-set containing both active and inactive molecules, with the scoring function

$$L_1(j) = \frac{\sum_{i \in \text{actives}} K_1(i, j)}{\sum_{i \in \text{inactives}} K_1(i, j)}$$

being used to rank the molecules in the test-set, using the optimum values of  $|$  and  $k$  found for the training set [36]. When just a set of active reference structures is available, we have shown that the characteristics of the inactives can be approximated with a fair degree of accuracy by the characteristics of the entire database that is to be searched: a training-set can hence be generated by taking the set of reference structures and adding to it 100 molecules randomly selected from the database [10]. The optimal values of  $|$  and  $k$  were found to vary across the various types of fingerprint studied here, and extensive preliminary testing was required to identify the values that were used to obtain the main results that are discussed in RESULTS AND DISCUSSION. This variation in parameter value is a clear limitation of the BKD approach

#### *Effectiveness of virtual screening*

The experiments involved simulated virtual screening searches on the MDL Drug Data Report (MDDR) database [19]. After removal of duplicates and molecules that could not be processed using local software, a total of 102,535 molecules was available for searching that were represented by each of the types of fingerprint described below. These molecules were searched using the eleven sets of active compounds that are listed in Table 5, which also contains the numbers of actives in each class and the numbers of active assemblies and frameworks (ring-system descriptors that are discussed in RESULTS) in each class. The table also contains a numeric estimate of the level of structural diversity in each of the chosen sets of bioactives. The diversity estimate was obtained by matching each compound with every other in its activity class, calculating similarities using the Unity fingerprint and Tanimoto coefficient and computing the mean of these intra-set similarities. The resulting similarity scores are listed in the right-hand part of Table 5, where it will be seen that the renin inhibitors are the most homogeneous and the cyclooxygenase inhibitors are the most heterogeneous.

For each of the 11 activity classes, ten active compounds were selected for use as reference structures. The selections were done at random, subject to the constraint that no pair-wise similarity in a group exceeded 0.80 (using Unity fingerprints and the Tanimoto coefficient). The set of reference structures was searched against the MDDR database using the data fusion and BKD search algorithms described above, with the search being carried out once for each of the different types of fingerprint. The procedure was then repeated using ten different sets of reference structures, and



in each search, a note was made of the *recall*, that is the percentage of the active molecules (i.e., those in the same class as those in the reference set) that occurred in the top-1% and the top-5% of the ranking resulting from that search. Formally, if a search retrieves the top- $x\%$  of a ranked database, and that this subset contains  $a$  of the  $A$  actives for that activity class, then the recall,  $R_x$ , is defined to be [20]

$$R_x = \frac{100 \cdot a}{A}.$$

The results presented in Tables 1-4 are the mean and standard deviations for these recall values, averaged over each set of ten searches.

### *Fingerprint types*

Having summarised the virtual screening environment and the two search algorithms, we now describe in some detail the four classes of fingerprint descriptors that we have studied. These are structural keys, hashed fingerprints, circular substructures and pharmacophores; in all, we evaluated ten different descriptors, of which seven are commercially available, two are used in-house at Novartis, and one was implemented from the literature description. Moreover, some of these descriptors were encoded in more than one way, to give a total of 15 types of fingerprint.

Structural keys have been used in chemoinformatics for many years, and are usually encoded by a binary array, each element of which denotes the presence or absence of a specific 2D fragment. A predefined fragment dictionary lists the various fragment substructures that are encoded in the fingerprint. This study used the 1052-bit Barnard Chemical Information (BCI) fingerprints, which encode the following types of fragment substructure: augmented atoms, atom sequences, atom pairs, ring composition and ring fusion substructural fragments [41].

Hashed fingerprints differ from structural keys in that they do not use a predefined dictionary. Instead, patterns are encoded in the fingerprint, where a pattern describes, for example, a path of length  $n$  bonds, i.e., (atom-bond-atom) $_n$  with the natures of the atoms and bonds defined. The set of patterns produced from any molecule of non-trivial size is obviously very large and differs from molecule to molecule. It is hence not possible to assign each potential pattern to a specific bit position in a fingerprint of predefined length; instead, the pattern is passed to a hashing function to generate a position (or positions) within the available length of the bit-string. The study used three different hashed fingerprints: 2048-bit Daylight fingerprints [42], 988-bit Unity fingerprints [37] and 2048-bit Avalon fingerprints. Daylight fingerprints encode each atom's type, all augmented atoms and all paths of length 2-7 atoms. Unity fingerprints encode paths of length 2-6 atoms, and also include 60 structural keys for common atoms and ring counts. Avalon fingerprints are used for similarity search in Novartis' corporate data warehouse and encode atoms, augmented atoms, atom triplets and connection paths.

A circular substructure is a fragment descriptor where each atom is represented by a string of extended connectivity values that are calculated using a modification of the Morgan Algorithm [43]. The study evaluated two different circular substructure descriptors from Scitegic's Pipeline Pilot Software [35]: Extended Connectivity Fingerprints (ECFPs) and Functional Connectivity Fingerprints (FCFPs). The initial code assigned to an atom is based on the number of connections, the element type,

the charge, and the mass for ECFPs and on six generalised atom-types - viz., hydrogen bond donor, hydrogen bond acceptor, positively ionisable, negatively ionisable, aromatic and halogen - for FCFPs. This code, in combination with the bond information and with the codes of its immediate neighbour atoms is hashed to produce the next order code, which is mapped into an address space of size  $2^{32}$ , and the process iterated until the required level of description has been achieved. The experiments here used the ECFP\_2, ECFP\_4, FCFP\_2 and FCFP\_4 fingerprints, where the numeric code denotes the diameter in bonds up to which features are generated.

The Scitegic software represents a molecule by list of integers, each describing a molecular feature and each in the range  $-2^{31}$  to  $2^{31}$ . These integer lists were normalised in two ways, referred to as method-A and method-B. In method-A, all the features present in the database were enumerated, so that each feature was given as its new code its rank in the sorted list of codes, with the length of the resulting fingerprints being the number of distinct features in the database. In method-B, the integers describing a molecule were hashed to a bit-string of length 1024 bits. This inevitably means that collisions occur, with the result that method-B loses some of the structural information that is retained by method-A; however, the latter representation is dependent on the precise database that is being processed.

Pharmacophore points are features (such as a heteroatom or the centre of an aromatic ring) that are thought to be required for a molecule to show bioactivity. Pharmacophore fingerprinting involves generating all of the patterns of three or four pharmacophore points in a molecule, together with the corresponding inter-point distances, and then using the resulting 3D structural codes as descriptors for similarity searching or diversity analysis (see, e.g., [17, 18, 30, 31]). When used with 2D, rather than 3D, structural representations, the inter-atomic distances can be replaced by through-bond distances, and this approach forms the basis of the two pharmacophore fingerprints studied here: Similog keys [9] and the Chemically Advanced Template Search (CATS) descriptor [28], both of which are based on generalised atom-types describing potential pharmacophores.

The Similog keys use a “DABE” atom-typing scheme based on the following four properties: hydrogen bond donor (D), hydrogen bond acceptor (A), bulkiness (B) and electropositivity (E). The presence or absence of these properties for an atom is encoded in a 4-bit string, and each triplet of atoms is represented by the three DABE strings and by the associated topological distances: in all, 8031 different codes were identified in the MDDR database. The Similog keys store the occurrence of each distinct code, and not just their presence or absence as in a conventional bit-string. A binning scheme was hence used to bin the occurrence data into 8-bit strings: the two binning schemes used (called Method-A and Method-B) are shown in Table 2

The CATS descriptor is based on counts of atom-pair topological distances, with the following generalised types of atom being considered in the generation of the descriptor: lipophilic, positive, negative, hydrogen-bond donor and hydrogen-bond acceptor. The occurrences of the 15 possible pairs of pharmacophores are determined for distances up to 10 bonds to give a 150-element (i.e.,  $15 \times 10$ ) vector. The vectors were generated using the description in Fechner *et al.* [29] and then converted to a binary fingerprint using Method-B in Table 6 (we only used Method-

B for CATS as a substantial fraction of the keys occurred more than seven times in a molecule).

Table 7 lists the abbreviated names that are used in the paper for each of the 15 types of fingerprints, where the A and B subscripts denote the type of normalisation scheme used for binning in the case of the ECFP, FCFP and Similog descriptors. The table also details statistical characteristics of each of these fingerprints: an inspection of the average numbers of bits and the densities (i.e., the mean number of bits that are set divided by the bit-string length and then expressed as a percentage) show a very wide range of levels of molecular description.

## ACKNOWLEDGEMENTS

We thank the following: Novartis Institutes for Biomedical Research for funding; MDL Information Systems Inc. for the provision of the MDDR database; and Barnard Chemical Information Ltd., Daylight Chemical Information Systems Inc., the Royal Society, Tripos Inc. and the Wolfson Foundation for software and laboratory support; and Drs Bernd Rohde and Philippe Floersheim for helpful discussions and for access to the Avalon and Similog descriptors, respectively. The Krebs Institute for Biomolecular Research is a designated biomolecular sciences centre of the Biotechnology and Biological Sciences Research Council.

## REFERENCES

1. H.-J. Böhm, and G. Schneider, eds., *Virtual Screening for Bioactive Molecules*, Wiley-VCH, Weinheim, 2000.
2. G. Klebe, ed. *Virtual Screening: an Alternative or Complement to High Throughput Screening*, Kluwer, Dordrecht, 2000.
3. P. Willett, J. M. Barnard and G. M. Downs, *J. Chem. Inf. Comput. Sci.*, 1998, **38**, 983-996
4. R. P. Sheridan and S. K. Kearsley, *Drug Discov. Today*, 2002, **7**, 903-911
5. M. A. Johnson and G. M. Maggiora, eds. *Concepts and Applications of Molecular Similarity*, Wiley, New York, 1990
6. D. E. Patterson, R. D. Cramer, A. M. Ferguson, R. D. Clark, and L. E. Weinberger, *J. Med. Chem.*, 1996, **39**, 3049-3059
7. Y. C. Martin, J. L. Kofron and L. M. Traphagen, *J. Med. Chem.*, 2002, **45**, 4350-4358
8. L. Xue, F. L. Stahura, J. W. Godden and J. Bajorath, *J. Chem. Inf. Comput. Sci.*, 2001, **41**, 746-753
9. A. Schuffenhauer, P. Floersheim, P. Acklin and E. Jacoby, *J. Chem. Inf. Comput. Sci.*, 2003, **43**, 391-405
10. J. Hert, P. Willett, D. J. Wilton, P. Acklin, K. Azzaoui, E. Jacoby and A. Schuffenhauer, *J. Chem. Inf. Comput. Sci.*, 2004, **44**, 1177-1185
11. A. R. Leach and V. J. Gillet *An Introduction to Chemoinformatics*, Kluwer, Dordrecht, 2003
12. J. Gasteiger, ed. *Handbook of Chemoinformatics*. Wiley-VCH, Weinheim, 2003

13. R. E. Carhart, D. H. Smith and R. Venkataraghavan, *J. Chem. Inf. Comput. Sci.*, 1985, **25**, 64-73
14. P. Willett, V. Winterman and D. Bawden, *J. Chem. Inf. Comput. Sci.*, 1986, **26**, 36-41
15. A. K. Ghose and V. N. Viswanadhan, eds., *Combinatorial Library Design and Evaluation: Principles, Software Tools and Applications in Drug Discovery*, Marcel Dekker, New York, 2001
16. G. M. Downs and J. M. Barnard, *Rev. Comput. Chem.*, 2002, **18**, 1-40
17. R. D. Brown and Y. C. Martin, *J. Chem. Inf. Comput. Sci.*, 1996, **36**, 572-584
18. R. D. Brown and Y. C. Martin, *J. Chem. Inf. Comput. Sci.*, 1997, **37**, 1-9
19. The MDL Drug Data Report database is available from MDL Information Systems Inc. at <http://www.mdli.com>
20. S. J. Edgar, J. D. Holliday and P. Willett, *J. Mol. Graph. Model.*, 2000, **18**, 343-357
21. J. E. Dubois, in *Chemical Applications of Graph Theory*, ed. A. T. Balaban, Academic Press, London, 1976, p. 161
22. M. Randic, *J. Chem. Inf. Comput. Sci.*, 1978, **18**, 101-107
23. P. Willett, *J. Chem. Inf. Comput. Sci.*, 1979, **19**, 159-162
24. W. Schubert and I. Ugi, *J. Amer. Chem. Soc.*, 1978, **100**, 37-41
25. W. Bremser, *Anal. Chim. Acta*, 1978, **103**, 355-365
26. A. Bender, H. Y. Mussa, R. C. Glen and S. Reiling, *J. Chem. Inf. Comput. Sci.*, 2004, **44**, 170-178
27. A. E. Klon, M. Glick, M. Thoma, P. Acklin and J. W. Davies, *J. Med. Chem.*, 2004, **47**, 2743-2749
28. G. Schneider, W. Neidhart, T. Giller and G. Schmid, *Angewandte Chem. Int.*, 1999, **38**, 2894-2896
29. U. Fechner, F. Lutz, S. Renner, P. Schneider and G. Scheider, *J. Comput.-Aid. Mol. Design*, 2003, **17**, 687-698
30. H. Matter and T. Pötter, *J. Chem. Inf. Comput. Sci.*, 1999, **39**, 1211-1225
31. J. S. Mason, I. Morize, P. R. Menard, D. L. Cheney, C. Hulme and R. F. Labaudiniere, *J. Med. Chem.*, 1999, **42**, 3251-3264
32. G. W. Bemis, M. A. Murcko, *J. Med. Chem.*, 1996, **39**, 2887-2893
33. Y. J. Xu, M. Johnson, *J. Chem. Inf. Comput. Sci.*, 2001, **41**, 181-185
34. Y. J. Xu, M. Johnson, *J. Chem. Inf. Comput. Sci.*, 2002, **42**, 912-926
35. Scitegic Inc. is at <http://www.scitegic.com>
36. D. J. Wilton, P. Willett, K. Lawson & G. Mullier, *J. Chem. Inf. Comput. Sci.*, 2003, **43**, 469-474
37. Tripos Inc. is at <http://www.tripos.com>
38. C. M. R. Ginn, P. Willett and J. Bradshaw, *Perspect. Drug Discov. Design*, 2000, **20**, 1-16
39. N. Salim, J. D. Holliday and P. Willett, *J. Chem. Inf. Comput. Sci.*, 2003, **43**, 435-442
40. G. Harper, J. Bradshaw, J. C. Gittins, D. V. S. Green and A. R. Leach, *J. Chem. Inf. Comput. Sci.*, 2001, **41**, 1295-1300
41. Barnard Chemical Information Ltd. is at <http://www.bci.gb.com/>
42. Daylight Chemical Information Systems Inc. is at <http://www.daylight.com>
43. H. L. Morgan, *J. Chem. Docum.*, 1965, **5**, 107-113

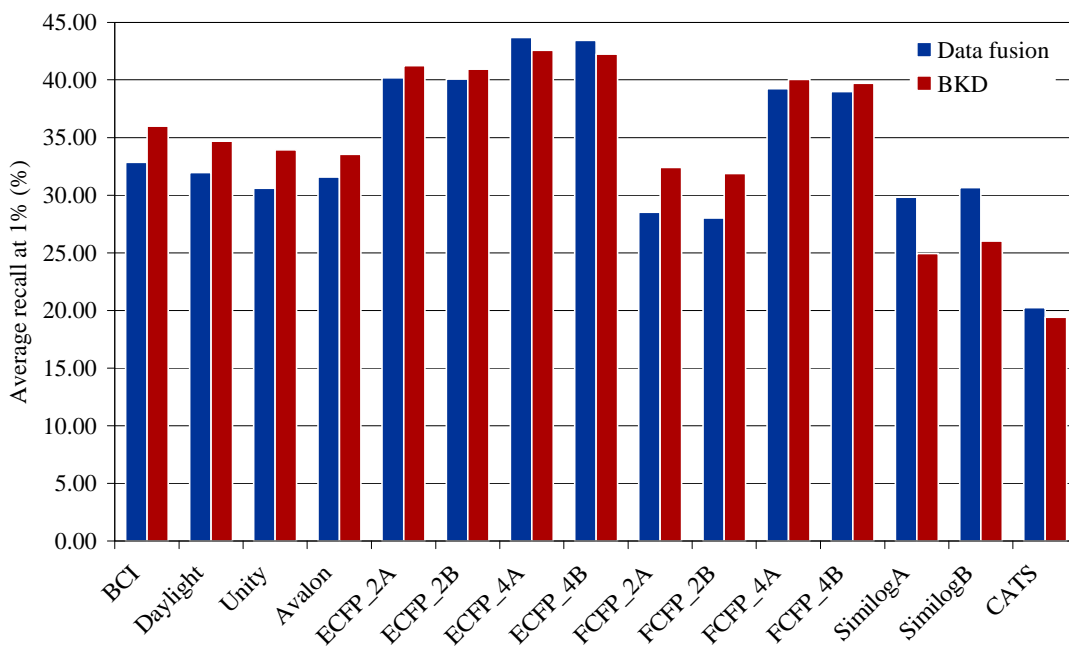


Figure 1. Comparison of the average recalls at 1% obtained using the BKD and the data fusion approaches

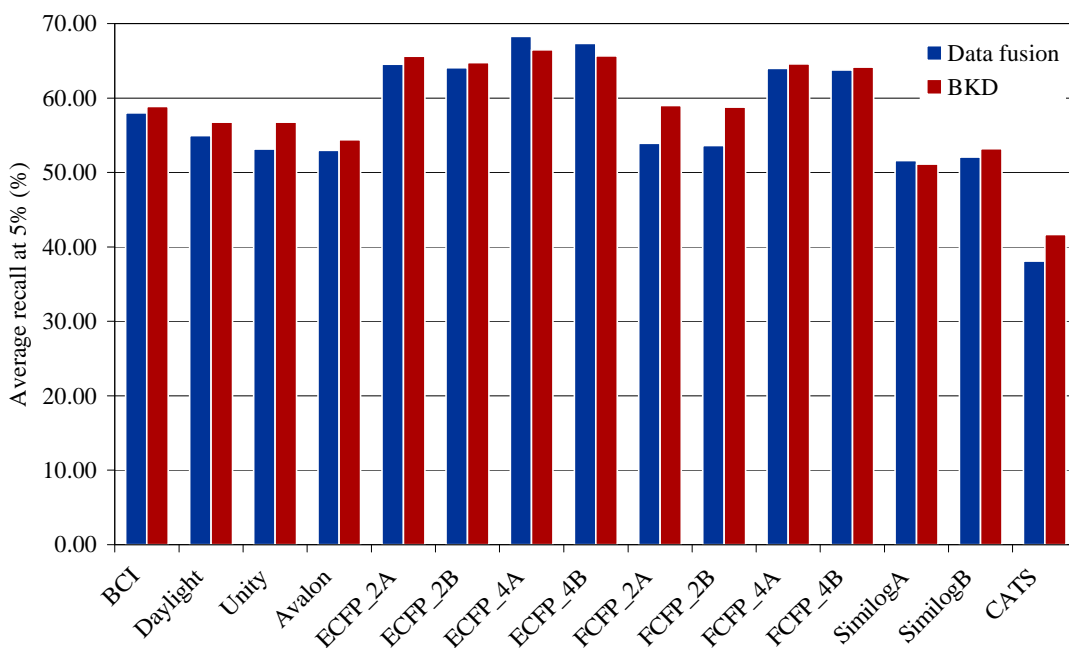


Figure 2. Comparison of the average recalls at 5% obtained using the BKD and the data fusion approaches

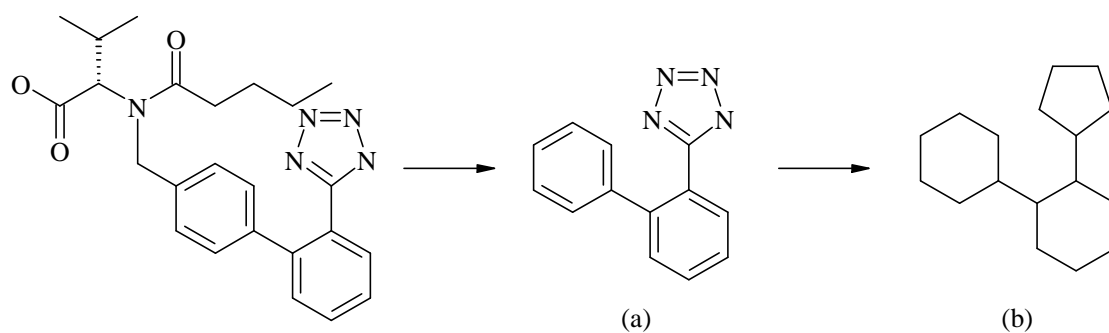


Figure 3. Example of (a) molecular framework (or cyclic system) and (b) framework (or skeletal cyclic system) of Diovan®

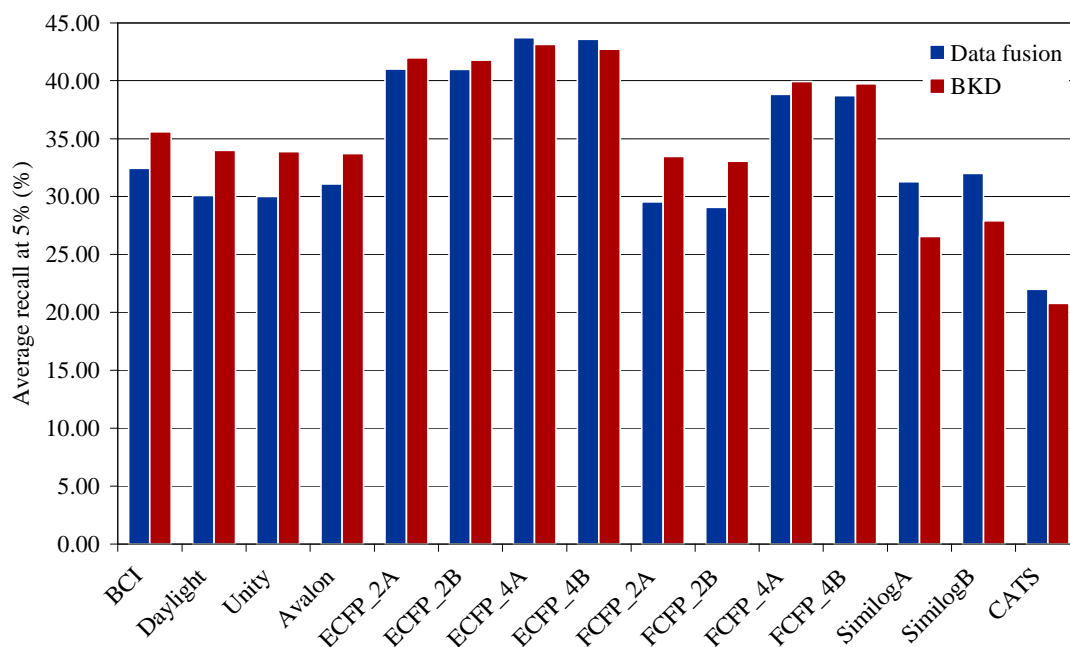


Figure 4. Comparison of the average percentage of molecular frameworks retrieved in the top 1% of the ranked test set obtained using BKD and data fusion.

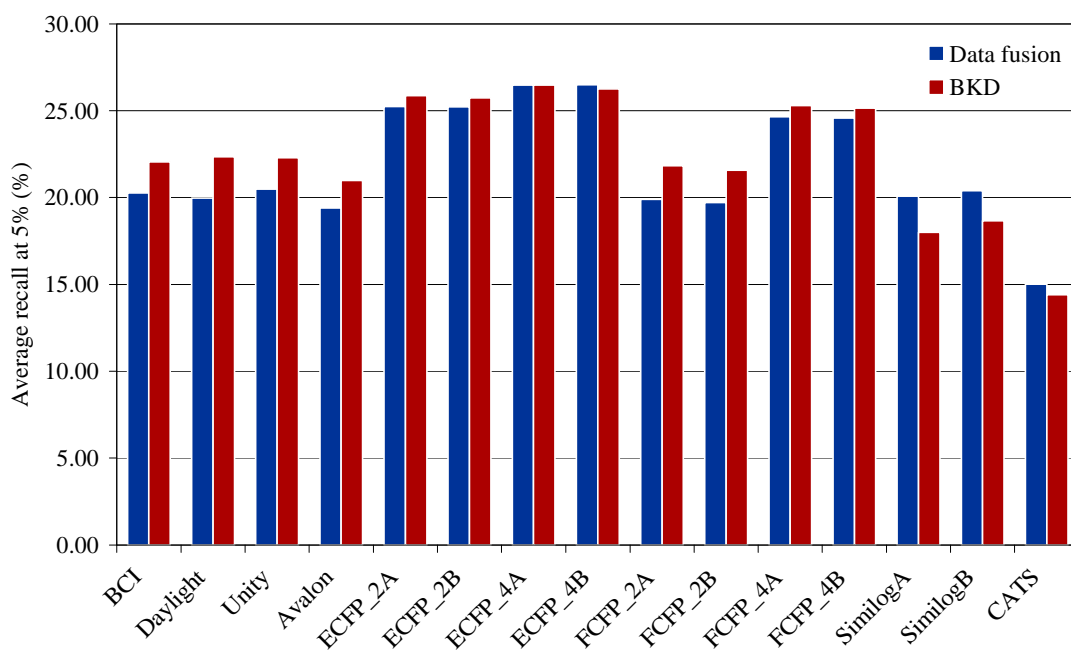


Figure 5. Comparison of the average percentage of frameworks retrieved in the top 1% of the ranked test set obtained using BKD and data fusion.

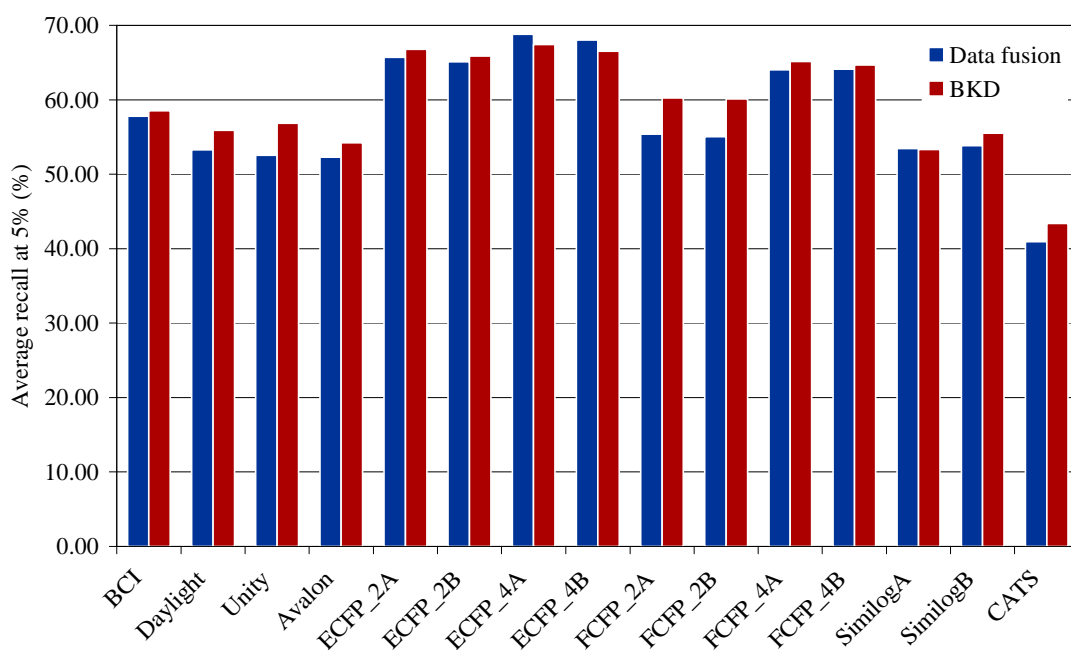


Figure 6. Comparison of the average percentage of molecular frameworks retrieved in the top 5% of the ranked test set obtained using BKD and data fusion.

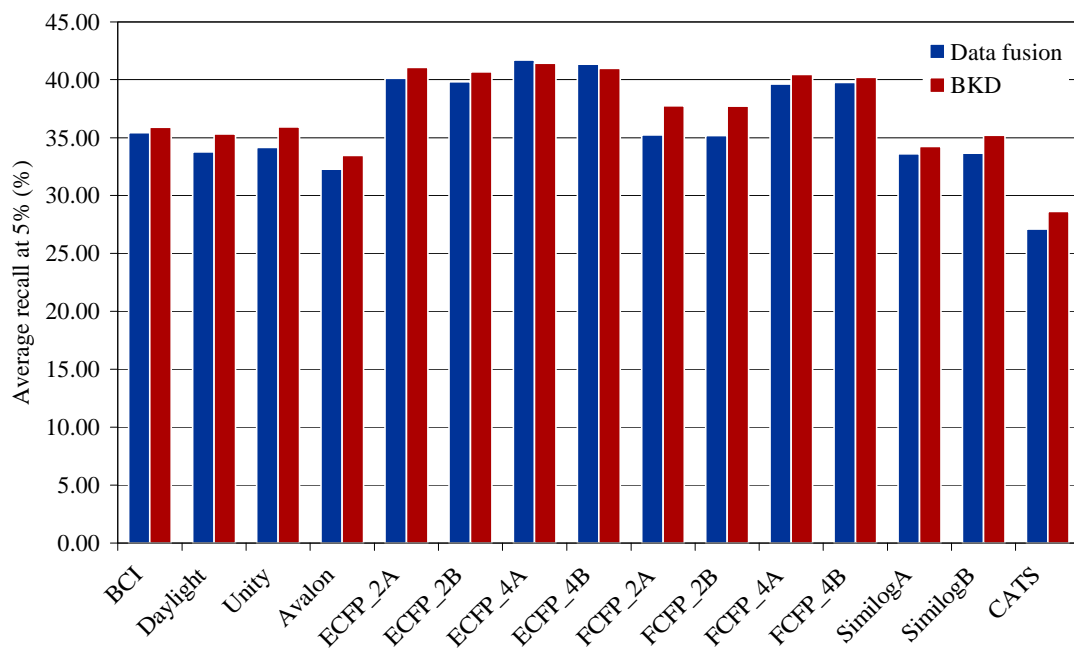


Figure 7. Comparison of the average percentage of frameworks retrieved in the top 5% of the ranked test set obtained using BKD and data fusion.



Activity Classes	BCI		Daylight		Unity		Avalon		SimilogA		SimilogB		CATS	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	35.51	3.32	27.55	4.97	28.99	6.12	34.91	4.84	34.58	11.25	41.28	10.60	11.79	2.50
5HT1A agonists	27.41	6.36	21.70	4.90	19.52	3.88	21.29	5.70	16.22	5.41	20.07	6.20	9.02	2.14
5HT Reuptake inhibitors	25.50	5.57	29.28	6.80	27.82	6.61	24.44	6.28	19.40	6.21	22.15	6.86	9.60	1.86
D2 antagonists	27.12	4.92	24.86	6.31	23.19	5.57	23.30	4.57	17.95	3.75	21.61	3.85	10.52	4.02
Renin inhibitors	65.14	8.73	61.60	7.25	62.29	5.50	62.98	4.42	32.35	18.85	32.32	15.73	43.00	10.35
Angiotensin II AT1 antagonists	47.12	2.31	47.73	2.58	47.19	2.05	45.89	2.58	36.78	6.68	36.10	4.56	38.55	6.39
Thrombin inhibitors	39.04	7.01	32.60	5.58	33.69	7.05	36.54	7.41	11.16	5.93	11.26	6.22	30.44	8.53
Substance P antagonists	31.98	6.09	33.01	5.47	31.68	3.98	22.56	4.44	15.70	6.51	15.39	6.90	7.98	3.15
HIV protease inhibitors	39.35	8.07	44.07	9.32	41.30	7.92	39.62	9.34	42.24	9.54	40.32	5.68	25.62	8.81
Cyclooxygenase inhibitors	24.58	5.10	21.23	5.27	20.80	4.22	20.99	2.82	16.21	2.66	16.52	3.64	10.45	2.52
Protein kinase C inhibitors	33.39	7.56	38.10	9.04	36.93	9.73	36.43	10.25	31.65	10.93	29.07	10.42	16.73	6.66
Average over all classes	36.01	5.91	34.70	6.14	33.95	5.69	33.54	5.70	24.93	7.97	26.01	7.33	19.43	5.18

Activity Classes	ECFP_2A		ECFP_2B		ECFP_4A		ECFP_4B		FCFP_2A		FCFP_2B		FCFP_4A		FCFP_4B	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	45.84	4.18	43.19	5.62	45.93	7.00	44.66	6.57	38.37	3.64	36.48	4.46	45.07	4.51	44.31	4.07
5HT1A agonists	29.42	6.25	29.24	5.29	33.45	6.15	33.44	5.55	18.87	2.29	18.02	2.54	28.52	4.27	27.98	3.98
5HT Reuptake inhibitors	31.78	5.88	31.20	4.89	32.81	6.09	31.00	6.36	29.57	4.03	28.94	4.40	33.58	6.90	33.52	7.06
D2 antagonists	28.60	6.18	28.68	5.33	30.36	6.70	30.47	6.59	20.65	4.84	20.70	4.49	29.14	6.88	28.42	6.55
Renin inhibitors	76.59	2.40	76.17	2.71	77.91	3.37	77.97	2.28	49.51	6.46	48.92	6.80	71.14	7.55	71.16	8.03
Angiotensin II AT1 antagonists	49.00	3.03	49.95	3.01	49.57	3.36	49.80	2.98	44.93	2.16	45.22	1.89	48.31	3.18	48.71	3.23
Thrombin inhibitors	43.91	11.88	43.46	11.99	41.98	9.36	41.41	9.26	34.89	8.12	34.07	7.50	39.89	6.68	39.71	6.30
Substance P antagonists	35.56	6.30	37.00	6.01	38.29	9.33	37.38	9.22	28.33	5.26	27.66	5.48	35.26	6.39	34.72	6.26
HIV protease inhibitors	52.41	8.49	51.59	7.64	55.28	8.63	56.28	7.89	34.93	7.04	35.34	7.97	44.77	8.25	44.64	8.52
Cyclooxygenase inhibitors	23.07	3.25	22.11	3.25	22.80	3.62	22.35	3.73	20.48	3.33	20.34	3.16	24.46	5.04	23.75	4.60
Protein kinase C inhibitors	37.45	10.93	37.67	9.95	39.95	11.89	40.14	9.75	35.94	7.42	34.94	9.65	40.34	11.50	40.07	11.31
Average over all classes	41.24	6.25	40.93	5.97	42.58	6.86	42.26	6.38	32.41	4.96	31.88	5.31	40.04	6.47	39.73	6.36

Table 1. Comparison of the average recalls at 1% obtained with BKD

Activity Classes	BCI		Daylight		Unity		Avalon		SimilogA		SimilogB		CATS	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	35.51	4.08	31.91	2.84	31.04	5.29	30.54	4.50	34.43	3.64	34.69	4.61	8.71	3.15
5HT1A agonists	30.98	5.50	22.63	3.88	20.06	4.33	19.41	5.28	26.18	4.16	27.31	4.62	10.11	2.58
5HT Reuptake inhibitors	28.74	4.38	30.86	5.48	29.03	4.83	31.78	3.37	27.88	5.92	27.99	5.49	9.37	2.87
D2 antagonists	27.87	3.66	25.09	4.89	22.65	4.48	23.40	4.81	30.23	4.14	28.96	4.61	10.44	3.75
Renin inhibitors	52.96	8.34	51.29	3.34	50.23	3.18	56.71	4.03	43.59	15.43	48.64	13.78	60.61	1.99
Angiotensin II AT1 antagonists	43.10	2.95	43.44	3.03	41.39	4.04	42.96	4.74	44.15	1.82	44.96	1.79	40.15	2.21
Thrombin inhibitors	36.25	9.31	27.59	7.37	27.74	9.23	34.15	9.85	17.65	9.82	20.38	9.90	22.47	14.04
Substance P antagonists	23.82	5.82	24.66	5.07	23.07	4.61	17.96	4.76	19.69	6.82	20.19	6.63	9.27	2.30
HIV protease inhibitors	25.57	6.77	33.14	8.17	33.78	7.21	33.07	7.86	39.51	5.24	39.84	4.66	27.89	7.36
Cyclooxygenase inhibitors	22.59	6.38	19.76	4.88	18.48	4.83	19.04	4.65	17.24	1.32	16.57	1.57	7.83	1.22
Protein kinase C inhibitors	33.91	8.87	41.31	8.24	39.23	6.61	38.42	6.30	27.72	9.41	27.79	9.32	16.03	4.76
Average over all classes	32.84	6.01	31.97	5.20	30.61	5.33	31.59	5.47	29.84	6.16	30.67	6.09	20.26	4.20

Activity Classes	ECFP_2A		ECFP_2B		ECFP_4A		ECFP_4B		FCFP_2A		FCFP_2B		FCFP_4A		FCFP_4B	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	47.33	3.40	45.98	3.06	52.18	4.99	50.77	4.99	32.82	2.81	30.70	3.02	44.16	3.82	44.00	3.45
5HT1A agonists	30.89	5.56	30.81	5.28	35.83	4.42	34.99	3.91	20.21	1.38	20.13	1.90	29.17	3.31	29.16	3.09
5HT Reuptake inhibitors	29.89	4.42	30.06	4.41	31.89	5.80	31.95	5.67	26.39	6.61	25.16	7.21	32.69	6.21	33.12	6.13
D2 antagonists	27.95	6.18	27.58	6.19	31.84	5.42	31.90	5.78	21.06	4.42	20.42	5.00	29.74	5.51	29.22	5.18
Renin inhibitors	72.21	4.79	72.11	3.81	75.10	4.59	75.01	3.96	45.79	5.32	44.85	5.56	65.64	6.51	65.58	6.92
Angiotensin II AT1 antagonists	47.83	3.80	47.71	4.37	49.99	3.95	51.14	3.72	40.96	3.48	41.91	3.69	49.25	3.28	49.51	3.29
Thrombin inhibitors	41.90	11.71	41.31	11.08	41.98	10.59	42.06	10.00	27.33	7.84	26.58	7.75	37.93	9.38	37.14	8.36
Substance P antagonists	32.81	6.25	32.80	6.38	36.51	7.65	36.05	8.25	19.96	3.98	19.72	3.61	30.44	6.64	30.05	6.42
HIV protease inhibitors	48.69	6.83	49.22	6.60	54.31	7.51	54.07	6.65	28.55	6.58	28.85	7.15	43.38	8.09	42.91	8.04
Cyclooxygenase inhibitors	21.61	3.92	21.39	3.86	24.30	3.91	23.74	4.47	17.56	3.28	17.17	3.29	24.03	4.52	23.31	4.41
Protein kinase C inhibitors	41.02	7.38	41.92	6.85	46.70	8.93	46.30	8.08	33.16	9.46	32.62	9.51	45.26	9.12	44.90	8.89
Average over all classes	40.19	5.84	40.08	5.63	43.69	6.16	43.45	5.95	28.53	5.02	28.01	5.25	39.24	6.03	38.99	5.83

Table 2. Comparison of the average recalls at 1% obtained with data fusion

Activity Classes	BCI		Daylight		Unity		Avalon		SimilogA		SimilogB		CATS	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	55.71	4.36	46.25	6.73	52.08	8.35	52.44	6.08	59.51	13.72	65.08	12.88	28.80	2.51
5HT1A agonists	49.99	10.31	43.35	7.61	38.24	7.05	40.16	7.25	45.68	10.67	52.04	9.78	22.94	5.66
5HT Reuptake inhibitors	42.61	8.85	44.56	10.50	45.36	8.66	40.63	10.98	38.25	9.14	44.99	8.29	25.73	5.15
D2 antagonists	46.47	8.49	42.36	9.68	38.62	7.54	35.92	6.04	45.69	7.45	52.03	5.02	21.71	6.53
Renin inhibitors	93.23	3.69	92.54	2.80	93.34	1.38	93.98	2.91	81.25	8.67	82.96	5.73	89.67	2.65
Angiotensin II AT1 antagonists	90.92	2.31	88.87	3.43	84.47	6.61	80.58	4.21	71.52	9.94	70.99	8.06	75.68	5.05
Thrombin inhibitors	68.98	5.31	61.30	8.12	63.01	7.60	69.19	7.30	37.35	13.03	37.12	11.53	60.83	8.43
Substance P antagonists	51.89	9.11	57.01	5.23	58.37	8.25	44.46	7.96	36.13	8.73	33.06	8.55	19.41	4.62
HIV protease inhibitors	66.47	6.73	67.32	9.76	68.41	8.34	60.92	12.82	70.39	8.31	69.93	6.80	60.91	9.18
Cyclooxygenase inhibitors	36.02	6.55	32.30	5.48	33.13	4.59	32.04	4.23	30.75	6.43	30.89	6.80	22.49	2.92
Protein kinase C inhibitors	45.15	7.85	48.35	9.48	49.16	10.98	47.90	8.12	45.67	9.58	46.00	8.33	30.02	7.33
Average over all classes	58.86	6.69	56.75	7.17	56.74	7.21	54.38	7.08	51.11	9.61	53.19	8.34	41.65	5.46

Activity Classes	ECFP_2A		ECFP_2B		ECFP_4A		ECFP_4B		FCFP_2A		FCFP_2B		FCFP_4A		FCFP_4B	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	67.55	6.15	63.05	8.52	65.31	6.74	63.36	7.55	62.70	6.92	61.56	7.41	68.64	6.71	67.67	6.57
5HT1A agonists	54.42	8.14	52.83	6.65	58.65	7.93	57.49	7.74	45.08	5.32	44.81	4.14	55.10	7.23	54.06	6.05
5HT Reuptake inhibitors	50.54	4.00	49.11	2.90	50.26	5.88	49.00	6.53	49.77	3.66	50.43	3.95	50.49	5.69	50.14	7.01
D2 antagonists	50.42	8.49	51.48	7.31	55.19	9.46	54.18	9.02	46.18	9.45	45.14	8.68	52.62	8.50	51.74	9.24
Renin inhibitors	97.58	0.58	97.60	0.54	96.74	0.77	96.99	0.77	92.09	4.76	92.94	3.90	97.57	1.28	97.59	1.08
Angiotensin II AT1 antagonists	97.02	1.66	96.86	1.72	97.97	0.80	97.81	0.58	86.84	5.48	87.33	5.15	94.97	4.66	95.09	4.14
Thrombin inhibitors	77.60	8.71	75.08	9.96	74.79	8.77	74.07	8.01	69.90	5.85	67.89	6.96	74.12	6.92	73.77	6.37
Substance P antagonists	61.97	8.17	62.74	7.66	67.31	9.50	65.51	10.24	48.78	6.49	48.09	8.10	59.80	7.93	59.10	9.06
HIV protease inhibitors	79.32	7.72	79.07	7.13	80.78	6.00	80.76	5.02	62.66	9.46	63.39	9.98	70.49	11.83	70.61	12.21
Cyclooxygenase inhibitors	35.99	5.42	34.55	5.25	34.35	5.44	32.67	6.20	36.71	4.66	35.62	4.78	36.93	6.20	36.47	6.79
Protein kinase C inhibitors	49.21	11.44	49.89	11.15	49.57	13.65	50.29	8.85	48.22	7.79	49.41	8.09	49.73	12.72	49.10	12.06
Average over all classes	65.60	6.41	64.75	6.25	66.45	6.81	65.65	6.41	58.99	6.35	58.78	6.47	64.59	7.24	64.12	7.33

Table 3. Comparison of the average recalls at 5% with BKD

Activity Classes	BCI		Daylight		Unity		Avalon		SimilogA		SimilogB		CATS	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	58.80	5.89	51.33	3.49	49.03	5.43	52.36	4.05	49.51	7.74	48.81	8.41	20.85	6.19
5HT1A agonists	54.71	5.77	40.88	5.59	37.16	4.05	34.71	4.74	48.75	6.14	49.82	7.12	23.01	3.64
5HT Reuptake inhibitors	45.36	4.66	46.85	5.43	49.66	5.48	47.59	5.22	46.62	6.07	46.73	6.03	23.15	6.42
D2 antagonists	48.26	4.38	42.42	6.60	37.40	4.92	33.30	6.03	53.14	7.68	50.29	7.23	21.77	5.14
Renin inhibitors	93.54	1.30	90.10	1.95	88.62	1.90	90.00	4.02	85.53	2.69	87.54	2.58	91.14	1.49
Angiotensin II AT1 antagonists	86.33	3.54	86.90	1.99	80.45	6.08	82.02	4.63	82.03	4.79	84.07	4.60	71.24	4.81
Thrombin inhibitors	66.58	5.57	56.47	7.56	58.56	8.98	63.25	8.60	35.70	10.90	39.71	10.87	43.30	16.02
Substance P antagonists	44.83	7.18	51.82	6.28	47.14	5.16	39.90	3.73	36.79	7.55	36.47	6.96	24.05	2.71
HIV protease inhibitors	58.95	4.60	58.69	6.97	61.62	7.85	56.08	7.43	63.54	4.53	63.45	4.56	56.70	10.08
Cyclooxygenase inhibitors	33.35	7.78	29.87	7.81	26.52	7.15	30.93	6.74	28.08	4.37	27.97	4.18	15.88	1.49
Protein kinase C inhibitors	47.25	9.39	48.87	8.29	48.01	8.99	52.26	5.21	37.83	8.55	38.04	8.77	27.99	5.24
Average over all classes	58.00	5.46	54.93	5.63	53.11	6.00	52.95	5.49	51.59	6.46	52.08	6.48	38.10	5.75

Activity Classes	ECFP_2A		ECFP_2B		ECFP_4A		ECFP_4B		FCFP_2A		FCFP_2B		FCFP_4A		FCFP_4B	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
5HT3 antagonists	69.66	4.84	68.54	4.88	72.24	6.05	70.38	6.42	54.49	3.32	53.21	3.21	66.46	5.68	65.80	5.17
5HT1A agonists	55.73	6.50	55.65	6.30	64.19	4.45	63.17	4.60	45.48	3.94	45.80	4.00	55.40	4.54	55.72	3.86
5HT Reuptake inhibitors	48.08	4.01	47.91	2.78	49.68	5.74	49.51	4.22	45.99	5.73	44.93	5.26	51.12	5.15	50.29	5.62
D2 antagonists	50.29	5.98	48.94	5.26	56.05	5.39	54.70	5.90	44.36	5.86	45.58	6.10	50.60	6.22	51.17	5.84
Renin inhibitors	97.27	0.72	97.35	0.63	96.76	0.72	96.92	0.87	88.61	3.09	87.60	3.52	97.22	0.85	97.04	0.74
Angiotensin II AT1 antagonists	95.55	2.11	95.14	2.33	97.36	0.73	97.13	0.66	83.65	4.05	83.82	3.82	94.86	2.46	94.99	2.62
Thrombin inhibitors	73.98	6.64	71.68	6.80	74.73	6.41	72.50	7.27	58.88	6.88	58.07	6.74	70.15	4.20	70.11	4.32
Substance P antagonists	55.16	7.22	54.40	8.07	62.22	7.25	61.42	7.61	37.25	6.44	36.93	5.99	53.45	7.04	53.01	6.69
HIV protease inhibitors	76.24	4.54	76.92	4.63	79.97	4.67	79.41	4.80	58.99	6.76	58.41	6.46	71.86	8.94	70.73	8.29
Cyclooxygenase inhibitors	34.84	6.11	34.11	6.05	40.06	6.21	38.77	5.72	29.76	5.39	29.97	5.36	37.78	6.69	36.98	5.75
Protein kinase C inhibitors	52.87	7.05	53.93	6.97	57.83	6.72	56.95	6.64	45.37	9.98	45.46	9.94	54.60	7.40	55.21	7.62
Average over all classes	64.52	5.07	64.05	4.97	68.28	4.94	67.35	4.97	53.89	5.58	53.62	5.49	63.95	5.38	63.73	5.14

Table 4. Comparison of the average recalls at 5% with data fusion

Activity class	Actives	Number Of		Similarity	
		Assemblies	Frameworks	Mean	SD
5HT3 antagonists	752	438	237	0.351	0.116
5HT1A agonists	827	478	271	0.343	0.104
5HT Reuptake inhibitors	359	193	126	0.345	0.122
D2 antagonists	395	270	187	0.345	0.103
Renin inhibitors	1130	595	339	0.573	0.106
Angiotensin II AT1 antagonists	943	496	285	0.403	0.101
Thrombin inhibitors	803	451	295	0.419	0.127
Substance P antagonists	1246	633	380	0.399	0.106
HIV protease inhibitors	750	475	331	0.446	0.122
Cyclooxygenase inhibitors	636	308	139	0.268	0.093
Protein kinase C inhibitors	453	190	134	0.323	0.142

Table 5. MDDR activity classes used in the study

Type A	Type B	8-bit string
1 occurrence	$2^0 \leq \text{occurrences} < 2^1$	10000000
2 occurrences	$2^1 \leq \text{occurrences} < 2^2$	11000000
3 occurrences	$2^2 \leq \text{occurrences} < 2^3$	11100000
4 occurrences	$2^3 \leq \text{occurrences} < 2^4$	11110000
5 occurrences	$2^4 \leq \text{occurrences} < 2^5$	11111000
6 occurrences	$2^5 \leq \text{occurrences} < 2^6$	11111100
7 occurrences	$2^6 \leq \text{occurrences} < 2^7$	11111110
8 and more occurrences	$2^7 \leq \text{occurrences}$	11111111

Table 6. Binning schemes to convert the occurrence of Similog keys to incidences.

Name	Type	Normalised	Abbreviation	Length	Mean	SD	Max	Min	Density
Barnard Chemical Information	Dictionary-based	-	BCI	1052	96.67	30.91	264	8	9.19
Daylight	Hashed	-	Daylight	2048	289.45	111.21	1046	24	14.13
Unity	Hashed	-	Unity	988	219.65	69.16	558	27	22.23
Avalon	Hashed	-	Avalon	2048	285.06	149.31	1076	16	13.92
ECFP_2	Circular substructure	A	ECFP_2A	7445	32.36	9.40	103	5	0.44
ECFP_2	Circular substructure	B	ECFP_2B	1024	31.82	9.12	98	5	3.11
ECFP_4	Circular substructure	A	ECFP_4A	142864	53.97	16.95	191	8	0.04
ECFP_4	Circular substructure	B	ECFP_4B	1024	52.43	15.95	177	8	5.12
FCFP_2	Circular substructure	A	FCFP_2A	600	20.88	5.00	47	5	3.48
FCFP_2	Circular substructure	B	FCFP_2B	1024	20.41	4.83	45	5	1.99
FCFP_4	Circular substructure	A	FCFP_4A	30267	40.63	11.14	122	7	0.13
FCFP_4	Circular substructure	B	FCFP_4B	1024	39.44	10.55	113	7	3.85
Similog	Pharmacophore	A	SimilogA	64248	1308.14	1437.17	14740	1	2.04
Similog	Pharmacophore	B	SimilogB	64248	863.52	900.10	10101	1	1.34
CATS	Pharmacophore	-	CATS	1200	95.99	36.18	453	1	8.00

Table 7. Comparison of the numbers of bits set in each of the 15 types of fingerprint evaluated in the study.