

Inverse Frequency Weighting of Fragments for Similarity-Based Virtual Screening

Shereena M. Arif,^{†,‡} John D. Holliday,[‡] and Peter Willett^{*,‡}

Faculty of Information Science and Technology, National University of Malaysia, 43600 UKM Bangi, Malaysia and Information School, University of Sheffield, Sheffield S10 2TN, United Kingdom

Received April 1, 2010

This paper discusses the weighting of two-dimensional fingerprints for similarity-based virtual screening, specifically the use of weights that assign greatest importance to the substructural fragments that occur least frequently in the database that is being screened. Virtual screening experiments using the MDL Drug Data Report and World of Molecular Bioactivity databases show that the use of such inverse frequency weighting schemes can result, in some circumstances, in marked increases in screening effectiveness when compared with the use of conventional, unweighted fingerprints. Analysis of the characteristics of the various schemes demonstrates that such weights are best used to weight the fingerprint of the reference structure in a similarity search, with the database structures' fingerprints unweighted. However, the increases in performance resulting from such weights are only observed with structurally homogeneous sets of active molecules; when the actives are diverse, the best results are obtained using conventional, unweighted fingerprints for both the reference structure and the database structures.

INTRODUCTION

Similarity searching using two-dimensional (2D) fingerprints is one of the most common techniques for ligand-based virtual screening, since it has been shown to offer high levels of screening effectiveness despite its very limited computational requirements.^{1–5} Molecules are here represented by fingerprint vectors encoding the presence of topological fragment substructures, with the similarity between a user-defined reference structure and each of the structures in a screening database computed using an association coefficient, most usually the Tanimoto coefficient.^{6,7}

The Tanimoto coefficient (and other such coefficients) is based on the numbers of fragments common and noncommon to two fingerprints. The coefficient normally takes no account of the natures of those fragments, since all fragments are considered to contribute equally to the measurement of intermolecular similarity. Weighting schemes that can distinguish between fragments hence represent an additional source of information that could be used to enhance the effectiveness of similarity-based screening with a high-weight fragment that is common to both a reference and database structure making a greater contribution to the overall degree of similarity between those two molecules than does a lower-weight fragment that is in common. Here, we report an evaluation of one type of information that can be used for weighted similarity searching. This approach, which we shall refer to as *inverse frequency weighting*, is described in the next section. We then report simulated screening experiments using the MDL Drug Data Report database (MDDR, from Symyx Technologies Inc. at <http://www.symyx.com>) and the World Of Molecular Bioactivity database (WOMBAT, from Sunset Molecular Discovery LLC at <http://sunsetmolecular.com>) and discuss the variations in screening effectiveness that are observed with different types of fingerprints.

INVERSE FREQUENCY WEIGHTING

The weighting approaches in which we are interested make use of information about how frequently fragments occur within molecules. The simplest type of information is that relating to fragment *incidence*, i.e., whether a fragment is either present or absent, so that the possible weights are either one or zero (as typically denoted by a bit being switched on or off, respectively, in a fingerprint). Instead of encoding merely the incidence of fragments, use may also be made of information relating to their *occurrence*, i.e., the number of times that a particular fragment occurs in a molecule. We have recently described a detailed analysis of the use of occurrence-based weighting schemes and shown that they can result in substantial increases in the effectiveness of screening, as compared to conventional binary weighting.⁸ In this paper, we extend our work by considering the occurrences of fragments within the entire database that is being screened, rather than within the individual molecules comprising a database as in our previous study.

The use of within-molecule fragment occurrences in similarity calculations is predicated on the basis that if two molecules have several occurrences of a fragment in common, then they are more similar than if they have just a single occurrence in common.⁸ The use of within-database fragment occurrences has a different basis. The assumption here is that if two molecules have in common a fragment that occurs only rarely in the database as a whole, then they should be regarded as being more similar than if they have in common a fragment that occurs very frequently (e.g., sharing an $\alpha\beta$ -unsaturated carboxylic acid versus sharing a phenyl ring). The assumption that there should be a premium associated with rarity seems intuitively reasonable; however, there have

* Corresponding author. E-mail p.willett@sheffield.ac.uk.

[†] National University of Malaysia.

[‡] University of Sheffield.

been only a few studies to date that have considered this type of fragment weighting.

Adamson and Bush reported a comparison of seven different similarity and dissimilarity measures in simulated property prediction experiments on a quantitative structure–activity relationship (QSAR) data set containing 39 local anesthetics.⁹ Following a suggestion by Harrison,¹⁰ three of the seven measures were based on the use of inverse frequency weighting; none of these performed as well as the measures that treated all of the fragments equally. Simulated property prediction was also used by Willett and Winterman in an extended comparison of similarity coefficients and weighting schemes applied to 16 QSAR and quantitative structure–property relationship (QSPR) data sets;¹¹ their findings mirrored those of Adamson and Bush in showing no advantage to the use of inverse weighting, when compared to other weighting approaches that they considered. Both of these studies involved only small numbers of molecules, so that there could consequently be only small variations in the within-database fragment occurrences and hence in the magnitudes of the fragment weights that were used. Mookk et al. tested an inverse fragment weighting scheme for similarity searching in a large reactions database (rather than in a database of individual molecules as in conventional virtual screening) and found that this gave far better results than did unweighted searching.¹² Finally, Downs et al. carried out searches on three Pfizer screening data sets, each containing ca. 10 000 structures represented by both 2D and 3D fingerprints. The search results were inconsistent, and it was not possible to draw unequivocal conclusions as to the relative merits of inverse frequency weighted and unweighted fingerprints.¹³ Finally, two very recent studies of the application of Bayesian inference networks to virtual screening have included inverse frequency counts as just one of several components in the complex scoring functions that are used to model molecular probabilities of activity.^{14,15}

Given the variable nature of the empirical results to date, it seems appropriate to revisit the use of inverse weighting schemes, especially as work in machine learning suggests that such schemes might prove to be of benefit. Machine learning methods are being increasingly used for virtual screening.¹⁶ Here, we focus on substructural analysis (or naïve Bayesian classifier) methods. Ormerod et al.¹⁷ and Cosgrove and Willett¹⁸ studied a class of fragment weights for substructural analysis that were based on techniques developed for the weighting of keywords in text searching.¹⁹ Two of these weights, R1 and R2, are shown in eqs 1 and 2, where it is assumed that there is a training set containing a total of N compounds. This total comprises the N_A compounds that are, and the N_I compounds that are not, active in the biological screen of interest. There are T_j training set compounds that contain the j -th fragment; of these, A_j compounds are, and I_j compounds are not, active. Then the weights assigned to the j -th fragment using the R1 and R2 weights are

$$R1(j) = \log\left(\frac{A_j/N_A}{T_j/N}\right) \quad (1)$$

$$R2(j) = \log\left(\frac{A_j/N_A}{I_j/N_I}\right) \quad (2)$$

Hert et al.²⁰ showed that R1 and R2 are closely related to the naïve Bayesian classifiers described by Xia et al.²¹ and

Bender et al.,²² respectively. We now show that these weights, which assume the availability of the activity information that characterizes machine learning methods, suggest the use of inverse weighting schemes when such information is not available (as is normally the case when similarity searching methods are used).

Consider the weight R1(j) (i.e., the weight for a fragment j in the reference structure that occurs also in a database structure) when a training set and the associated activity data are not available. Both N and N_A are constants that are independent of j , and hence N_A/N is also a constant (here C), while T_j is the frequency of the j -th fragment in the database as a whole (i.e., the number of binary fingerprints in which it occurs). The value of A_j is unknown but is at least one, if it is assumed (as is normally the case) that an active is being used as the reference structure for the similarity search. Hence

$$R1(j) \geq \log\left(\frac{1}{CT_j}\right) \quad (3)$$

Consider now R2, where N_A/N_I is also a constant (here D) and where A_j is again at least one. I_j is unknown but is at most $T_j - 1$ (since we know that there is at least one active containing j in the shape of the reference structure), and hence

$$R2(j) \geq \log\left(\frac{1}{D(T_j - 1)}\right) \quad (4)$$

In both cases, we obtain an expression that involves the reciprocal of T_j and that provides at least some support for the use of inverse frequency weighting in similarity searching.

METHODS

Our experiments have used the conventional model for similarity-based virtual screening. The similarity is computed between a reference structure of known biological activity and each of the molecules in a database, and a check is then made as to the activity or otherwise of the molecules that are most similar to the reference structure. We have employed the MDDR and WOMBAT data sets that were used in our previous study of frequency-based weighting⁸ and that are detailed in Table 1. In all, there were 102 535 and 138,127 molecules in the MDDR and WOMBAT data sets, respectively. In Table 1, the counts of ring systems contained in active molecules are based on the Murcko scaffolds facility implemented in the Pipeline Pilot software, and the diversity figures are the mean intraclass similarities computed for each activity class using the standard Tripos Unity 2D binary fingerprint and the Tanimoto coefficient.

The MDDR and WOMBAT molecules were represented by four different types of binary fingerprints, all of which are widely used in modern chemoinformatics systems: BCI keys (1052 bits and available from Digital Chemistry Ltd.) are selected to maximize discrimination in substructure searching using a frequency-based selection algorithm; Pipeline Pilot ECFP_4 circular substructures (1024 bits and available from Accelrys Inc.) encode circular substructures; MDL keys (166 bits and available from Symyx Technologies Inc.) encode common fragment substructures; and Sunset keys (560 bits and available from Sunset Molecular Discov-

Table 1. Activity Classes Used in the Virtual Screening Experiments Chosen from the MDDR and WOMBAT Databases

activity class (abbreviation)	active molecules	active ring systems	mean pairwise similarity
MDDR			
5HT3 antagonists (5HT3)	752	417	0.35
5HT1A agonists (5HT1)	827	450	0.34
5HT reuptake inhibitors (5HT)	359	181	0.35
D2 antagonists (D2)	395	258	0.35
renin inhibitors (REN)	1125	554	0.57
angiotensin II AT1 antagonists (ANG)	943	464	0.40
thrombin inhibitors (THR)	803	425	0.42
substance P antagonists (SUBP)	1246	586	0.40
HIV protease inhibitors (HIV)	750	461	0.45
cyclooxygenase inhibitors (COX)	636	282	0.27
protein kinase C inhibitors (PKC)	453	171	0.32
WOMBAT			
5HT3 antagonists (rat)	220	117	0.38
5HT1A antagonists (rat)	592	224	0.40
D2 antagonists (rat)	910	324	0.37
renin inhibitors (human)	474	253	0.59
angiotensin II AT1 antagonists (rat)	724	253	0.44
thrombin inhibitors (human)	421	196	0.42
substance P antagonists (human)	558	186	0.43
HIV protease inhibitors (human)	1128	473	0.44
cyclooxygenase inhibitors (human)	965	220	0.32
protein kinase C inhibitors (rat)	142	31	0.57
acetylcholine esterase inhibitors (human)	503	220	0.37
factor Xa inhibitors (human)	842	328	0.39
matrix metalloprotease inhibitors (human)	694	280	0.44
phosphodiesterase inhibitors (human)	596	270	0.36

ery LLC) combine chemical substructure recognition with topologically relevant pharmacophore patterns based on atom pairs.

A MaxMin selection procedure²³ was used to identify 10 diverse reference structures from each of the activity classes. Cut-offs were applied to the ranked lists resulting from each similarity search to identify the top 1% and 5% of the rankings, and the mean numbers of active molecules in these subsets were obtained by averaging over the 10 searches for each activity class. We refer to this mean number of actives as the *recall*. We also noted the numbers of distinct Murcko scaffolds in the active molecules that were retrieved, rather than just the number of active molecules, to assess the effectiveness of the various weighting schemes for scaffold hopping.^{24,25}

Inverse frequency weighting has been intensively studied for the weighting of keywords in text retrieval.^{26,27} Its use was first reported by Spärck Jones in 1972,²⁸ and inverse document frequency (IDF) weighting is a key component of the scoring functions used in modern search engines. Various forms of the IDF weight have been suggested in the text retrieval literature, and we have used three of these here. We have also considered a further inverse frequency weight that has previously been used specifically for chemoinformatics searching.

Each of the fingerprints (which will be described subsequently as BCI, ECFP_4, MDL, or Sunset) can be considered as a vector, **X**, where the *j*-th element, x_j , denotes the weight that the *j*-th fragment has in that molecule. Thus the conventional, 0/1 binary weight is represented by:

$$W1:x_j = 1$$

Using the notation of the previous section, we shall assume that this *j*-th fragment occurs T_j times ($T_j \geq 0$) in the database

(which contains a total of N molecules). Then we consider in this study the following three closely related weights, all of which have been used for IDF weighting in information retrieval (we use natural logarithms in our experiments):

$$W2:x_j = \log\left(\frac{N}{T_j + 1}\right)$$

$$W3:x_j = \log\left(\frac{N}{T_j}\right) + 1$$

$$W4:x_j = \log\left(\frac{N + 0.5}{T_j + 0.5}\right)$$

There are clearly only very slight differences between these three formulations, with W2 and W4 in particular giving comparable weights for all but the least frequently occurring fragments. W3 can be regarded as a direct modification of W1, with the conventional binary match (i.e., unity being added to the number of common fragments during the calculation of the Tanimoto coefficient) being augmented by an inverse frequency contribution. The final weight is rather different and is suggested by Moock et al.¹² for similarity searching in reaction databases:

$$W5:x_j = \sqrt{\frac{\text{Max}\{T_j\}}{T_j}}$$

The similarity S_{XY} between two fragment vectors **R** and **D**, representing the reference structure and the database structure, respectively, was computed in all cases using the full form of the Tanimoto coefficient, i.e.,

$$S_{RD} = \frac{\sum r_j d_j}{\sum r_j^2 + \sum d_j^2 - \sum r_j d_j}$$

where r_j and d_j denote the weight of the j -th fragment for the reference and database structures, and where the summations are over all of the elements in each fingerprint (i.e., 1052, 1024, 166, or 560 elements for BCI, ECFP_4, MDL, or Sunset, respectively). Experiments were also carried out using the modified Tanimoto coefficient that has been suggested by Wang and Bajorath.²⁹ The results were found to be inferior to those obtained with the conventional form of the coefficient and are hence not discussed further.

RESULTS

Details of the characterizations resulting from use of the weights W1–W5 are shown in Table 2. The bottom part of the table lists the mean weights, when averaged over the nonzero elements for all of the fingerprints. It will be seen that the MDL and Sunset inverse weights are noticeably less than the corresponding ECFP_4 and BCI weights. This is due to the much greater frequencies of occurrence for the fragments comprising the former two types of fingerprint, as exemplified by the figures in the fingerprint density row of Table 2.

The frequency distributions for the four sets of weights are illustrated in Figure 1. These distributions are for W3 and the MDDR data set, but similar behavior is observed for W2, W4, and W5 and for the WOMBAT data set. The MDL and Sunset fingerprints have markedly skewed distributions in which most of the fragments have low weights, whereas the ECFP_4 fingerprints exhibit a much more even distribution of weights, with the BCI fingerprints intermediate between these two extremes. The distributions ECFP_4, Sunset, BCI, and MDL have skewness values of 0.91, 3.24, 2.13, and 3.37, respectively. Thus, 21.3% of the ECFP_4 fragments have W3 weights in the range of 1.00–1.99, with the corresponding values for the BCI, Sunset, and MDL fingerprints being 49.6, 74.9, and 80.0%, respectively. The MDL fingerprint has the most skewed distribution of fragment frequencies, with both the largest fraction of high-frequency fragments and the smallest fraction of low-frequency fragments.

Each of the five different weighting schemes can be applied to the reference and to each of the database structures, giving a total of 25 possible similarity measures for the searches using a given type of fingerprint. Here, we have considered a subset of the possible combinations: those where either the reference structure or the database structures use conventional, binary-weighted fingerprints, and those

where both the reference and database structures use the same type of weighted fingerprint. This gives a total of 13 different types of fingerprint to be considered. For brevity in what follows, we refer to each similarity measure by Mab , where a denotes the weight applied to the database structures' fingerprints, and b denotes the weight applied to the reference structure's fingerprint so that, e.g., M15 refers to the set of searches (ten searches for each of the chosen activity classes) in which the database structure fingerprints use W1 (the conventional binary weight) and the reference structure fingerprints use W5 (the reaction searching weight).

A typical set of results is summarized in Table 3, which lists the averaged results for the top 5% searches of the 11 MDDR activity classes using the ECFP_4 fingerprints. Each column lists the recall for the top 5% of a sorted ranking when averaged over 10 reference structures for each activity class (as denoted by the abbreviated form of the class name from Table 1a). The penultimate column on the right is the mean value for that similarity measure when averaged over the 11 activity classes. The weighting scheme with the best mean recall in each column is strongly shaded, and the recall value is bold-faced; any scheme with an average recall within 5% of the value for the best weighting scheme is shown lightly shaded. Thus, M13 is the best performing measure across the 11 activity classes in terms of mean recall, with M11 and M12 also performing well.

Table 3 has been included to exemplify the searches that were carried out. Similar sets of results were obtained, as shown in Table 4, using the following conditions: MDDR or WOMBAT databases; BCI, ECFP_4, MDL or Sunset fingerprints; and recall for the top 1% or 5% of the search rankings. For example, the average results shown in the penultimate right-hand column of Table 3 form the fifth column of results in Table 4. In addition, all of these experiments were repeated using the number of distinct scaffolds (specifically, the Murcko scaffolds in the Pipeline Pilot software) present in the retrieved active molecules. These sets of results are summarized in Table 5.

Use of the mean number of actives retrieved means that data sets with large numbers of actives and with high retrieval rates (such as the renin data set in Table 3) could contribute disproportionately to the overall results. We have hence adopted a further evaluation criterion, which is to count the total number of shaded cells for each similarity measure across the full set of activity classes, as shown in the right-hand column of Table 3. These shaded cell results are listed in Tables 6 and 7 for the retrieval of active molecules and for the retrieval of active scaffolds, respectively. For example, the right-hand column in Table 3, labeled 'Shaded cells' forms the upper half of the results in Table 6.

Table 2. Statistical Data Describing the MDDR and WOMBAT Fingerprints

		MDDR				WOMBAT			
		ECFP_4	Sunset	BCI	MDL	ECFP_4	Sunset	BCI	MDL
number of nonzero elements		5 375 756	20 454 179	10 488 516	5 641 174	6 950 009	26 853 131	14 087 156	8 221 890
mean number of nonzero bits		52.43	199.48	102.29	55.01	50.32	194.41	101.99	59.52
per fingerprint									
fingerprint density		0.05	0.36	0.10	0.33	0.05	0.35	0.10	0.36
mean value of the	W1	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
nonzero elements	W2	2.34	0.71	1.28	0.68	2.33	0.72	1.43	0.70
	W3	3.34	1.71	2.28	1.68	3.33	1.72	2.43	1.68
	W4	2.34	0.71	1.28	0.68	2.33	0.72	1.43	0.70
	W5	3.70	1.35	2.09	1.32	3.82	1.51	2.43	1.50

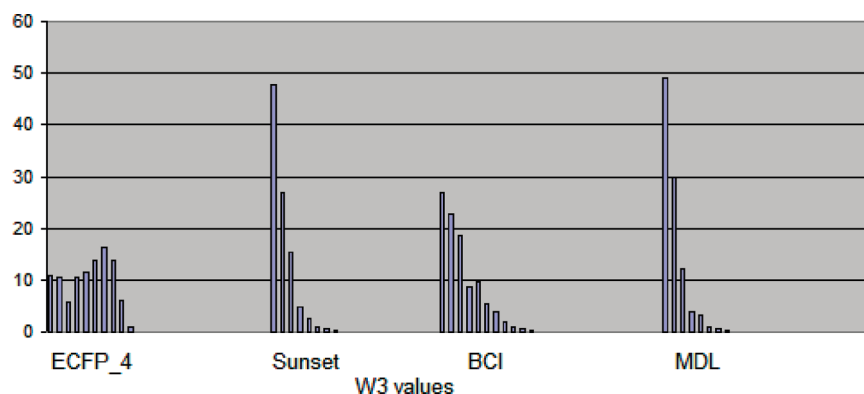


Figure 1. Distribution (as a percentage) of W3 values for the four fingerprints in the MDDR database. Each column represents a step-size of 0.5 units on the x-axis; thus the first column for each fingerprint is for the range 1.00–1.49, the second for 1.50–1.99, etc.

Table 3. Mean Recall of Actives for the Top 5% Searches of MDDR using ECFP_4 Fingerprints

Measure	Activity Class											Mean	Shaded cells
	SHT3	SHT1	SHT	D2	REN	ANG	THR	SUBP	HIV	COX	PKC		
M11	152.70	146.90	51.90	73.10	728.80	502.70	134.40	217.90	178.60	70.70	72.60	211.85	6
M12	113.60	126.40	40.00	61.40	743.90	508.80	154.30	226.70	177.20	47.30	60.60	205.47	3
M13	112.40	124.10	40.30	60.80	792.10	526.60	154.80	235.40	189.20	43.60	61.50	212.80	5
M14	113.50	125.90	39.70	60.20	666.10	507.60	155.50	226.40	177.10	47.20	61.30	198.23	3
M15	107.40	117.80	38.30	57.30	751.40	486.10	152.30	234.80	182.00	43.30	57.10	202.53	3
M21	139.40	148.80	46.20	65.80	631.30	405.00	115.20	196.00	139.30	72.60	63.10	183.88	1
M22	117.70	132.50	40.00	59.90	666.40	389.90	137.50	211.10	153.50	62.30	56.20	184.27	0
M31	148.50	151.30	45.20	65.00	574.90	367.50	99.90	186.10	131.90	77.00	64.50	173.80	3
M33	127.10	139.40	43.30	64.30	696.00	449.30	140.20	216.70	160.00	67.90	62.50	196.97	0
M41	139.30	148.30	46.00	66.00	631.30	405.30	115.00	195.30	139.20	72.30	64.90	183.90	2
M44	118.20	131.70	40.20	60.20	666.10	390.20	137.30	211.50	153.30	62.00	56.10	184.25	0
M51	138.90	147.50	44.80	63.70	581.70	331.20	100.20	175.70	123.80	69.00	56.90	166.67	1
M55	111.50	122.20	36.20	55.80	631.60	325.30	133.20	213.80	148.80	54.30	51.20	171.26	0

Table 4. Mean Recall of Actives Using Different Fingerprints: MDDR and WOMBAT Top 1% and 5%

measure	ECFP_4	Sunset	BCI	MDL	ECFP_4	Sunset	BCI	MDL
MDDR Top 1%					MDDR Top 5%			
M11	109.68	68.23	58.57	50.10	211.85	162.02	154.09	136.28
M12	111.42	70.79	58.61	39.48	205.47	165.21	148.06	108.59
M13	113.05	68.81	59.41	46.75	212.80	166.09	150.44	126.12
M14	111.50	71.15	58.48	41.05	198.23	165.36	148.01	113.23
M15	108.60	70.04	56.75	46.40	202.53	166.94	145.24	125.62
M21	96.29	72.71	45.11	43.70	183.88	169.14	125.51	121.39
M22	100.36	70.05	48.83	31.87	184.27	163.93	133.45	98.16
M31	79.94	45.22	35.16	40.13	173.80	123.59	103.24	115.59
M33	105.41	72.47	53.70	42.05	196.97	170.78	144.13	118.24
M41	89.63	69.07	45.11	43.66	183.90	169.05	125.59	121.34
M44	100.37	70.03	48.81	32.04	184.25	164.00	133.41	98.11
M51	76.58	46.57	29.82	40.43	166.67	124.33	91.85	115.45
M55	90.63	68.24	44.53	37.06	171.26	164.42	128.15	112.00
WOMBAT Top 1%					WOMBAT Top 5%			
M11	103.63	73.02	73.26	65.29	188.16	157.18	154.99	147.45
M12	101.21	82.47	75.29	60.90	184.86	169.24	153.81	145.88
M13	100.77	80.21	74.36	66.11	187.76	170.17	155.74	157.93
M14	101.01	82.52	75.26	60.88	184.61	169.32	153.81	145.87
M15	98.49	81.51	74.40	66.44	181.83	171.39	154.14	156.73
M21	85.57	77.26	58.56	60.20	166.24	161.99	129.88	138.71
M22	95.88	82.34	71.03	53.16	171.11	166.27	144.86	125.80
M31	75.56	54.17	46.06	45.65	155.06	122.16	107.31	115.77
M33	100.44	82.89	74.46	62.46	180.34	166.88	151.07	143.68
M41	85.28	75.97	58.46	60.18	166.58	162.03	129.88	138.64
M44	96.06	82.32	70.91	53.11	171.22	166.28	144.86	125.77
M51	68.45	52.94	36.86	46.27	143.49	121.03	89.14	117.62
M55	85.80	80.68	63.60	58.16	153.57	165.04	138.74	136.59

DISCUSSION

Visual inspection of the results in Tables 4 and 5 suggests the following. The first, unsurprising generalization is that the ECFP_4 fingerprints are consistently superior to the

Table 5. Mean Recall of Active Scaffolds Using Different Fingerprints: MDDR and WOMBAT Top 1% and 5%

measure	ECFP_4	Sunset	BCI	MDL	ECFP_4	Sunset	BCI	MDL
MDDR Top 1%				MDDR Top 5%				
M11	60.01	41.56	32.06	30.33	114.43	95.16	84.38	81.57
M12	61.34	39.86	32.15	23.09	113.81	93.17	81.49	63.32
M13	62.75	39.93	33.25	27.56	117.70	96.41	84.30	74.34
M14	61.43	40.06	32.17	24.25	113.92	93.05	81.49	66.34
M15	60.13	40.62	31.45	27.45	112.80	96.85	81.33	74.17
M21	48.39	42.57	25.16	27.62	100.45	98.92	70.26	73.35
M22	54.68	39.75	26.73	19.82	101.83	93.23	74.69	58.71
M31	42.90	26.86	19.42	25.73	94.88	73.11	57.20	71.63
M33	57.43	41.95	29.15	25.94	108.18	97.15	80.30	70.88
M41	48.38	39.96	25.17	27.56	100.45	93.84	70.26	73.31
M44	54.68	39.71	26.76	19.88	101.88	93.27	74.63	58.63
M51	41.17	28.15	17.13	26.02	91.44	73.84	51.38	71.84
M55	49.28	39.33	24.63	23.15	95.15	94.66	73.16	67.42
WOMBAT Top 1%				WOMBAT Top 5%				
M11	44.87	31.20	30.07	27.96	81.78	69.51	64.09	65.80
M12	44.11	35.46	30.76	25.90	81.34	74.36	64.94	62.99
M13	44.16	35.29	30.59	29.01	82.83	76.94	66.26	69.89
M14	44.01	35.45	30.77	25.86	81.25	74.38	64.81	62.98
M15	43.31	34.66	30.47	29.22	80.43	76.61	65.79	69.25
M21	36.65	37.83	23.01	26.91	71.25	76.59	52.90	63.12
M22	41.83	36.08	29.23	23.59	74.86	72.81	60.86	55.20
M31	32.11	22.06	17.78	19.85	66.28	52.30	42.28	51.16
M33	43.88	30.26	30.86	27.54	78.72	72.20	63.45	63.71
M41	36.65	32.86	22.96	26.93	71.63	71.74	52.90	63.04
M44	41.93	36.08	29.17	23.56	74.91	72.84	60.81	55.19
M51	28.96	22.56	13.94	20.44	61.51	53.11	35.04	51.96
M55	37.71	35.57	25.92	25.44	68.11	73.97	58.56	60.83

others. This fingerprint differs markedly from the other three in the distribution of weights that is observed, as shown in Figure 1, which plots the distribution of the calculated M3 weights for the fragments comprising each fingerprint. It will

Table 6. Numbers of Shaded Cells for Mean Recall of Actives Using Different Fingerprints: MDDR and WOMBAT Top 1% and 5%

measure	ECFP_4	Sunset	BCI	MDL	ECFP_4	Sunset	BCI	MDL
MDDR Top 1%				MDDR Top 5%				
M11	6	2	4	2	6	4	4	2
M12	4	2	4	2	3	1	3	0
M13	6	1	4	3	5	3	4	4
M14	4	3	4	1	3	1	3	0
M15	4	2	2	3	3	3	3	1
M21	2	2	1	2	1	3	2	2
M22	0	3	0	0	0	3	1	1
M31	1	2	2	3	3	2	2	4
M33	1	3	2	0	0	3	2	0
M41	2	3	1	2	1	3	2	2
M44	0	3	0	0	0	0	1	1
M51	1	2	0	2	1	2	2	3
M55	0	1	0	1	0	2	1	0
WOMBAT Top 1%				WOMBAT Top 5%				
M11	10	2	5	6	9	1	6	5
M12	7	3	6	1	6	4	5	4
M13	6	5	3	6	7	5	7	8
M14	7	4	6	1	6	4	5	4
M15	2	7	7	7	6	6	7	7
M21	3	2	3	3	3	2	4	0
M22	6	5	4	1	2	1	3	0
M31	2	2	1	0	3	1	3	2
M33	7	5	6	2	3	1	2	0
M41	2	3	3	3	3	3	3	1
M44	6	5	4	1	2	2	3	0
M51	1	2	0	1	0	2	0	2
M55	1	2	2	1	1	1	2	0

Table 7. Numbers of Shaded Cells for Mean Recall of Active Scaffolds Using Different Fingerprints: MDDR and WOMBAT Top 1% and 5%

measure	ECFP_4	Sunset	BCI	MDL	ECFP_4	Sunset	BCI	MDL
MDDR Top 1%				MDDR Top 5%				
M11	6	2	4	2	5	4	3	2
M12	4	1	2	0	3	1	2	0
M13	5	2	4	3	5	4	4	4
M14	4	2	2	0	3	1	2	0
M15	4	2	2	3	4	4	3	4
M21	1	3	0	2	0	4	1	1
M22	0	3	1	0	0	3	1	1
M31	1	1	2	3	2	2	2	5
M33	0	3	1	1	0	1	1	0
M41	1	3	0	1	1	3	1	1
M44	0	2	2	0	0	3	0	1
M51	0	2	1	3	0	2	2	6
M55	0	0	2	1	0	1	2	1
WOMBAT Top 1%				WOMBAT Top 5%				
M11	9	0	5	7	8	1	7	6
M12	8	3	6	2	6	3	5	4
M13	7	4	3	4	8	6	6	8
M14	6	3	6	2	7	3	5	4
M15	4	3	6	5	6	6	7	7
M21	2	8	3	4	3	4	2	2
M22	5	1	4	1	2	2	1	1
M31	1	2	1	1	2	1	0	0
M33	8	0	3	3	5	0	2	1
M41	4	4	3	4	4	3	1	1
M44	5	1	3	1	2	2	2	1
M51	1	3	0	1	0	3	1	2
M55	1	1	2	1	1	1	3	1

be seen that there is a much larger fraction of high-value weights for ECFP_4 than for the other fingerprints, all of which exhibit skewed distributions reflecting relatively larger numbers of low-weight, i.e., high-frequency fragments. This

difference in the distributions is reflected in the large average values in the lower part of Table 2, and there is thus a much greater scope for inverse frequency weighting to contribute to the similarity scores. However, the nature of the distributions may not be a controlling factor since ECFP_4 has been shown to offer a consistently high level of screening performance in our previous comparative studies of fingerprints.^{30–32} Second, the pattern of behavior for the recall of active molecules (in Table 4) is very similar to that for the recall of active scaffolds (in Table 5), i.e., a similarity measure's ability to retrieve active molecules is a fair guide as to its ability to retrieve active ring scaffolds. Turning to the shaded cell results in Tables 6 and 7, which take no account of the precise magnitudes of the recall results, a further generalization is possible. It will be seen that the higher values (i.e., corresponding to greater number of high-effectiveness searches) tend to occur in the upper parts of these tables (corresponding to the measures M11–M15) than in the remainder, i.e., these measures give better screening performance than the others listed here. However, if statements are to be made about the relative merits of the different approaches, then a quantitative evaluation is required, rather than simple visual inspection. We have hence used Kendall's *W* test of statistical significance to assess the significance, if any, of the differences in screening performance resulting from the various similarity measures.

The Kendall *W* test is used to evaluate the consistency of *k* different sets of ranked judgements of the same set of *N* different objects.^{8,33} In the present context, each fingerprint type acts as a judge of the effectiveness of the various similarity measures, where the effectiveness is taken to be the mean recall for that fingerprint, i.e., *k* = 4 and *N* = 13. Taking the MDDR top 1% results in Table 4 as an example, the mean recall figures are converted to ranks, and the computed value for *W* is 0.76. The significance of this value can be tested using the χ^2 distribution, giving a value of 36.46 for χ^2 with 12 degrees of freedom; this value is highly significant (*p* ≤ 0.001). When a value is statistically significant, Siegel and Castellan suggest that it is possible to give an overall ranking of the *N* objects that are being judged,³³ i.e., the 13 different similarity measures in the present case. Doing this for the MDDR top 1% search results in Table 4, the derived ranking is: M13 > M14 > M11 = M12 = M15 > M33 > M21 > M22 > M41 > M44 > M55 > M31 = M51.

Table 8 gives the level of statistical significance and the associated ranking for the eight sets of results in Tables 4 and 5. It will be seen that there is a very high degree of commonality in the rankings, with M13 in particular providing a level of performance that is generally superior to that of M11, i.e., the conventional binary situation, and to the other measures tested here (of which the best overall performance would appear to be M15). The only obvious exception is the WOMBAT top 1% searches, where M13 is ranked fifth. The related M31 weight performs consistently poorly when compared to M13; and similar comments can be made for the pairs M12 and M21, M14 and M41, and M15 and M51 (the last of which vies with M31 as the worst overall measure). Table 9 shows the comparable rankings based on the shaded cell results in Tables 6 and 7. It should be noted that both of the MDDR 1% rankings are omitted from this table as the computed value of *W* in the Kendall

Table 8. Rankings of Similarity Measures Based on Kendall W Analyses of Numbers of Actives (Tables 4 and 5)

data set	recall type	significance	ranking
MDDR	1%	$p \leq 0.001$	M13 > M14 > M11 = M12 = M15 > M33 > M21 > M22 > M41 > M44 > M55 > M31 = M51
	5%	$p \leq 0.05$	M13 > M11 > M15 > M33 > M12 = M14 > M21 = M41 > M22 > M44 > M55 > M31 > M51
	1% scaffolds	$p \leq 0.01$	M13 > M11 > M14 = M15 > M33 > M21 > M12 > M41 > M44 > M22 > M55 > M51 > M31
	5% scaffolds	$p \leq 0.01$	M13 = M11 > M15 > M33 > M21 > M14 > M12 > M41 > M55 > M44 = M22 > M51 = M31
WOMBAT	1%	$p \leq 0.001$	M12 > M14 = M33 > M15 > M13 > M11 > M22 > M55 > M44 > M21 > M41 > M31 > M51
	5%	$p \leq 0.001$	M13 > M15 > M12 > M14 > M11 > M33 > M44 > M22 > M41 > M21 > M55 > M31 > M51
	1% scaffolds	$p \leq 0.05$	M13 > M12 > M11 = M14 = M15 > M33 > M21 = M22 > M44 > M55 > M41 > M31 = M51
	5% scaffolds	$p \leq 0.001$	M13 > M15 > M12 > M14 > M11 > M33 > M21 > M22 = M44 > M55 > M41 > M31 = M51

Table 9. Rankings of Similarity Measures Based on Kendall W Analyses of Shaded Cells (Tables 6 and 7)^a

data set	recall type	significance	ranking
MDDR	5%	$p \leq 0.01$	M11 > M13 > M15 > M31 > M21 = M41 > M51 > M12 = M14 > M33 > M22 > M55 > M44
	5% scaffolds	$p \leq 0.01$	M13 > M11 = M15 > M31 > M51 > M21 = M41 > M12 = M14 > M22 > M55 > M44 > M33
	1%	$p \leq 0.001$	M15 > M33 > M11 > M13 > M14 > M12 > M22 = M44 > M41 > M21 > M55 > M31 = M51
WOMBAT	5%	$p \leq 0.001$	M13 > M15 > M12 = M14 > M11 > M41 > M21 > M31 > M44 > M51 > M22 = M33 > M55
	1% scaffolds	$p \leq 0.05$	M12 > M15 > M11 = M13 > M14 > M21 = M41 > M33 > M22 > M44 > M51 > M31 > M55
	5% scaffolds	$p \leq 0.001$	M13 > M15 > M11 > M14 > M12 > M21 > M41 > M51 > M44 > M33 > M55 > M22 > M31

^a Note that the two MDDR 1% results have been omitted since the computed value of W is not significant at the 0.05 level of statistical significance.

test was not statistically significant ($p \geq 0.05$). While M11 and M13 are again the best performers, M15 also figures strongly at the head of the rankings, and there is generally a less marked difference in the performance of pairs of measures, such as M13 and M31. Indeed, the poorest performance here is generally observed with nonbinary measures where both the reference and database structures are weighted in the same way, i.e., M22, M33, M44, and M55.

We have seen that measures of the form M1*b* (i.e., those with binary database structure fingerprints) are often better than those of the form Ma1 (i.e., those with binary reference structure fingerprints). When both types of fingerprints are weighted, the measures seem to come in around the middle of the rankings (Table 8) or in the lower half of the rankings (Table 9). That such variations might occur can be demonstrated by consideration of the various components of the Tanimoto coefficient when a similarity is computed; this we shall exemplify using the M13 and M31 measures. As defined previously, the Tanimoto coefficient has the general form:

$$\frac{\sum r_j d_j}{\sum r_j^2 + \sum d_j^2 - \sum r_j d_j}$$

Consider the ECFP₄ fingerprint. Reference to Table 2 shows that the mean W3 weight (3.34 or 3.33 for MDDR and WOMBAT, respectively) is considerably larger than the W1 weight (1.0). When the M13 measure is being used, i.e., when the reference structure fingerprint is weighted and the database structure fingerprint is unweighted, $\sum r_j^2$ will be much

greater than $\sum d_j^2$, and the Tanimoto coefficient will hence be approximately

$$\frac{\sum r_j d_j}{\sum r_j^2 - \sum r_j d_j}$$

In like manner, the M31 Tanimoto coefficient will be approximately

$$\frac{\sum r_j d_j}{\sum d_j^2 - \sum r_j d_j}$$

In a similarity search, the reference structure is matched against each of the database structures in turn, and hence the $\sum r_j^2$ term in the M13 formulation is a constant, C , i.e., the Tanimoto coefficient is approximately

$$\frac{\sum r_j d_j}{C - \sum r_j d_j}$$

When M31 is used, conversely, the $\sum d_j^2$ term in the denominator will vary from molecule to molecule.

The difference in the coefficient values could manifest itself in different search behaviors. This is indeed the case, as exemplified in Figure 2a, which summarizes the similarity values obtained in a search of MDDR for one of the renin reference structures using the two measures with the ECFP₄ fingerprint. The database structures were ranked in decreasing order of the M13 values, the mean similarity computed for each successive set of 1000 structures using both the M13

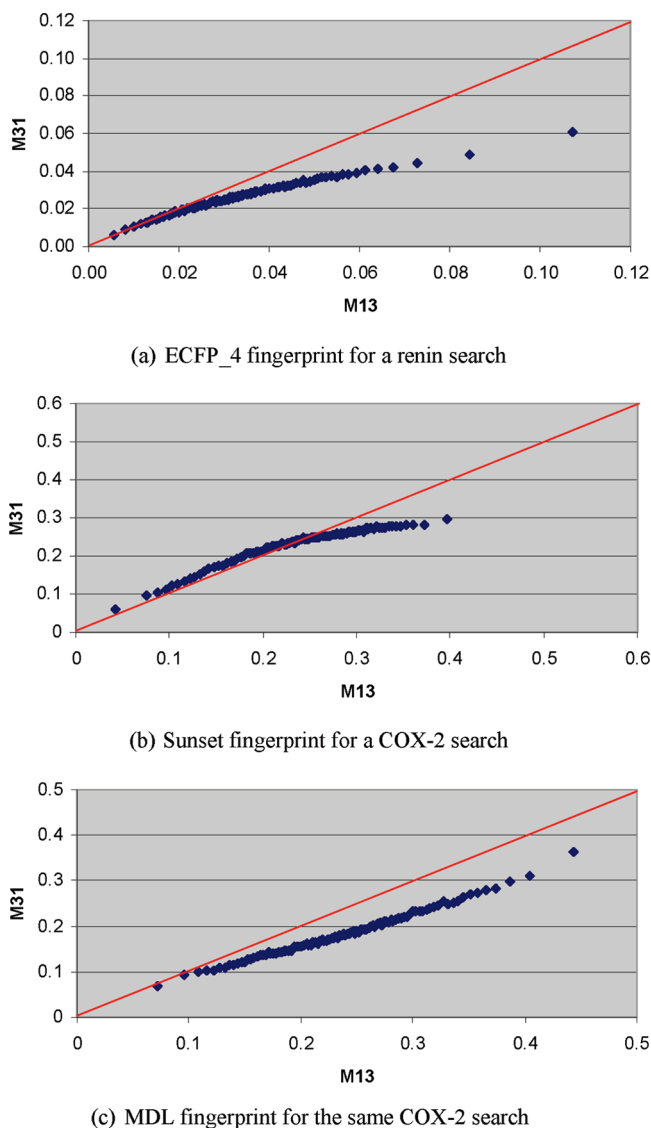


Figure 2. Mean similarity values for M13 and M31 similarity measures. Each point is averaged over a set of 1000 successive structures when the MDDR database is ranked in decreasing similarity order in a similarity search.

and M31 measures and then the two sets of mean values plotted (an entirely comparable plot is obtained if the ranking is in decreasing order of the M31 values). If the two weighted coefficients gave equal values, then the points would lie along the diagonal; instead, it will be seen that the points often correspond to relatively higher M13 values at high similarities, i.e., the molecules at the top of the ranking that determine screening effectiveness and relatively lower M13 values at low similarities. This behavior is observed across all four fingerprints for a wide range of searches, e.g., the Sunset COX-2 search shown in Figure 2b. That said, there are many individual exceptions, e.g., the same COX-2 search but using MDL fingerprints shown in Figure 2c, where the curvature is in the opposite direction.

The analysis of the Tanimoto coefficient given for M13 and M31 will also apply to M15 and M51, since the W3 and W5 weights are by far the largest in the lower section of Table 2. The W2 and W4 weights are closer in magnitude to W1, and hence it might be expected that the differences between M1b and Ma1 could be less well-marked. This is indeed the case. While M12 and M14 tend to perform better

than M21 and M41 overall, there are many exceptions to this general behavior (see Tables 8 and 9).

Figure 2 illustrates a renin search, and inspection of Table 1 reveals that the renin activity class is the least diverse in terms of the mean pairwise similarity for the class. However, analogous plots to those shown in Figure 2 are also obtained in searches for the substance P and protein kinase C classes, which are of intermediate and high structural diversity, respectively. However, while diversity does not appear to affect the relationship between M13 and M31, it certainly does affect that between M11 and M13. Inspection of Table 1a and Table 3 shows that the activity classes where M11 performs well (i.e., where there is shaded cell in the corresponding row of Table 3) have notably lower mean-pairwise similarities than the classes where M13 performs well. The difference in diversity manifests itself in the numbers of retrieved actives, since it is clearly easier to achieve high recall with structurally homogeneous sets of actives than with structurally heterogeneous sets of actives. We have hence divided the MDDR and WOMBAT activity classes into two sets: those with similarities in Table 1 ≥ 0.40 (homogeneous) and those with similarities < 0.40 (heterogeneous). Since there are 10 searches for each activity class, this means that there are 50 and 60 MDDR homogeneous and heterogeneous searches, respectively, and 70 WOMBAT searches for both homogeneous and heterogeneous. We have then taken all the top 1% searches using M11 and M13 and compared the number of actives retrieved to see which of the two similarity measures performed better. For each set of searches (M11 and M13), we have compared the results using both the Sign and Wilcoxon tests. The Sign test simply compares the number of searches where one of the measures was superior to the other, while the Wilcoxon test additionally takes account of the magnitude of the difference in each case. The significance of the differences in each case is assessed using a Z test.³³ The statistically significant results are detailed in Table 10 for each of the four fingerprints inspection of the table reveals a marked difference in that—where the differences are significant—M11 is superior to M13 for the diverse data sets, whereas the converse applies for the homogeneous data sets.

That there might be differences can again be deduced from consideration of the behavior of the Tanimoto coefficient. For M11 the coefficient has the usual form, i.e.,

$$\frac{\sum r_j d_j}{\sum r_j^2 + \sum d_j^2 - \sum r_j d_j}$$

As discussed previously, the $\sum r_j^2$ term in a similarity search will be a constant, call it c in this context. With a diverse set of actives, the similarities are unlikely to be large on average. The contribution from the matching fragments, i.e., the $\sum r_j d_j$ term, is hence likely to be small so that the M11 coefficient is very approximately

$$\frac{\sum r_j d_j}{c + \sum d_j^2}$$

However, with an homogeneous set the reference structure and the active database structures are likely to have many fragments in common so that $\sum r_j d_j$ term is likely to be large

Table 10. Top 1% Searches of Homogenous and Heterogeneous Subsets of MDDR and WOMBAT Using the M11 and M13 Similarity Measures^a

description	searches	M13	M11	ties	test		summary
					Wilcoxon	Sign	
MDDR							
ECFP4-homo	50	34	14	2	$p \leq 0.01$	$p \leq 0.01$	M13 > M11
ECFP4-hetero	60	6	44	10	$p \leq 0.001$	$p \leq 0.001$	M11 > M13
Sunset-hetero	60	12	43	5	$p \leq 0.001$	$p \leq 0.001$	M11 > M13
BCI-homo	50	29	17	4	$p \leq 0.01$	$p \geq 0.05$	<i>M13 > M11</i>
BCI-hetero	60	5	48	7	$p \leq 0.001$	$p \leq 0.001$	M11 > M13
MDL-hetero	60	20	37	3	$p \leq 0.01$	$p \leq 0.05$	M11 > M13
WOMBAT							
ECFP4-hetero	70	22	40	8	$p \geq 0.05$	$p \leq 0.05$	<i>M11 > M13</i>
BCI-hetero	70	15	38	17	$p \geq 0.05$	$p \leq 0.01$	<i>M11 > M13</i>
MDL-hetero	70	19	38	13	$p \leq 0.05$	$p \leq 0.05$	M11 > M13

^a Italicized entries in the summary column denote cases where only one of the two statistical tests indicates a significant difference in performance.

for the top-ranked database structures. Indeed, it may be comparable in magnitude to the $\sum d_j^2$ term, in which case, the M11 coefficient would be very approximately

$$\frac{\sum r_j d_j}{c}$$

We have seen previously that the M13 weight is approximately

$$\frac{\sum r_j d_j}{C - \sum r_j d_j}$$

and this is the expected form with an homogeneous set of actives. However, when they are structurally diverse, the reference and active database structures will generally have fewer fragments in common. The C term in the denominator will then be much larger than $\sum r_j d_j$ and the coefficient will hence be approximately

$$\frac{\sum r_j d_j}{C}$$

It must be emphasized that the approximations here are gross; however, they do suggest that the similarities from the M11 and the M13 schemes may be differentially affected by changes in the diversity of the sets of active molecules that are being sought in a similarity search. That said, marked differences are hardly obvious from curves such as those shown in Figure 3 (which are analogous to those used previously when discussing the relationship between M13 and M31). Figure 3a shows the similarity plot for M11 and M13 ECFP₄ searches of an homogeneous set of actives (the MDDR renin search used in Figure 2), and Figure 3b shows the analogous plot for a search of an heterogeneous set (the MDDR COX-2 search used in Figure 2). There are few obvious differences between the two, and broadly comparable curves are seen with other searches and fingerprints. It may be that curves such as these are at too low a level of granularity to reflect the observed differences in screening effectiveness.

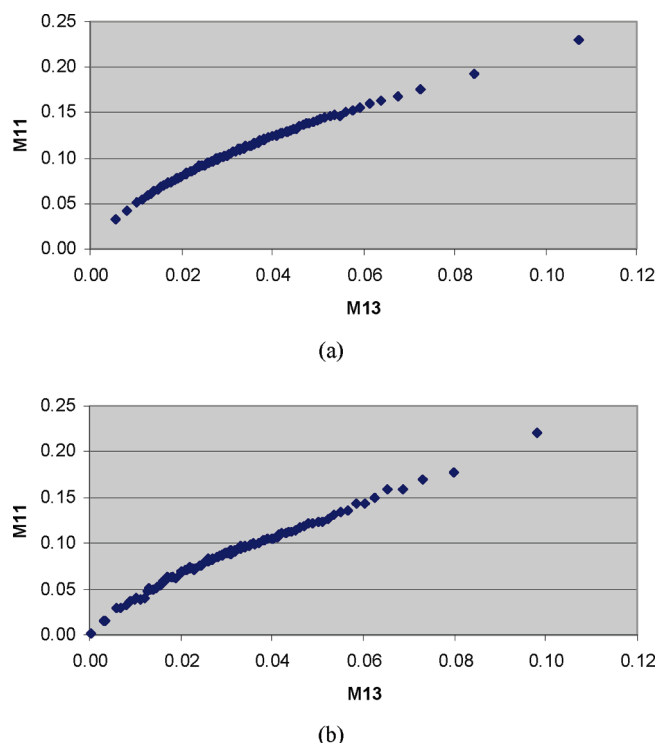


Figure 3. Relationship between M11 and M13 MDDR similarity values using ECFP₄ fingerprints for the renin (a) and COX-2 (b) searches illustrated in Figure 2.

CONCLUSIONS

In this paper, we have carried out a detailed study of the use of fragment occurrence data in similarity-based virtual screening. The results and discussion above suggest that the M13 weight, where the reference structure fingerprint is weighted using the scheme W3 and where each database structure fingerprint uses binary weighting (W1), can provide significantly more effective screening, when averaged over large numbers of similarity searches, than will a conventional fingerprint search in which both types of fingerprint use binary weighting (the M11 measure). The use of inverse frequency weighting might appear to provide a very simple way of increasing the screening effectiveness of operational similarity searching systems, since the weights are calculable from the database as a whole and are applied only to the fingerprint of the reference structure. Unfortunately, however,

this enhanced performance is not observed when structurally diverse sets of actives are being searched. In these more challenging, and more realistic, circumstances, unweighted M11 fingerprints are superior.

In conclusion, our experiments have demonstrated the multiple factors involved in the use of inverse frequency weighting, and hence it is perhaps not surprising that the previous studies summarized in the introduction yielded equivocal conclusions as to the efficacy of this approach to fragment weighting. On the basis of the work reported here, we believe that unweighted fingerprints are the most generally appropriate when the Tanimoto coefficient is used to quantify the degree of resemblance between two fingerprints. It will be of interest to see if this were also the case if alternative types of similarity coefficients were to be used.

ACKNOWLEDGMENT

We thank a referee for valuable comments on an initial draft of this manuscript, Accelrys Inc., Digital Chemistry Limited, Sunset Molecular Discovery LLC, Symyx Technologies Inc., and Tripos Inc. for software and data, the Royal Society and the Wolfson Foundation for laboratory support, and the Government of Malaysia for funding.

REFERENCES AND NOTES

- (1) Willett, P. Similarity-Based Virtual Screening Using 2D Fingerprints. *Drug Discovery Today* **2006**, *11*, 1046–1053.
- (2) Sheridan, R. P. Chemical Similarity Searches: When Is Complexity Justified? *Expert Opin. Drug Discovery* **2007**, *2*, 423–430.
- (3) Eckert, H.; Bajorath, J. Molecular Similarity Analysis in Virtual Screening: Foundations, Limitation and Novel Approaches. *Drug Discovery Today* **2007**, *12*, 225–233.
- (4) Willett, P. Similarity Methods in Chemoinformatics. *Annu. Rev. Inf. Sci. Technol.* **2009**, *43*, 3–71.
- (5) Geppert, H.; Vogt, M.; Bajorath, J. Current Trends in Ligand-Based Virtual Screening: Molecular Representations, Data Mining Methods, New Application Areas, and Performance Evaluation. *J. Chem. Inf. Model.* **2010**, *50*, 205–216.
- (6) Gasteiger, J. *Handbook of Chemoinformatics*; Wiley-VCH: Weinheim, Germany, 2003.
- (7) Leach, A. R.; Gillet, V. J. *An Introduction to Chemoinformatics*. 2nd ed.; Kluwer: Dordrecht, The Netherlands, 2007.
- (8) Arif, S. M.; Holliday, J. D.; Willett, P. Analysis and Use of Fragment Occurrence Data in Similarity-Based Virtual Screening. *J. Comput.-Aided Mol. Des.* **2009**, *23*, 655–668.
- (9) Adamson, G. W.; Bush, J. A. A Comparison of the Performance of Some Similarity and Dissimilarity Measures in the Automatic Classification of Chemical Structures. *J. Chem. Inf. Comput. Sci.* **1975**, *15*, 55–58.
- (10) Harrison, P. J. A Method of Cluster Analysis and Some Applications. *App. Stat.* **1968**, *17*, 226–236.
- (11) Willett, P.; Winterman, V. A Comparison of Some Measures of Inter-Molecular Structural Similarity. *Quant. Struct.-Act. Relat.* **1986**, *5*, 18–25.
- (12) Moock, T. E.; Grier, D. L.; Hounshell, W. D.; Grethe, G.; Cronin, K.; Nourse, J. G.; Theodosiou, J. Similarity Searching in the Organic Reaction Domain. *Tetrahedron Comput. Methodol.* **1988**, *1*, 117–128.
- (13) Downs, G. M.; Poirrette, A. R.; Walsh, P.; Willett, P. Evaluation of Similarity Searching Methods Using Activity and Toxicity Data. In *Chemical Structures 2. The International Language of Chemistry*; Warr, W. A., Ed.; Springer Verlag: Berlin, Germany, 1993; pp 409–421.
- (14) Abdo, A.; Salim, N. Similarity-Based Virtual Screening with a Bayesian Inference Network. *ChemMedChem* **2009**, *4*, 210–218.
- (15) Chen, B.; Mueller, C.; Willett, P. Evaluation of a Bayesian Inference Network for Ligand-Based Virtual Screening. *J. Cheminf.* **2009**, *1* (5), DOI: 10.1186/1758-2946-1-5.
- (16) Goldman, B. B.; Walters, W. P. Machine Learning in Computational Chemistry. *Annu. Rep. Comput. Chem.* **2006**, *2*, 127–140.
- (17) Ormerod, A.; Willett, P.; Bawden, D. Comparison of Fragment Weighting Schemes for Substructural Analysis. *Quant. Struct.-Act. Relat.* **1989**, *8*, 115–129.
- (18) Cosgrove, D. A.; Willett, P. SLASH: A Program for Analysing the Functional Groups in Molecules. *J. Mol. Graphics Modell.* **1998**, *16*, 19–32.
- (19) Robertson, S. E.; Spärck Jones, K. Relevance Weighting of Search Terms. *J. Am. Soc. Inf. Sci.* **1976**, *27*, 129–146.
- (20) Hert, J.; Willett, P.; Wilton, D. J.; Acklin, P.; Azzaoui, K.; Jacoby, E.; Schuffenhauer, A. New Methods for Ligand-Based Virtual Screening: Use of Data-Fusion and Machine-Learning Techniques to Enhance the Effectiveness of Similarity Searching. *J. Chem. Inf. Model.* **2006**, *46*, 462–470.
- (21) Xia, X. Y.; Maliski, E. G.; Gallant, P.; Rogers, D. Classification of Kinase Inhibitors Using a Bayesian Model. *J. Med. Chem.* **2004**, *47*, 4463–4470.
- (22) Bender, A.; Mussa, H. Y.; Glen, R. C.; Reiling, S. Molecular Similarity Searching Using Atom Environments: Information-Based Feature Selection and a Naive Bayesian Classifier. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 170–178.
- (23) Snarey, M.; Terrett, N. K.; Willett, P.; Wilton, D. J. Comparison of Algorithms for Dissimilarity-Based Compound Selection. *J. Mol. Graphics Modell.* **1997**, *15*, 372–385.
- (24) Brown, N.; Jacoby, E. On Scaffolds and Hopping in Medicinal Chemistry. *Mini-Rev. Med. Chem.* **2006**, *6*, 1217–1229.
- (25) Martin, Y. C.; Muchmore, S. Beyond QSAR: Lead Hopping to Different Structures. *QSAR Comb. Sci.* **2009**, *28*, 797–801.
- (26) Robertson, S. E. Understanding Inverse Document Frequency: On Theoretical Arguments for IDF. *J. Doc.* **2004**, *60*, 503–520.
- (27) Harman, D. K. The History of IDF and Its Influences on IR and Other Fields. In *Charting a New Course: Natural Language Processing and Information Retrieval. Essays in Honour of Karen Sparck Jones*; Tait, J. I., Ed.; Springer: Dordrecht, The Netherlands, 2005; pp 69–79.
- (28) Spärck Jones, K. A Statistical Interpretation of Term Specificity and Its Application in Retrieval. *J. Doc.* **1972**, *28*, 11–21.
- (29) Wang, Y.; Bajorath, J. Development of a Compound-Class Directed Similarity Coefficient That Accounts for Molecular Complexity Effects in Fingerprint Searching. *J. Chem. Inf. Model.* **2009**, *49*, 1369–1376.
- (30) Hert, J.; Willett, P.; Wilton, D. J.; Acklin, P.; Azzaoui, K.; Jacoby, E.; Schuffenhauer, A. Comparison of Topological Descriptors for Similarity-Based Virtual Screening Using Multiple Bioactive Reference Structures. *Org. Biomol. Chem.* **2004**, *2*, 3256–3266.
- (31) Hert, J.; Keiser, M. J.; Irwin, J. J.; Oprea, T. I.; Shoichet, B. K. Quantifying the Relationship among Drug Classes. *J. Chem. Inf. Model.* **2008**, *48*, 755–765.
- (32) Gardiner, E. J.; Gillet, V. J.; Haranczyk, M.; Hert, J.; Holliday, J. D.; Malim, N.; Patel, Y.; Willett, P. Turbo Similarity Searching: Effect of Fingerprint and Dataset on Virtual-Screening Performance. *Stat. Anal. Data Mining* **2009**, *2*, 103–114.
- (33) Siegel, S.; Castellan, N. J. *Nonparametric Statistics for the Behavioural Sciences*; 2nd ed.; McGraw-Hill: New York, 1988.

CII001235