

A Structural Analogue Approach to the Prediction of the Octanol–Water Partition Coefficient

Aleksandr Y. Sedykh[†] and Gilles Klopman^{*,‡}

Department of Chemistry, Case Western Reserve University, 10900 Euclid Ave, Cleveland, Ohio 44106, and
MULTICASE Inc., 23811 Chagrin Blvd, Suite 305, Beachwood, Ohio 44122

Received December 5, 2005

A new strategy for the calculation of *n*-octanol/water partition coefficients is presented. Log P calculations of unknown chemicals are based on their closest structural analogues from a database of molecules with known experimental log P values. The contribution of the differing molecular parts is then estimated from a compilation of fragment contributions. Such a strategy is found to be superior to conventional group contribution methods and promises an overall enhancement of the prediction's accuracy.

INTRODUCTION

The importance of the *n*-octanol/water partition coefficient (log P) as an indicator of lipophilicity in biological and medical studies has been well established.^{1,2} Since the pioneering work by Hansch et al.,³ the need for accurate and fast calculation of its value for new molecules has been recognized. Initially, substituent constants were employed to calculate log P values of derivative compounds from a parent structure.^{4,5} However, such a strategy could not handle chemicals with entirely different structures, for which a precursor or appropriate substituent constants are not available. Group contribution methods were thus introduced to allow the calculations to be done for a wider variety of chemicals.⁴ Such methods consisted of breaking a molecule into small fragments, each having an assigned hydrophobic value, equal to the average change produced by the presence of the fragment in the molecules used in the development of the model. Thus, the molecule's log P value could be calculated from these contributions providing that the values of each of the fragments of the molecule have been tabulated. Owing to their simplicity and virtual applicability to any chemical, group contribution methods became extremely popular, and presently, a great number of log P models are based on them (such as ACD/LogP,⁶ KOWWIN,⁷ and KlogP⁸), many of which include both molecular fragments and correction factors representing interactions between them. Such models produce good results within a variety of molecules, for which they were trained, but their predictions for unusual chemicals (for example, molecules with many identical functional groups, whose contributions often fall out of the additivity range of the model and may require complicated corrections) can still be considerably biased.⁹ Because the number of chemicals with experimentally measured log P values is ever increasing, the original idea of calculating log P starting from some known chemical has been recognized as one which could mitigate the shortcomings of conventional group contribution models. The Ex-

perimental Value Adjusted feature of KOWWIN⁷ is a good example of such a procedure. Here, the user supplies a precursor compound with its experimental log P value; then, the log P of the test chemical (normally, a derivative of the precursor) will be predicted by adjusting the precursor's log P value with the contributions of the differing groups. Other efforts to enhance the accuracy of the group contribution models involve optional training with representative chemicals of the structural classes that were not known to the model,¹⁰ resulting in modified contribution values of groups to provide better predictions of those particular classes of chemicals. These techniques, however, require the user's help in supplying relevant experimental data to help find appropriate adjustments for the test chemical. No exploration was apparently undertaken to build, based on diverse sets of chemicals, a model that would automatically select the best precursors for making its predictions.

The obvious advantage of predicting from a known similar chemical is that most of the important functional groups and their interactions are already included in the precursor's experimental log P value, and only few corrections are needed to account for the structural differences. In the past, this was demonstrated in our group¹¹ for substituted benzenes and some other congeneric series. The present work expands this approach to any class of chemicals and offers a new strategy for making high-quality log P predictions without the need of the user's supplied data.

METHODS

The data set which was used to train the log P model reported in our previous work⁸ was also used as a data source for the current study. It consists of over 8000 publicly available chemicals with experimental log P values and is available upon request from the authors.

To build a model out of such data, two steps were necessary:

1. Within the training data set, all pairs of "structurally similar" molecules must be found. This was devised as a two-part process not unlike Stahl and Mauser's clustering strategy for chemical databases:¹²

* Corresponding author fax: 1-216-831-3742; phone: 1-216-831-3740; e-mail: klopman@multicase.com.

[†] Case Western Reserve University.

[‡] MULTICASE Inc.

A. A fragment-based molecular similarity search¹³ is applied to obtain a group of candidates (see section 1A for further discussion).

B. Each candidate is evaluated by finding its maximal common substructure (MCS; the largest structural part shared, by the Bron–Kerbosch clique-detection algorithm¹⁴) with the test chemical¹⁵ (see section 1B for more details).

2. On the basis of the pairs identified in part 1, we also need to optimize the coefficients of the atomic and structural parameters of our model. The adjusted group contribution equation for each such pair of molecules can be written as follows:

$$V(X) = V(Y) + \sum w_i G_i(X-Y) - \sum w_j G_j(Y-X) \quad (1)$$

where $V(X)$ and $V(Y)$ are the log P values of molecules X and Y, respectively (X is being predicted and Y is used as a reference in this case), $G_i(X-Y)$ is the i th structural descriptor present in X and not in Y, $G_j(Y-X)$ is the j th structural descriptor present in Y and not in X, and w_i and w_j are the corresponding coefficients of these descriptors, which are being optimized on the basis of experimental log P values for each XY pair.

Each of these steps (1A, 1B, and 2) requires a special definition of atomic-type (-fingerprint) space. The following atomic features were available as the basis for constructing such atomic fingerprints: nuclear charge, valency, hybridization, coordination number, number of hydrogens, presence of aromaticity, presence of resonance, and various topological information (cyclic, spirocyclic, ring fusions, etc.).

Moreover, step 1A requires a choice of the range and kind of fragments to be used, while step 1B implies the need of a cutoff criteria for what is to be considered a proper “structural analogue”. All of these decisions were subject to trial runs with various configurations and the subsequent selection of the optimum ones.

1A. Fragment-based Similarity Search. To find structurally similar chemicals for a given molecule (X), it needs to be split into all possible fragments (i.e., small substructures of allowed atomic types), which are then checked against the database of fragments, generated likewise out of the training data set. As a result, each molecule (Y) of the training database that shares fragments with the test molecule X can be retrieved. The Dice association coefficient¹⁶ is then used to calculate a similarity score between X and Y:

$$2N_{XY}/(N_X + N_Y) \quad (2)$$

where N_{XY} is the number of fragments occurring both in X and Y, while N_X and N_Y are the total numbers of fragments of X and Y, respectively.

Naturally, the definitions of atomic and fragment types used in the search drastically affect the speed and accuracy of the process, the tendency being that the longer and more-branched fragments (as well as the greater number of atomic features) increase the accuracy but decrease the speed (as the number of fragments grows), which, at its extreme, may render the screening of large groups of compounds impossible. Consequently, many alternative runs were tried to find the optimal definitions. The results are presented in Table 1.

Table 1. Some of the Best Trial Runs of the Fragment Similarity Search

trial ID	atomic fingerprint space ^a	fragment type (sizes/branches) ^b	fragments stored ^c	efficiency, % ^d
9	three-member cycles, ring fusions	3–5/1	22 000	94.36
11	ring fusions, cycles	6/2	43 000	98.29
12		3–6/2	35 000	95.81
13		7/2	46 000	92.87
14		2–7/2	82 000	98.93
15		5–7/2	78 000	98.86
18		2–6/2	36 000	95.16

^a Nuclear charge, valency, hybridization, aromaticity, and resonance were used as atomic features in all of the shown trials. No extra features were used for trials 12–18. ^b “3–5/1” defines a fragment as a chain of three to five atoms and one possible branch. ^c Total number of fragments generated from the training data set under the given definitions of atomic and fragment types. ^d A ratio of the number of “valid” best-scored similarity candidates (generated out of the training data set by the corresponding trial) to the maximum number of such valid best-scored candidates (obtained during these trials). The validity was defined as the need of the maximum common substructure (MCS) between similar chemicals to be more than a certain percentage of their mean size; 80% (strict) and 60% (moderate) MCS values were used to calculate two respective ratios, which were then averaged.

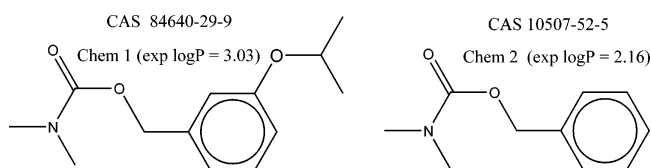


Figure 1. Sample analogues with their experimental log P values.

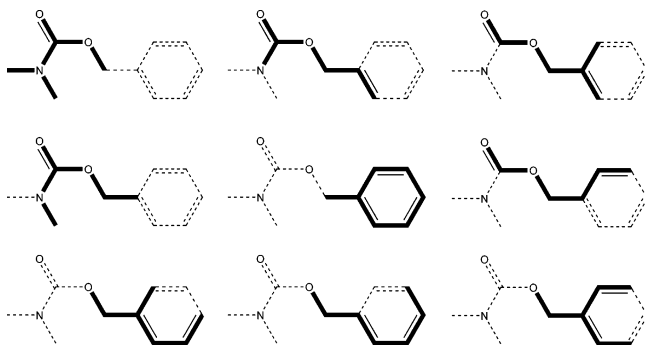


Figure 2. All fragments of size 7 (bold lines, including implicit hydrogens) shared by Chem 1 and Chem 2 of Figure 1.

As can be seen, while the settings of trial 14 yielded the most accurate results, those of trial 11 were only slightly worse but were computed significantly faster, as seen by the efficiency and number of stored fragments for these trials (Table 1).

For a better understanding, the predicting procedure will be demonstrated for the sample chemicals presented in Figure 1.

The molecules of Figure 1 have the following fragments in common (on the basis of trial 14, Table 1), as shown in Figure 2: any shared fragment of smaller size is a part of one of these fragments, with the total number of shared fragments (N_{XY} , eq 2) being 50, the total number of fragments being 100 in Chem 1 (N_X , eq 2) and 51 in Chem 2 (N_Y), and the similarity coefficient being equal to $2 \times 50/(100 + 51) = 0.662$.

1B. Maximal Common Substructure (MCS) Search. An adjusted clique-detection algorithm by Bron and Kerbosch¹⁴

Table 2. Similarity Pairs Generated for Trial 14 under Various Threshold Conditions

I/II ^a	2	4	6	8	no limit
60%	904 [1312]	7116 [5566]	12 418 [6611]	15 350 [6935]	19 342 [7430]
70%	904 [1312]	6934 [5478]	11 873 [6479]	13 904 [6747]	15 684 [6997]
80%	874 [1283]	6527 [5269]	9246 [5991]	9968 [6111]	10 257 [6165]

^a I is the requirement on the MCS size (vertical); II is the number of differing atoms (horizontal). The value reported is the number of unique pairs of two molecules that satisfies all specified similarity conditions (I and II). The value in the brackets is the number of molecules of the training data set involved in these pairs.

was used to obtain the MCSs while validating candidates of the fragment similarity search. Because it is a relatively slow algorithm (it's performance in the worst cases being proportional to $N \times N$, where N is the number of atoms of the two molecules being superimposed¹⁷), all of the available atomic features were used in defining atomic types to produce maximum differentiation. Also, each similarity search requires only a few MCS evaluations because the first few best-scoring candidates exhaust virtually all possible structural analogues. For example, for the five best candidates of trial 14 (Table 1), the first one yields 94.5%, the second 3.2%, the third 1.3%, the fourth 0.6%, and the fifth 0.4% of all successfully validated structures.

Each chemical of the training database was then searched for its analogues (excluding itself), and all unique pairs (XY) of similar molecules (X and Y) were stored for further model training. Two criteria were used in selecting such pairs: the size of the shared part, which is the MCS (in percent of their mean molecular size, a relative parameter; a 60–80% range was tried as a threshold value), and the size of the differing parts (in number of atoms, an absolute parameter; values from 2–8 were employed in the trials).

The choice of these conditions affects both the quality of the log P predictions (large structural differences decrease the accuracy) and the coverage (fewer of the test molecules will have a chance to be predicted if only a small structural deviation is allowed). The outcomes of various configurations of these parameters for trial 14 (Table 1) are shown in Table 2.

2. Adjusted Group Contribution Model. As an initial basis for modeling, we used the same list of atomic and fragment descriptors as that reported in our earlier work.⁸ It consisted of 153 basic parameters and 41 correction factors. However, as our initial runs have shown, many atomic features (in the definition of these parameters) were redundant, such as most of the topological information. A newly redefined atomic-fingerprint space was, thus, selected, consisting only of the nuclear charge, valency, hybridization, coordination number, number of hydrogens, aromaticity, resonance, and presence in a three-membered ring. As a result, the overall number of variables was reduced to 102 basic parameters (of which 61 were atomic and 41 were fragment descriptors) and 36 correction factors, whose coefficients were optimized on the basis of the sets of similarity pairs (Table 2) produced out of the learning data set in the previous stage.

RESULTS AND DISCUSSION

The best modeling results were obtained for the pool of similarity pairs generated with the requirements of 80% for

Table 3. Self-Validation Results of the log P Models

models	R^2	s^a
Current Work		
basic model (138 descriptors)	0.939	0.29
138 descriptors with steric factor	0.940	0.29
Previous Work ⁸		
Model 1 (208 descriptors)	0.915	0.49
Model 2 (208 descriptors with steric factor)	0.922	0.46

^a Standard deviation of differences between calculated and experimental log P values.

Table 4. Prediction of Unknown Chemicals

tests	R^2	s	coverage (%)
Similarity Model			
I ^a	0.930	0.54	51.86
II ^b	0.944	0.51	32.87
III ^c	0.918	0.75	65.75
REF ^d	0.781	0.97	51.86
Previous Work ⁸			
Model 1 ^e	0.880	0.69	100.00
Model 2	0.890	0.65	100.00

^a I: Maximum allowed corrections up to eight atoms or 20% of the structure size, whichever is stricter. ^b II: Molecules with unknown fragments (in the fragment-based similarity search) were not predicted. ^c III: Larger corrections were allowed when predicting: up to 10 atoms or 40% of the structure size, whichever is stricter. ^d REF: No corrections applied, reference molecule's log P was used directly. ^e See Table 3 for these models' description.

the MCS and no more than eight different atoms (whichever was stricter; 9968 pairs, Table 2). The results of the self-validation and a comparison with our earlier-reported group contribution model⁸ are shown in Table 3.

The self-validation of the similarity model was based on the pairs of analogues which were identified from the learning set (Table 2) and then were used for the training of the model. The log P value of each molecule of every such pair was calculated on the basis of its pair partner as the precursor.

As can be seen, the steric hindrance factor used in the previous work⁸ does not bring any significant improvement in the new model. This finding fits the natural expectation that in predicting log P by analogy various molecular property descriptors and topological indices may not be particularly advantageous, because most of such information is already present in the structural analogue.

For the estimation of the model's predictive ability, about 1700 chemicals unknown to the model were tested; this was the same test set used for the validation in the previous study.⁸ The comparative results of the tests are presented in Table 4.

Because the test chemicals were not involved in the creation of the model, a considerable number of them did not have a suitable analogue to be used for prediction. Nevertheless, reduced coverage, but enhanced accuracy, may be the preferable choice, because the chemicals that were left out can still be predicted by a conventional group-contribution method. Noteworthy is the information provided by the REF test in Table 4, which gives an estimate of how much valuable information is already present in the properly chosen structural analogue.

To better demonstrate the predicting procedure, the calculation of Chem 1 from Figure 1 is given in Figure 3.

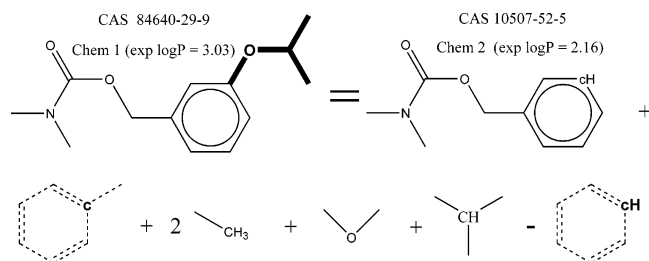


Figure 3. Sample prediction of log P from the structural analogue.

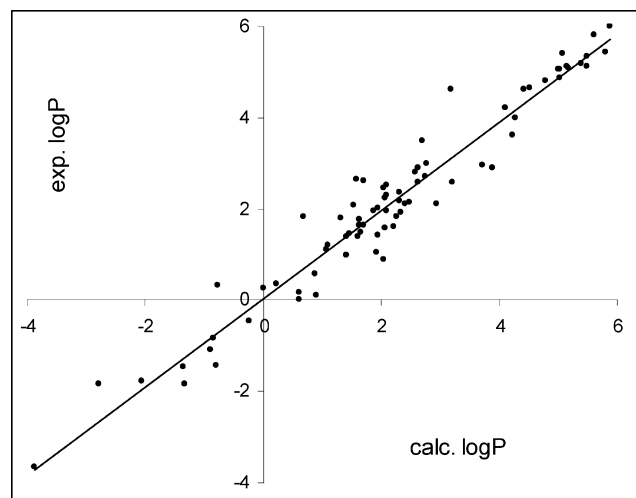


Figure 4. Prediction of “heavy outliers” test set⁹ ($n = 78$, $R^2 = 0.944$, $s = 0.50$).

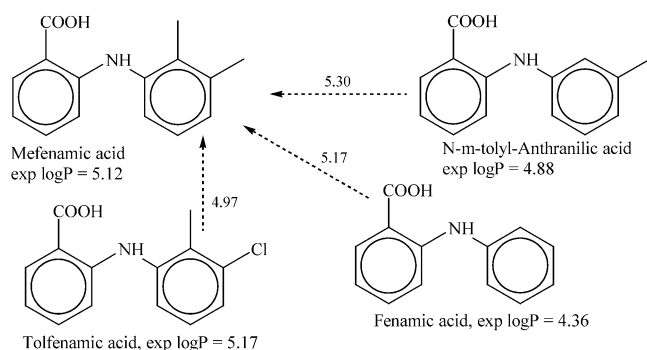


Figure 5. Log P calculation of mefenamic acid. Numbers on the dashed lines are log P values calculated from the corresponding precursors.

The predicted log P of Chem 1 = $2.16\{\text{exp}(\log P) \text{ of Chem 2}\} - 0.267\{\text{cH aromatic}\} + 2 \times 0.452\{-\text{CH}_3\} + 0.359\{-\text{CH}<\} + 0.222\{\text{c-aromatic}\} - 0.358\{-\text{O}- \text{ in resonance}\} = 3.02$.

When modeling, it is necessary to consider the range of possible experimental error present in the data. While the measurement protocol for log P is simple and should yield highly consistent values for simple chemicals, it may not be the case for more complex entities, especially when alternative measuring techniques are employed.¹⁸ A good example of the latter is the work of Benfenati et al.,¹⁹ where log P data for over 200 pesticides were collected from various sources. The average experimental error for their data set was 0.22 log units, which is considerable. This reflects the fact that log P models claiming higher accuracy within such classes of chemicals may be overfitted. To explore such possibility and further validate our model, we employed a specially gathered test set from the work of Eros et al.,⁹ which

Table 5. Chemicals of the “Heavy Outlier Test Set” That Did Not Receive Prediction

Name	Structure	Exp LogP	Best analogue	MCS (%) ^a	Exp ^b LogP
H-2545		1.58		49	4.69
Ascorbic acid		-1.84	N/A	N/A	N/A
Clobazam		2.12	N/A	N/A	N/A
Dextromoramide		3.61		35	4.78
Flumequine		1.6		57	-1.03
KHL 8430		5.8		49	3.00
Papa-verine		2.9		55	0.97
Pilocarpine		0.16		49	2.71
Propafenone		4.63		61	2.48
Ranitidine		0.27		56	2.06
Riboflavin		-1.46	N/A	N/A	N/A
Timolol		1.83		50	1.72
LSD		2.95	N/A	N/A	N/A
Carocainide		1.38	N/A	N/A	N/A
Nicainoprol		1.63		62	0.32

^a Shows (in percentage of the molecular size) how much structure is in common between the test chemical and its best analogue found (see Methods, section 1B). ^b Experimental log P of the analogues reported.

consisted of 78 chemicals reported as “heavy outliers” in the literature. The results are presented in Figure 4.

Table 6. Results of Different Prediction Methods on “Heavy Outlier” Test Set^a

	Clog P	KOWWIN	AUTOQSAR/MLR	AUTOQSAR/PLS	AUTOQSAR/NN	Model 2 ^b	similarity model ^c
I. <i>n</i> = 63 Chems							
<i>R</i> ²	0.914	0.921	0.908	0.899	0.908	0.890	0.964
<i>s</i>	0.63	0.62	0.63	0.67	0.85	0.70	0.40
II. <i>n</i> = 78 Chems (Including Those in Table 5)							
<i>R</i> ²	0.903	0.918	0.902	0.893	0.906	0.884	0.944 ^d
<i>s</i>	0.66	0.61	0.65	0.69	0.86	0.71	0.50 ^d
<i>s</i> ^e	0.66	0.61	0.60	0.60			

^a Compiled from the data of Table 6 of the original work.⁹ I: The subset that was covered by the similarity model. II: The entire data set. ^b Our group contribution model from Table 3. ^c The prediction was based on the same configuration as I of Table 4; exact matches during the search of a precursor were eliminated. ^d Model 2 was used to supply log P values for the chemicals not covered by the similarity model (Table 5). ^e Standard deviation reported in the original work, Table 7.⁹

The predictions were made with the basic model (Table 3) that was allowed to operate on both learning (Table 3) and test data sets (Table 4) of chemicals in order to have the best chance of finding suitable analogues (the requirements used were those of I, Table 4). Exact matches were removed, and if a test chemical had several estimations of its log P value from different but equally satisfactory structural analogues, the average value was used. One such prediction is shown in Figure 5 for mefenamic acid as an example of the procedure.

Out of these 78 chemicals, 63 received predictions from the similarity model (suitable precursors were found), that is, a little over 80% of the test set. The remaining 15 were predicted by our conventional group contribution model (Model 2 in Table 3). These 15 remaining chemicals are shown in Table 5 together with their nearest analogues, wherever possible.

To compare the results of our methodology with each of the models reported in the original Eros publication,⁹ we extracted the literature predictions for the 63 chemicals for which our model found a proper analogue and recalculated the statistical data for them (Table 6, section I). We also calculated the log P value of the 15 chemicals that were outside the domain of our new methodology by our group contribution model (Model 2, Table 3) and combined the results in Table 6, section II.

As can be seen, our new methodology performed better than any of the compared methods, even if it is complemented by our previous group contribution model as shown in Table 6, section II. Also, on the reduced test set (63 compounds, Table 6, section I), all of the compared models with the exception of KOWWIN perform slightly better than on the whole test set of 78 chemicals (Table 7 of the original publication⁹). Nevertheless, KOWWIN shows the best results among the conventional group contribution models both for the full and the reduced “heavy outlier” test sets but trails substantially when compared to our current results.

CONCLUSION

A new approach to make log P predictions by analogy has been developed, and the results (Table 6) were compared with those of known models⁹ and those (Tables 3 and 4) obtained from our previously published group contribution model,⁸ which employed the same training and test sets as well as most of the structural parameters. Superior accuracy makes possible an overall enhancement of the log P predictions.

Another important consideration is that the model is easily extendable to new classes of chemicals by simply adding the representative structures with their experimental log P values to the model’s database. These will serve as structural precursors, and basically, no adjustment of the model’s parameters should be required.

ACKNOWLEDGMENT

The authors thank Dr. Suman Chakravarti and Dr. Hao Zhu for the useful criticism and helpful suggestions in the progress of this work.

REFERENCES AND NOTES

- Leo, A. Octanol/Water Partition Coefficients. *Handbook of Property Estimation Methods for Chemicals*; Boethling, R. S., Mackay, D., Eds.; Lewis Publishers: Boca Raton, FL, 2000; pp 89–114.
- Leo, A.; Hansch, C.; Elkins, D. Partition Coefficients and Their Uses. *Chem. Rev.* **1971**, *71* (6), 525–616.
- Fujita, T.; Iwasa, J.; Hansch, C. A New Substituent Constant, π , Derived from Partition Coefficients. *J. Am. Chem. Soc.* **1964**, *86* (23), 5175–5180.
- Leo, A.; Jow, P. Y.; Silipo, C.; Hansch, C. Calculation of Hydrophobic Constant (log P) From π and f Constants. *J. Med. Chem.* **1975**, *18* (9), 865–868.
- Hansch, C.; Leo, A.; Nikaitani, D. Additive–Constitutive Character of Partition Coefficients. *J. Org. Chem.* **1972**, *37* (20), 3090–3092.
- Petrauskas, A. A.; Kolovanov, E. A. ACD/Log P Method Description. *Perspect. Drug Discovery Des.* **2000**, *19* (Hydrophobicity and Solvation in Drug Design, Pt. 3), 99–116.
- Meylan, W. M.; Howard, P. H. Estimating log P with Atom/Fragments and Water Solubility with log P. *Perspect. Drug Discovery Des.* **2000**, *19* (Hydrophobicity and Solvation in Drug Design, Pt. 3), 67–84.
- Zhu, H.; Sedykh, A.; Chakravarti, S. K.; Klopman, G. A New Group Contribution Approach to the Calculation of log P. *Curr. Comput.-Aided Drug Des.* **2005**, *1*(1), 3–9.
- Eros, D.; Kovesdi, I.; Orfi, L.; Takacs-Novak, K.; Acsady, G.; Keri, G. Reliability of logP Predictions based on Calculated Molecular Descriptors: A Critical Review. *Curr. Med. Chem.* **2002**, *9* (20), 1819–1829.
- Walker, M. J. Training ACD/LogP with Experimental Data. *QSAR Comb. Sci.* **2004**, *23* (7), 515–520.
- Venegas, R. E. Ph. D. Thesis, Case Western Reserve University, Cleveland, OH, 1989.
- Stahl, M.; Mauser, H. Database Clustering with a Combination of Fingerprint and Maximum Common Substructure Methods. *J. Chem. Inf. Model.* **2005**, *45* (3), 542–548.
- Chen, X.; Reynolds, C. H. Performance of Similarity Measures in 2D Fragment-Based Similarity Searching: Comparison of Structural Descriptors and Similarity Coefficients. *J. Chem. Inf. Comput. Sci.* **2002**, *42* (6), 1407–1414.
- Bron, C.; Kerbosch, J. Finding All Cliques of an Undirected Graph. *Commun. Assoc. Comput. Mach.* **1973**, *16* (9), 575–577.
- Rhodes, N.; Willett, P.; Calvet, A.; Dunbar, J. B.; Humblet, C. CLIP: Similarity Searching of 3D Databases Using Clique Detection. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (2), 443–448.
- Kuhns, J. L. The Continuum of Coefficients of Association. *Proceedings of the Symposium on Statistical Association Methods For*

- Mechanized Documentation*; Stevens, M. E., Giuliano, V. E., Heilprin, L. B., Eds.; National Bureau of Standards: Washington, DC, 1965; pp 33–39.
- (17) Bunke, H.; Guidobaldi, C.; Foggia, P.; Sansone, C.; Vento, M. *A Comparison of Algorithms for Maximum Common Subgraph on Randomly Connected Graphs*; Caelli, T., Amin, A., Duin, R., Kamel, M., de Ridder, D., Eds.; Structural, Syntactic, and Statistical Pattern Recognition, Joint IAPR International Workshops SSPR 2002 and SPR 2002, Windsor, Ontario, Canada, August 6–9, 2002; Springer: Berlin, 2002; pp 123–132.
- (18) Paschke, A.; Neitzel, P. L.; Walther, W.; Schueuermann, G. Octanol/Water Partition Coefficient of Selected Herbicides: Determination Using Shake-Flask Method and Reversed-Phase High-Performance Liquid Chromatography. *J. Chem. Eng. Data* **2004**, *49* (6), 1639–1642.
- (19) Benfenati, E.; Gini, G.; Piclin, N.; Roncaglioni, A.; Vari, M. R. Predicting logP of Pesticides Using Different Software. *Chemosphere* **2003**, *53* (9), 1155–1164.

CI0505269