

## Using Molecular Fingerprint as Descriptors in the QSPR Study of Lipophilicity

Ruifeng Liu\* and Diansong Zhou

Department of Chemistry and Department of Drug Metabolism and Pharmacokinetics,  
AstraZeneca, 1800 Concord Pike, Wilmington, Delaware 19850

Received October 18, 2007

Using SciTegic's extended connectivity fingerprint as raw descriptors, a robust partial least-squares model for logP prediction was developed. The PLS model is based on 39 latent variables. An additional 8 correction factors are employed to account for effects such as intramolecular hydrogen bonding. The model performs similarly to ClogP for compounds with molecular weight in the 250–400 range but significantly better than ClogP for molecules with molecular weight over 400. Considering modern drug discovery tends to generate larger candidate compounds, the PLS model is better suited for drug discovery applications. The good performance of the simple PLS model indicates that the molecular fingerprints encode detailed structure information. When used properly they outperform conventional descriptors in QSPR model development.

### INTRODUCTION

Lipid solubility of a compound is of special importance to drug discovery and development, because it is directly related to the transport abilities of a drug candidate to cross biological membranes.<sup>1,2</sup> The requirement is that drug molecules must be soluble enough in lipid to get into membranes but cannot be so soluble that they become trapped in the membranes. A simple parameter for evaluating lipid solubility is *n*-octanol/water partition coefficient. Many quantitative structure–activity/property relationship (QSAR/QSPR) studies have shown that the logarithm of partition coefficient (logP) is one of the most important descriptors.<sup>3</sup> There are many methods for logP measurement,<sup>4</sup> but the experimental approaches may be time-consuming, costly, and unreliable for highly lipophilic or highly hydrophilic compounds.<sup>4</sup> More importantly, in modern drug discovery it is important to know the lipophilicity of a carefully designed compound before it is synthesized. As a result, computational models that can give reliable logP predictions are invaluable.

Since the pioneering work of Hansch and Fujita<sup>5</sup> in early 1960s, many computational models for logP prediction were developed.<sup>6</sup> Among them, the ClogP model<sup>7</sup> of Leo and Hansch is regarded as the “gold” standard. This is because the ClogP model was based on a large and carefully evaluated data set. In a recent release, ClogP version 4 was parametrized on experimental data of over 12 000 compounds, namely the MedChem StarList.<sup>8</sup> The model was based on over 150 carefully selected molecular fragments, whose contributions to logP are assumed to be additive, and around 600 dependably measured molecular descriptors. Leo and Hoekman reported that the model reproduces the experimental logP values of the StarList with a mean deviation of  $\pm 0.31$  log unit.<sup>7</sup> Many other logP models were either trained and tested on much smaller data sets or trained on data sets of similar size but could not achieve the same level of accuracy.

To derive a robust QSPR model, one needs to employ both a sound statistical method and a set of molecular descriptors

that are easy to calculate and give a detailed description of molecular structures. Over the years, many statistical methods, ranging from simple multilinear regression to much more complicated artificial neural networks and support vector machines, have been applied in QSAR/QSPR studies. Meaningful and relevant molecular descriptor selection remains a very crucial step for successful model development.

Many logP prediction models used predefined atom types and/or molecular fragments as molecular descriptors.<sup>7,9–12</sup> The contribution of each fragment to the logP of a molecule is determined, often from experimental values of small molecules. To predict the logP value of a new molecule, the number and types of predefined fragments present in the molecule are counted. The predicted value is the sum of the fragment contributions. The performance of the ClogP<sup>7</sup> and AlogP<sup>9</sup> models indicates that the fundamental assumption of this approach, the atom/fragment additivity of logP contributions, is valid. However it is very tedious and requires a lot of expertise in devising and developing a model based on this approach. If too small basic fragments are defined, atomic environments within a molecule and long-range interactions cannot be properly accounted for. When larger fragments are defined, the number of distinct fragments will be high, which will significantly increase the difficulty of deriving reliable fragment contributions.

Over the years, many generic descriptors were developed for various QSAR/QSPR applications. These include molecular surface area-based descriptors,<sup>13,14</sup> molecular connectivity and shape-based descriptors,<sup>15,16</sup> and molecular electrotopological state indices,<sup>17</sup> etc. These descriptor values can be easily calculated from molecular structures, yet in most cases the abstracted descriptors cannot give a sufficiently accurate description of molecular structure. For example there are cases of molecules with significantly different chemical bond connectivity but having similar descriptor values. As a result, these descriptors do not perform well as a basis for 2D molecular similarity search retrieving structurally similar compounds.

\* Corresponding author e-mail: Ruifeng.Liu@astrazeneca.com.

Molecular fingerprints can be considered as another class of molecular descriptors. They have been successfully used in molecular similarity searches,<sup>18</sup> indicating they can give more accurate description of molecular structures. However they are not generally used as descriptors in QSAR/QSPR studies due to a few reasons. First, traditional molecular fingerprints encode the presence or absence of a structural feature only, not how many times a structural feature is present in a molecule. Second, many high performance molecular fingerprints are based on a large number of overlapping structural features. Even for a set of very small molecules, the number of basic structural features represented by the on-bits of molecular fingerprint could be much higher than available experimental observations. Third, because high performance molecular fingerprints are defined as overlapping structural features, fingerprint bits are strongly correlated. Due to these difficulties, few published QSAR/QSPR studies used molecular fingerprints as descriptors.

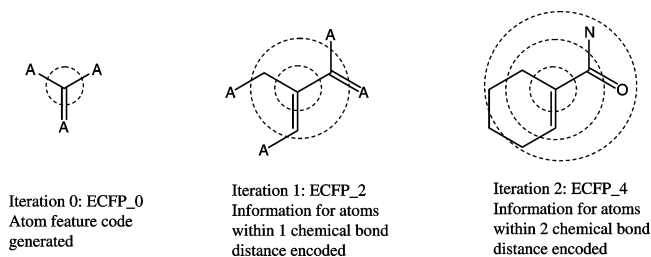
The present work is aimed at exploring the feasibility of using modified molecular fingerprints as generic descriptors for improving logP prediction. This was based on the thinking that since molecular fingerprints give more accurate description of molecular structures, they should be better descriptors for QSAR/QSPR studies. We were encouraged by a recent paper of Sun<sup>12</sup> who devised a scheme, based on Daylight fingerprints,<sup>19</sup> for selecting relevant atom types for logP prediction. He selected a total of 218 atom types and 26 correction factors to account for long-range effects. To overcome descriptor correlation, Sun applied a partial least-squares (PLS) technique and derived a PLS model with nine latent variables. His model was able to predict the logP values of a training set of 10 850 StarList compounds with an  $r^2$  of 0.912 and an rms error of 0.51 log unit. This is impressive as the total number of independent variables (latent variables) is much fewer than those used in any of the commercial packages, yet it achieved similar performance in terms the correlation coefficient between the calculated and experimental results.

## COMPUTATIONAL DETAILS

**1. Statistical Method and Descriptor Selection.** A fundamental assumption of QSAR/QSPR is that molecular properties are a function of molecular structure. Thus once a molecular structure is known, in principle, one should be able to predict its property/activity. In this sense, molecular fingerprints should be good descriptors because they give more detailed description of molecular structure than generic descriptors usually used in QSAR/QSPR studies. In the present study, we choose to explore using the scores of SciTegic extended connectivity fingerprint<sup>20</sup> (ECFP) features as molecular descriptors for logP prediction. The score of a molecular fingerprint feature is defined as the count of how many times the fingerprint feature presents in a molecule. The ECFP fingerprint was chosen because it was found to perform better when used in similarity searching to retrieve molecules with similar biological activity.<sup>21</sup>

Briefly speaking, ECFP fingerprints are generated iteratively to encode features that represent each atom in larger and larger structural neighborhoods. Figure 1 illustrates this process.

At iteration 0 (ECFP\_0), the information of individual atoms is encoded. It includes the number of connections

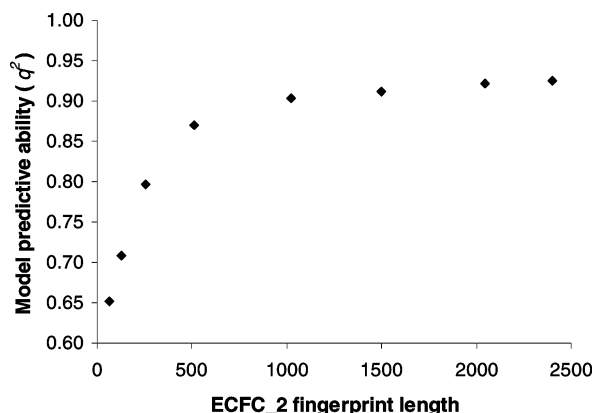


**Figure 1.** Schematic illustration of ECFP<sub>n</sub> fingerprint generation. Each iteration adds bits that represent larger and larger structural features.

(bonds) to the atom, element type, charge, and atom mass. At iteration 1, the information of all atoms directly bonded to the atom (within a diameter of 2 chemical bonds, and hence termed ECFP\_2) is encoded. At iteration 2 (ECFP\_4), the information of all atoms within a diameter of 4 chemical bonds is also encoded. When the desired neighborhood size is reached, the process is complete, and the set of bits representing all features for the atom is returned as part of the molecular fingerprint. This process is repeated for all the atoms in a molecule. The molecular ECFP<sub>n</sub> fingerprint is a collection of all the bits representing atoms in their molecular neighborhoods. Similar to other fingerprints, ECFP<sub>n</sub> in this form is not appropriate for our purpose as it encodes whether a structural feature is present in a molecule only, not how many times this structural feature is present in the molecule. An extension of the fingerprint, the score of ECFP<sub>n</sub>, is used in this study as molecular descriptors. It is termed the ECFC<sub>n</sub> fingerprint.

As one can imagine, with increasing  $n$ , the number of unique ECFP<sub>n</sub> features increases dramatically. To make the computational task manageable, the molecular fingerprint is usually folded into a computationally manageable length by the logical OR function. However, the folding leads to a loss of structure information to a certain degree. In the current study, we used the ECFC\_2 fingerprint folded to a fixed length of 1024 bits as raw molecular descriptors. This fingerprint length was found to be satisfactory for our purpose because the structural features represented by the ECFP\_2 fingerprint bits is small; therefore, the total number of unfolded ECFP\_2 fingerprint bits are small. For most of the molecules in the StarList, the number of on-bits of the unfolded ECFP\_2 fingerprint is well below 100. However, due to the rich structural diversity of the StarList, the total number of unique on-bits of the whole molecular set is over 1024. We performed model building to search for an appropriate fingerprint length. The results (shown in Figure 2 and details given in the Results and Discussion section) indicate that 1024 bits is a good compromise between the computational cost and the prediction reliability. We also tried using the ECFC\_4 fingerprint folded to the same length, but it did not lead to improvement over the ECFC\_2 fingerprint, probably because after folding more severely, too much structural information is lost. Since the unfolded ECFC\_2 fingerprint is much shorter (fewer unique fingerprint features than the ECFC\_4), it is folded less than ECFC\_4 to achieve the 1024 bits.

To solve problems associated with the large number of descriptors and descriptor correlation, the kernel-based PLS<sup>22</sup> method as implemented in the R packages<sup>23</sup> was used to generate regression models. Cross-validation was used to find



**Figure 2.** Predictive ability of the PLS model as a function of fingerprint length, derived from leave 10% out cross-validation on the training set.

the optimal number of latent variables. The final model was derived with the optimal number of latent variables. The PLS approach is known to handle noise in the data set very well.<sup>24</sup>

**2. Experimental Data Set.** To build a quality QSPR model, one needs a quality training set of significant size. Thanks to the efforts of Hansch and Leo, a large collection of measured logP data, the MedChem StarList, was assembled. The logP values in the StarList were carefully evaluated by Leo et al. and deemed to be reasonably reliable. ClogP and some other commercial logP prediction models used the StarList as the training set. In the current study, we used the StarList distributed by Daylight Chemical Information Systems in its Toolkit version 4.82 release. There are a total of 10 974 SMILES strings with logP values. However not all of them are organic compounds. For example, it includes gases such as hydrogen, oxygen, nitrogen, carbon monoxide, carbon dioxide, helium, argon, etc., and water as well as mercury. It also includes some organometallic compounds. In order to build a logP model for organic compounds, we removed all of the nonorganic compounds. We also removed the organometallic compounds, as the number of them is too few to reliably establish the contribution of fingerprint features of the organometallic moieties. A total of 119 entries were removed from the StarList due to these reasons.

The StarList was further cleaned to standardize structure presentation. First, for organic salts (disconnected cations and anions), we retained the organic cationic parts and removed the anionic counterions (mostly  $\text{Cl}^-$ ,  $\text{Br}^-$ , and  $\text{I}^-$ ). This is because some of the entries in the StarList are organic cations only; the counterions are not given. Second, the charges of all zwitterions except amino acids were neutralized. This is because the ECFC\_n fingerprint features are dependent on the molecular structure presentation. Some compounds presented as a zwitterion may not have the extent of charge separation as represented by the zwitterionic form. All nitro groups in the Starlist were represented by  $[\text{N}^+](=\text{O})(\text{O}^-)$ , but in reality the two oxygen atoms are equivalent. Even though the nitrogen–oxygen bond is polar, it is not as polar as a complete positive–negative charge separation. In the structure standardization process, the nitro group was converted into the neutral  $\text{N}(=\text{O})(=\text{O})$  form. Third, amino acids and all compounds that have both a carboxylic group and an amino group bonded to a nonaromatic carbon are ionized into the zwitterion form. This is because

the zwitterionic form is very likely the dominant form of these compounds in aqueous solution. We also similarly ionized amino sulfonic acids and amino sulfuric acids.

After structure standardization, the 10 855 compounds were clustered using the ECFC\_2 fingerprints. The clustering was based on maximal Tanimoto dissimilarity partitioning. Compounds in each cluster are at most within 0.20 Tanimoto distance from its cluster center. The cluster centers were selected sequentially as the compounds having the highest number of neighbors within 0.20 Tanimoto distance. Once a cluster center is identified, all members of the cluster are excluded from further selection of cluster centers. This process is repeated until no more clusters can be identified. This generated 1994 clusters with 2 to 21 compounds each and 5257 singletons. All the cluster centers and singletons were selected into the training set. Another 2518 compounds were randomly selected from the remaining compounds to the training set so that the size of the training set is 90% of the data set. The remaining 10% of the compounds (1086) were kept away from the training and used as a test set. It should be noted that neither the clustering algorithm nor the ECFC\_2 fingerprint is ideal for meaningful clustering of compounds. The purpose of the clustering here is to select a structurally diverse training set.

## RESULTS AND DISCUSSION

All calculations were performed using SciTegic's Pipeline Pilot. The first step is to explore the appropriate ECFC\_2 fingerprint length that balances accuracy and computational cost. To achieve this, we calculated ECFC\_2 fingerprints of increasing length and applied the partial least-squares (PLS) technique to the training set compounds. PLS combines principal component analysis (PCA) and least-squares regression to extract information relevant to the property under investigation and reduce dimensionality. This is done by combining raw descriptors via a linear combination to produce orthogonal latent descriptors. Statistical analysis is applied to identify the most relevant latent descriptors. The most relevant descriptors are then used to build a regression model. In the present study, leave 10% out cross-validation was performed. This was done by splitting the training set randomly into 10 equal-size data sets. Nine of them were used to build a model. The model was used to predict logP of the tenth set left out in model building. The deviation between the predicted and the experimental logP values were calculated for the left-out set. This process is repeated until each and every set is left out once. Root-mean-square error (RMSE) of prediction and estimate of predictive ability,<sup>25</sup>  $q^2$ , were calculated. The value of  $q^2$  is a good statistical criterion of the quality of a predictive model. For the training set of 9769 compounds,  $q^2$  values of the PLS models with increasing length of the ECFC\_2 fingerprint as raw descriptors are given in Figure 2. It shows that with too short of a fingerprint length, the ECFC\_2 fingerprint descriptors do not encode sufficient structural information for a high predictive ability model to be developed, because of information loss due to heavy folding. However further increasing the fingerprint length beyond 1024 bits does not improve model predictive ability significantly as  $q^2$  apparently reached a plateau. On the other hand, the computational cost increases tremendously with the increasing size of the fingerprint. As



**Table 1.** Summary of Statistical Parameters in PLS Model Building

number of latent variables	%descriptor variance covered	%logP variance explained	RMSE	$q^2$	number of latent variables	%descriptor variance covered	%logP variance explained	RMSE	$q^2$
1 LV's	22.3	26.9	1.408	0.267	36 LV's	79.9	92.7	0.510	0.904
2 LV's	28.0	54.3	1.116	0.539	37 LV's	80.1	92.8	0.510	0.904
3 LV's	34.2	64.9	0.979	0.645	38 LV's	80.3	92.8	0.509	0.904
4 LV's	37.5	71.0	0.892	0.706	39 LV's	80.5	92.9	0.509	0.904
5 LV's	40.9	75.3	0.826	0.748	40 LV's	80.8	92.9	0.509	0.904
6 LV's	45.5	78.1	0.779	0.775	41 LV's	81.1	92.9	0.509	0.904
7 LV's	48.7	80.0	0.747	0.794	42 LV's	81.3	93.0	0.510	0.904
8 LV's	52.7	81.5	0.719	0.809	43 LV's	81.5	93.0	0.510	0.904
9 LV's	56.2	83.4	0.686	0.826	44 LV's	81.7	93.0	0.511	0.904
10 LV's	59.4	84.5	0.664	0.837	45 LV's	81.9	93.1	0.511	0.904
26 LV's	76.1	91.8	0.521	0.900	46 LV's	82.1	93.1	0.512	0.904
28 LV's	77.0	92.1	0.517	0.901	47 LV's	82.3	93.1	0.512	0.904
30 LV's	77.8	92.3	0.514	0.903	48 LV's	82.6	93.1	0.512	0.903
32 LV's	78.5	92.4	0.512	0.903	49 LV's	82.8	93.1	0.513	0.903
34 LV's	79.2	92.6	0.511	0.904	50 LV's	83.0	93.2	0.513	0.903

a result, we decided to use a fixed ECFC\_2 fingerprint length of 1024 bits.

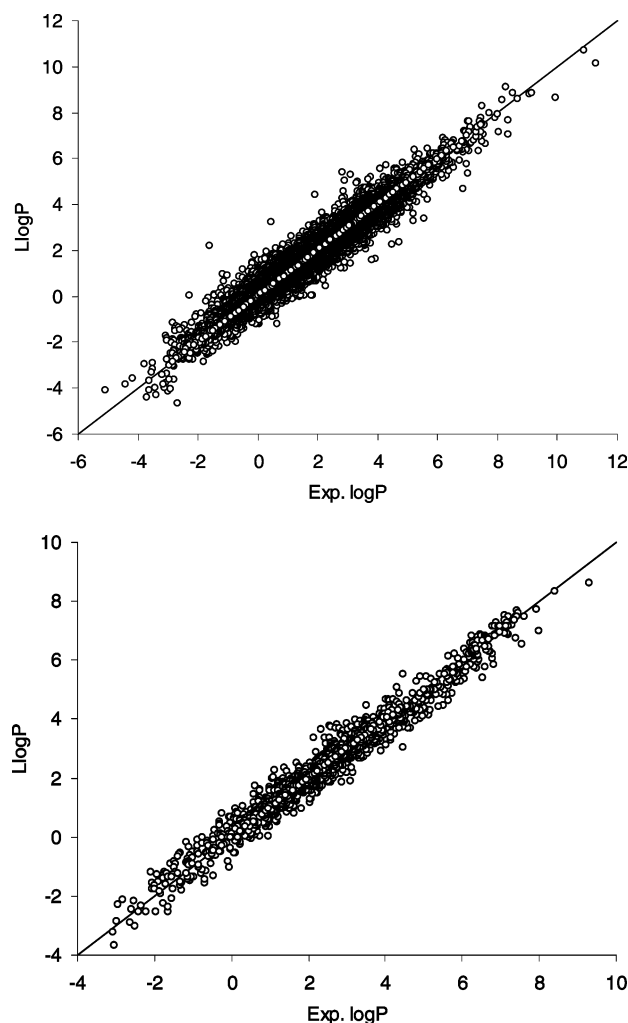
For the training set of 9769 compounds, the calculation indicates that 43 of the 1024 fingerprint bits are off for all the molecules (i.e., the structural features these bits represent do not appear in any of the compounds). These 43 bits were removed because their values are zeros for all the molecules (zero variance descriptors). Thus the total number of raw descriptors used in this study is 981.

Table 1 is a summary of RMSE and  $q^2$  from the cross-validation calculation with the training set. It also shows percent variance of descriptors covered by the latent variables and percent variance of logP data the latent variables explain. With an increasing number of latent variables, the %variance in raw descriptors covered by the latent variables increases, so does the %variance in measured logP explained by the latent variables. RMSE of prediction decreases with an increasing number of latent variables, reaching a bottom with 38 latent variables, and starts to increase from 41 latent variables and on. In the meantime,  $q^2$  increases with an increasing number of latent variables, reaching a plateau with 34 latent variables, and starts to drop from 47 latent variables. When all factors are considered, it seems the top 39 latent variables are an optimal set of descriptors for model building. With these 39 latent variables, nearly 93% of the variance in experimental logP data is explained, and RMSE is at minimum while  $q^2$  is at maximum.

On the basis of  $q^2$ , RMSE, and coverage of logP variance, a PLS model using the top 39 latent variables was constructed based on the training set compounds. For the training and the test set compounds, the logP calculated by the partial least-squares model (LlogP) are compared with the experimental values in Figure 3. It shows that the 39-descriptor PLS model performs reasonably well. The correlation coefficients between the calculated and the experimental logP values are 0.96 and 0.98 for the training and the test sets, respectively. The RMSE of prediction are 0.45 and 0.41 log units, respectively. These statistical parameters compare favorably with the logP calculated by ClogP version 4.3 and significantly better than those calculated by AlogP98<sup>9</sup> as shown in Table 2. To ascertain statistical soundness of the PLS model, we performed y-scrambling tests on the training set. With randomly scrambled experimental logP values, a PLS model with top 39 latent descriptors was built, and the correlation coefficient between the model predicted and

training set logP values is calculated. For five y-scrambling tests, the squared correlation coefficients range from 0.08 to 0.10. This contrasts sharply to the squared correlation coefficient of the unscrambled training set, 0.92, indicating convincingly the statistical soundness of the PLS model.

Careful examination of the calculated and experimental logP values indicates some systematic deviations. Most of the observed systematic deviations seem to occur for molecules with high likelihood of intramolecular hydrogen



**Figure 3.** Calculated (LlogP) versus experimental logP of the training (upper) and the test set (lower) compounds.

**Table 2.** Statistical Parameters of the Training and the Test Set Compounds

	training set (9769 comps)				test set (1086 comps)		
	LlogP	ClogP	AlogP		LlogP	ClogP	AlogP
R	0.96	0.97	0.92	R	0.98	0.98	0.96
RMSE	0.45	0.40	0.64	RMSE	0.41	0.46	0.62
max. dev.	3.83	11.44	7.39	max. dev.	1.44	2.46	2.35
# dev. > 2	14	21	100	# dev. > 2	0	4	3

**Table 3.** Difference between the Experimental and the Calculated logP of Compounds with Certain Structure Features

substructure <sup>a</sup>	comps	low diff.	high diff.	mean diff.	# diff. > 0	correction
SS1	103	-1.56	2.41	0.65	92	0.93
SS2	59	-0.55	1.17	0.36	50	0.54
SS3	25	-0.66	0.95	0.25	21	0.35
SS4	54	-0.50	2.22	0.55	48	0.83
SS5	29	0.21	0.99	0.71	29	0.86
SS6	94	-0.78	1.19	0.41	83	0.63
SS7	126	-0.87	1.31	0.44	117	0.56
SS8	28	0.35	2.22	0.97	28	1.13

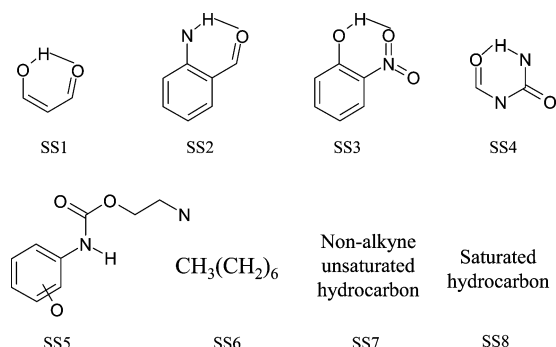
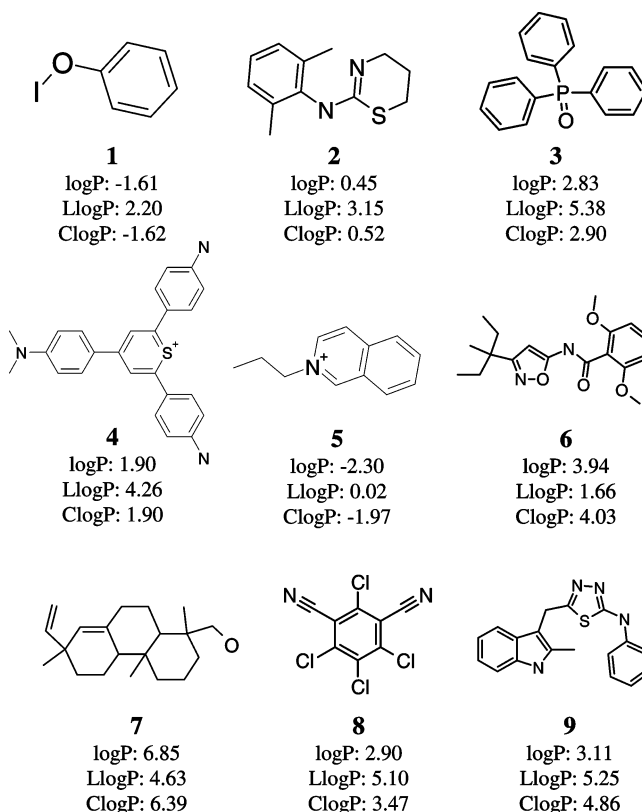
<sup>a</sup> The substructural features are defined in Figure 3.

bonding. The result is that our calculated values are systematically too low. It was also found that logP of some hydrocarbons are systematically underestimated. Details of systematic deviations for compounds with some characteristic structural features are given in Table 3. The characteristic structural features are shown in Figure 4.

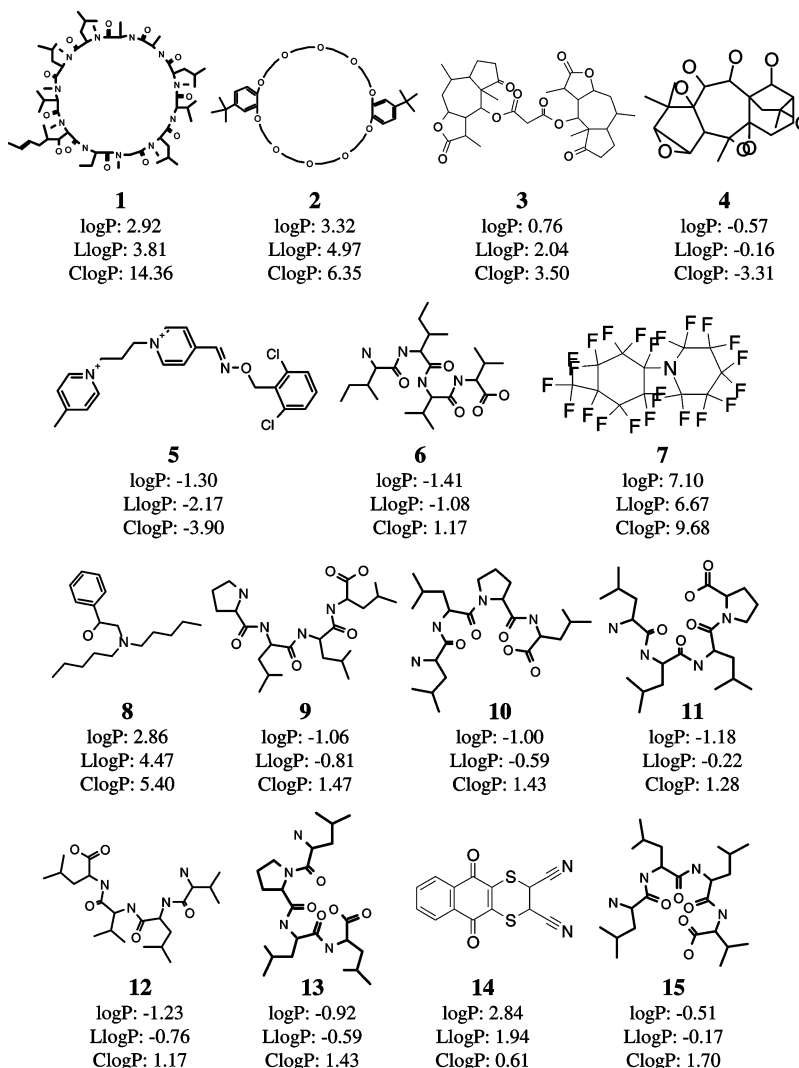
SS1 is a carbonyl group *cis* to a hydroxyl group. It matches a carbonyl ortho to a hydroxyl group on an aromatic ring. Table 3 shows that in the StarList, there are 103 compounds with this substructure feature. The difference between the experimental and our calculated logP values ranges from -1.56 to 2.41 log units. But there are only 11 of the 103 compounds with calculated logP higher than the experimental values, in sharp contrast to 92 compounds being underestimated by the PLS model. For all compounds with this structural feature, the PLS model underestimated logP by an average of 0.65 log units. Results for compounds with other substructure features in Figure 4 are similar. For SS1 to SS4, it is highly likely that intramolecular hydrogen bonding is responsible for the observed systematic deviation. The hydrogen bond donors and acceptors in these substructures are ideally situated for hydrogen bond formation. Intramolecular hydrogen bonding may also contribute significantly to the systematic deviations observed for compounds with SS5 when the oxygen atom is ortho to the HN group (15 compounds), but 14 compounds with oxygen atoms meta to the NH group were also found to be similarly underestimated by the PLS model. SS6 is a long saturated hydrocarbon chain substructure. A total of 94 compounds in the StarList have this structural feature, the PLS model underestimates the logP values of 83 of them. SS7 and SS8 are saturated hydrocarbons and non-alkyne unsaturated hydrocarbons.

Interestingly, for alkynes without any heteroatoms, there seems to be no systematic deviation between the calculated and experimental logP values. However, there are a total of only 9 such compounds in the StarList, which may not be sufficient for a statistically sound conclusion.

The structural features in Figure 4 suggest that most of the observed systematic deviations are associated with unique molecular substructural features that are large in size (mostly

**Figure 4.** Structural features associated with systematic deviations between the experimental and the calculated logP values.**Figure 5.** Molecular structures of compounds with LlogP error over 2 log units.

involving at least five heavy non-hydrogen atoms) that cannot be properly described by the small ECFC\_2 fingerprint. For example, the main contributing factors of SS1 to SS5 are intramolecular hydrogen bonding which requires the proton and the acceptor to be five bonds away. The ECFC\_2 fingerprint covers a distance of only two chemical bonds. The deviation associated with SS6 seems to be related to similarity to n-octanol tail. Based on the principle of “like” dissolves “like”, it is understandable that compounds with long saturated CC chains have higher n-octanol affinity. Deviations associated with SS7 and SS8 are likely due to the overall polarity of the molecules. Most of the compounds in StarList have heteroatoms and therefore are polar molecules. The hydrocarbons represented by SS7 and SS8 are nonpolar compounds. As water is a highly polar solvent, hydrocarbons may be more hydrophobic than can be predicted as a sum of atomic fragment contributions because of overall nonpolarity of the compounds. To properly take into account of these effects, one needs to use a fingerprint



**Figure 6.** Molecular structures of compounds with high ClogP errors.

at least the size of ECFC\_4 (diameter of 4 chemical bonds). However the ECFC\_4 fingerprint has too many distinct bits for the StarList compound set. In order to make the computational task manageable, one has to fold the fingerprint significantly. The folding loses too much structural information. As a result, our test calculation using ECFC\_4 folded to 1024 bits did not show improvement over the ECFC\_2 fingerprint.

As the ECFC\_2 fingerprint descriptors are not able to account for the effects associated with the structural features described in Figure 4, and the least-squares procedure fits all the data together, we estimate that the effects on logP are at least higher than the mean deviations given in Table 3. The exact values of these effects are unknown. To approximately account for these effects under the PLS framework without increasing the fingerprint size, we tried the following procedure. First, generate approximate special effect free "experimental" data by subtracting the respective mean deviations from the logP values of compounds with the corresponding structural features. Second, retrain the PLS model using the modified experimental logP data with the hope of deriving a more robust model for compounds without the characteristic structural features in Figure 4. Note that the new model should underestimate the logP values of compounds with the structural features in Figure 4 even more

than before. Third, compare the calculated and original experimental logP of compounds with the structural features. Use the new mean deviations between the experimental and calculated values as correction factors for compounds with the corresponding structural features. Fourth, the model predicts logP of a compound as a sum of logP calculated from the PLS model and the correction factors where applicable. The correction factors derived from this procedure are given in Table 3 under column name "correction".

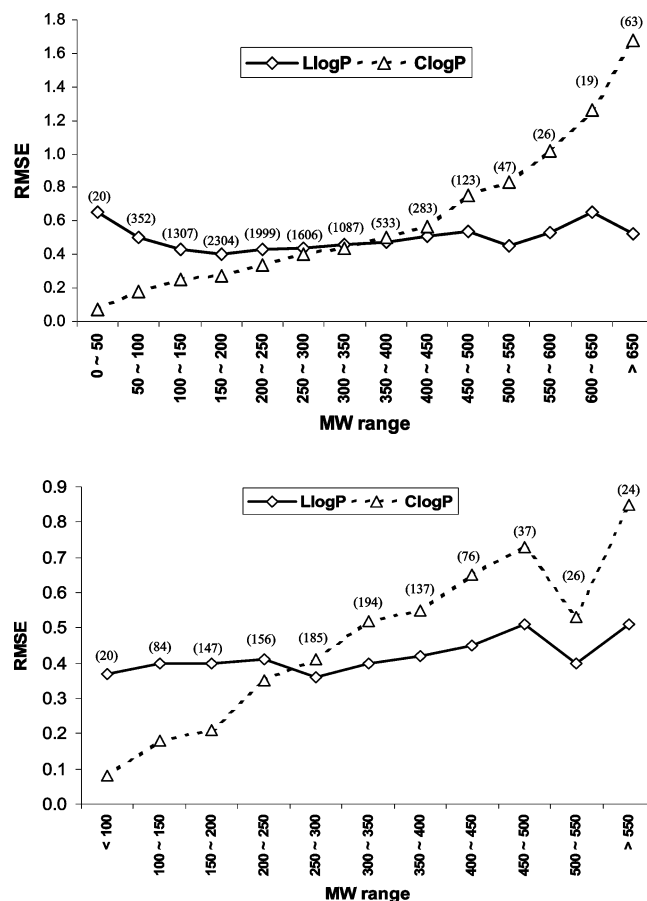
With the above-mentioned corrections for the special effects, the correlation coefficient between the calculated and the experimental logP values of the training set compounds is 0.97. The RMSE is 0.43 log units. For the test set compounds, the correlation coefficient and RMSE are 0.98 and 0.43, respectively. These statistical parameters are similar to those of ClogP. However our model is based on 39 latent descriptors and 8 correction factors. The ClogP model, on the other hand, is based on over 150 molecular fragments and over 600 other molecular descriptors.

It is interesting to compare molecular structures with high prediction errors by our method and ClogP. With corrections for the long-range effects, there are a total of 9 compounds for which logP predictions with our method are in error by more than 2 log units. Molecular structures of the 9 compounds are given in Figure 5.

Our highest prediction error is 3.81 log units (for compound **1** in Figure 5). In the StarList of nearly 11 000 compounds, this is the only compound with the I–O moiety. The first step in our PLS model building is to remove descriptors with zero variance. It is followed by descriptor and observation variance coverage analyses. As the I–O moiety appears only once in the StarList, the variance of this fingerprint descriptor for the training set (nearly 10K compounds) is nearly zero. As a result, this structural feature does not contribute to the 39 latent variables we used to build the PLS model. The calculated logP for this compound is simply a sum of contributions from five aromatic cH and an O–c(c) fragments (the lower case “c” represents an aromatic carbon atom). The closest analog is benzene with a predicted logP of 2.29 and an experimental value of 2.13. They are very close to the predicted value, 2.20, for compound **1**. The ECFC<sub>2</sub> features representing I–O and I–O–c did not contribute to the prediction by our model, but unfortunately they should be responsible for the low experimental value, –1.61, of this compound. The ClogP value for this compound is –1.615, surprisingly close to the experimental value. However, it is very likely that ClogP defined a parameter for this compound only to produce this surprisingly good agreement with experiment. In our PLS approach, we do not have the freedom of defining a special descriptor from a single experimental value. This is exactly the same for compounds **3** and **4** in Figure 5; for them, the central structural moieties appear in a single molecule only in the StarList. For compound **4**, the ClogP value is exactly the same as the experimental value, nearly too good to be true if ClogP did not define a special parameter from the experimental logP of this compound only.

Table 2 shows that there are a total of 25 StarList compounds for which the ClogP values are in error by over 2 log units. Molecular structures of 15 compounds with the highest ClogP errors are shown in Figure 6. An interesting observation is that most molecules in Figure 5 are small, some with unique or rare structural features. This contrasts to compounds in Figure 6, as most of them are relatively large molecules with common but repeating structural units. The highest ClogP error is found for compound **1** in Figure 6, for which the ClogP value, 14.36, is over 11 log units higher than the experimental value of 2.92. The value predicted by our model for this compound, 3.81, is much closer to the experimental value. It is likely that the high ClogP error for large molecules with repeating structural units is due to the fact that most ClogP fragment parameters were assigned from the experimental values of small molecules. When the fragment values were used to predict the logP for large molecules with many repeating structural units, the errors in fragment values were amplified. That is, our logP model tends to fail for compounds with rare or unprecedented structural features, while the ClogP failure is more obvious for large molecules with repeating precedent structural units.

The deteriorating performance of ClogP with increasing molecular size was noted before. To further examine this issue, we calculated RMSE of prediction by our model and ClogP in different MW ranges. The results are shown in Figure 7. The number of compounds in each MW range is given in the parentheses in the figure. The figure indicates that for compounds with molecular weight less than 250,



**Figure 7.** Comparison of root-mean-square error between the predicted and the experimental logP of the training (upper) and the test set (lower) compounds in different MW ranges. The numbers in parentheses are the number of compounds in the specific MW ranges.

ClogP performs better. For compounds with MW in the 250–400 range, ClogP and LlogP perform similarly. For compounds with MW over 400, the LlogP model performs significantly better than ClogP. In fact, over the full MW range, the performance of LlogP is much more consistent and stable compared to ClogP. Note that the performance of ClogP deteriorates monotonically with increasing MW. The ClogP error of prediction increases from nearly zero for compounds with MW less than 50 to over 1.6 log units for compounds with MW 650 or higher. These observations contrast sharply to the LlogP error of prediction, which runs between 0.4 and 0.6 log units over the full MW spectrum. Considering the fact that the molecular size of hits, leads, and candidate drugs tend to be bigger in modern drug discovery paradigm, it is clear that the LlogP model is superior for drug design applications.

## CONCLUDING REMARKS

This study demonstrates that the scores of molecular fingerprint can be used as high performance descriptors for QSPR studies. However the fingerprint features are highly correlated, and it requires regression techniques that can handle descriptor correlation for robust model development. The LlogP model built on 39 latent descriptors and 8 correction factors outperforms ClogP, especially for large molecules. Considering the large number of carefully designed molecular fragments and hundreds of other molecular



descriptors upon which ClogP was built, the ClogP model is much more complex. To explain the complexity of the ClogP model, Hansch and Leo stated that 'given our present state of knowledge there appears no way to reduce the complexity in logP calculations'.<sup>26</sup> The present study shows that now much simpler, more robust, and better performing models can be developed for logP prediction.

#### ACKNOWLEDGMENT

We are grateful to many colleagues for their suggestions/comments in the course of this study and their critical reading of the manuscript, especially those from Drs. James Damewood and Cristobal Alhambra. We are also grateful to members of the SciTegic Technical Support team for their timely assistance in resolving many issues encountered in the course of this study.

#### REFERENCES AND NOTES

- (1) Hansch, C.; Bjorkroth, J. P.; Leo, A. Hydrophobicity and central nervous system agents: on the principle of minimal hydrophobicity in drug design. *J. Pharm. Sci.* **1987**, *76*, 663–687.
- (2) Leo, A. *Environmental Health Chemistry*; McKinney, J. D., Ed.; Ann Arbor Science: Ann Arbor, MI, 1981.
- (3) Jain, N.; Yalkowsky, S. H. Estimation of the Aqueous Solubility 1: Application to Organic Nonelectrolytes. *J. Pharm. Sci.* **2001**, *90* (2), 234–252.
- (4) *OECD Guidelines for the Testing of Chemicals, Test No. 107*; Organisation for Economic Co-operation and Development (OECD): Paris, 1995. *Product Properties Test Guidelines, OPPTS 830.7570*; U.S. Environmental Protection Agency, U.S. Government Printing Office: Washington, DC, 1996. *OECD Guidelines for the Testing of Chemicals, Test No. 117*; Organisation for Economic Co-operation and Development (OECD): Paris, 2004. Valko, K.; Bevan, C.; Reynolds, D. *Anal. Chem.* **1997**, *69*, 2022–2029.
- (5) Fujita, T.; Iwasa, J.; Hansch, C. A New Substituent Constant,  $\pi$ , Derived from Partition Coefficients. *J. Am. Chem. Soc.* **1964**, *86*, 5175–5180.
- (6) Mannhold, R.; Dross, K. Calculation Procedures for Molecular Lipophilicity: A Comparative Study. *Quant. Struct.-Act. Relat.* **1996**, *15*, 403–409.
- (7) Leo, A.; Hoekman, D. Calculating logP(oct) with no missing fragments; The problem of estimating new interaction parameters. *Perspect. Drug Discovery Des.* **2000**, *18*, 19–38.
- (8) Daylight Chemical Information Systems Inc. ClogP Reference Manual. <http://www.daylight.com/dayhtml/doc/clogp/index.html> (accessed Dec 17, 2007).
- (9) Ghose, A. K.; Viswanadhan, V. N.; Wendoloski, J. Prediction of Hydrophobic (Lipophilic) Properties of Small Organic Molecules Using Fragmental Methods: An Analysis of ALOGP and CLOGP Methods. *J. Phys. Chem. A* **1998**, *102*, 3762–3772.
- (10) Klopman, G.; Li, J.-Y.; Wang, S. Dimayuga, J. Computer Automated log P Calculations Based on an Extended Group Contribution Approach. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 752–781.
- (11) Meylan, W.; Howard, P. H. Estimating log P with atom/fragments and water solubility with log P. *Perspect. Drug Discovery Des.* **2000**, *19*, 67–84.
- (12) Sun, H. A universal molecular descriptor system for prediction of logP, logS, logBB, and absorption. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 748–757.
- (13) Stanton, D.; Jurs, P. Development and use of charged partial surface area structural descriptors in computer-assisted quantitative structure-property relationship studies. *Anal. Chem.* **1990**, *62*, 2323–2329.
- (14) Ertl, P.; Rohde, B.; Selzer, P. Fast calculation of molecular polar surface area as a sum of fragment based contributions and its application to the prediction of drug transport properties. *J. Med. Chem.* **2000**, *43*, 3714–3717.
- (15) Hall, L. H.; Kier, L. B. The Molecular Connectivity Chi Indices and Kappa Shape Indices in Structure-Property Modeling. *Reviews of Computational Chemistry*; Boyd, D., Lipkowitz, K., Eds.; VCH Publishers Inc.: 1991; pp 367–422.
- (16) Hall, L. H.; Kier, L. B. The Nature of Structure-Activity Relationships and Their Relation to Molecular Connectivity. *Eur. J. Med. Chem.* **1997**, *12*, 307–312.
- (17) Kier, L. B.; Hall, L. H. *Molecular Structure Description: The Electrotological State*; Academic Press: 1999.
- (18) Kogej, T.; Engkvist, O.; Blomberg, N.; Muresan, S. Multifingerprint Based Similarity Searches for Targeted Class Compound Selection. *J. Chem. Inf. Model.* **2006**, *46*, 1201–1213.
- (19) Daylight Chemical Information Systems Inc. Fingerprint – Screening and Similarity. <http://www.daylight.com/dayhtml/doc/theory/theory.finger.html> (accessed Dec 17, 2007).
- (20) Pipeline Pilot Basic Chemistry Component Collection, SciTegic Inc., 9655 Chesapeake Drive, Suite 401, San Diego, CA 92123.
- (21) 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.
- (22) Wehrens, R.; Mevik, B. PLS: Partial Least Squares Regression (PLSR) and Principal Component Regression (PCR). R package version 2.0-1. <http://mevik.net/work/software/pls.html> (accessed Dec 17, 2007).
- (23) R Development Core Team (2007). *R: A language and environment for statistical computing*; R Foundation for Statistical Computing: Vienna, Austria, ISBN 3-900051-07-0. URL <http://www.R-project.org>. (accessed Dec 17, 2007).
- (24) Clark, M.; Cramer, R. D., III The Probability of Chance Correlation Using Partial Least Squares (PLS). *Quant. Struct.-Act. Relat.* **1993**, *12*, 137–145.
- (25) Cramer, R. D., III; Patterson, D. E.; Bunce, J. D. Comparative Molecular Field Analysis (CoMFA). 1. Effect of Shape on Binding of Steroids to Carrier Proteins. *J. Am. Chem. Soc.* **1988**, *110*, 5959–5967.
- (26) Hansch, C.; Leo, A.; Hoekman, D. *Exploring QSAR: Hydrophobic, Electronic, and Steric Constants*; American Chemical Society: Washington, DC, 1995.

CI700372S