

In Silico Prediction of Aqueous Solubility: The Solubility Challenge

M. Hewitt,* M. T. D. Cronin, S. J. Enoch, J. C. Madden, D. W. Roberts, and J. C. Dearden

School of Pharmacy and Biomolecular Sciences, Liverpool John Moores University, Byrom Street, Liverpool, L3 3AF, England

Received August 3, 2009

The dissolution of a chemical into water is a process fundamental to both chemistry and biology. The persistence of a chemical within the environment and the effects of a chemical within the body are dependent primarily upon aqueous solubility. With the well-documented limitations hindering the accurate experimental determination of aqueous solubility, the utilization of predictive methods have been widely investigated and employed. The setting of a solubility challenge by this journal proved an excellent opportunity to explore several different modeling methods, utilizing a supplied dataset of high-quality aqueous solubility measurements. Four contrasting approaches (simple linear regression, artificial neural networks, category formation, and available *in silico* models) were utilized within our laboratory and the quality of these predictions was assessed. These were chosen to span the multitude of modeling methods now in use, while also allowing for the evaluation of existing commercial solubility models. The conclusions of this study were surprising, in that a simple linear regression approach proved to be superior over more complex modeling methods. Possible explanations for this observation are discussed and also recommendations are made for future solubility prediction.

INTRODUCTION

The dissolution of a chemical into water is a process that is fundamental to both chemistry and biology. The persistence of a chemical (agrochemicals or pollutants) within the environment and the effects of a chemical (xenobiotic) within the body are dependent primarily upon aqueous solubility. Because the human body is ~60% water,¹ the latter point is of particular importance in both drug design and toxicology. For a xenobiotic to exert a physiological effect, it must first be transported to the site of action, almost exclusively in aqueous solution. Therefore, aqueous solubility determines this uptake, and also the distribution, metabolism, and elimination (ADME) characteristics of a chemical. Thus, in terms of pharmaceutical drug design, aqueous solubility is considered to be the most fundamental and important of physicochemical properties.^{2–4} Drugs that are not sufficiently water-soluble are open to an entire host of potential problems, including poor uptake,^{5,6} lack of pharmacological effect, and manufacturing and storage problems.⁷ As a result, it is estimated that, historically, up to 40% of drug discovery programs were abandoned because of insufficient solubility and associated poor pharmacokinetics under physiological conditions.^{8,9} However, current modeling methods now allow for the early screening of compounds and this figure has been reduced considerably.¹⁰

Given its role in ADME properties, it is clear why aqueous solubility is so important. Unfortunately, aqueous solubility is inherently difficult to measure experimentally, because of the difficulty in obtaining reproducible results. This is especially difficult in the case of compounds with very low solubilities. In addition, experimental measurements can be

very time-consuming.^{11,12} Other factors such as chemical stability and hydrolysis are also limiting factors in these experimental techniques. However, despite the development of numerous methodologies,¹³ experimental solubility determinations remain prone to high experimental error and thus are often unreliable.

Given these difficulties, and the increased utilization of combinatorial chemistry and high throughput screening, another approach is the estimation of solubility from more easily obtainable properties using quantitative structure–activity relationships (QSARs). The first example of aqueous solubility prediction using computational methods in the literature appeared in 1924 when Fühner observed that the solubility of a homologous series decreased with the addition of successive methylene groups.¹⁴ Of course, this addition of methylene groups was a measure of increasing molecular size, which became a crucial property at the core of solubility prediction. A detailed review of the historical modeling of aqueous solubility is given by Dearden¹⁵ and is not covered in any great detail here.

Over the next 80 years, the prediction of aqueous solubility slowly took form, really coming into its own in the 1990s. The predictive models developed in this time became more complex, as did the descriptors contained within them. Molecular size continued to be dominant, often being expressed in the form of molar surface area, because aqueous solubility is controlled by interactions between water and the surface of a molecule.¹⁶ Additional work also considered the negative correlation of solubility with the octanol:water partition coefficient (*P*), with this latter item becoming a benchmark relationship.^{17,18} Other properties that influence solubility also became apparent, including melting point,¹⁹ hydrogen bonding,²⁰ various atom/group contributions,^{21,22} and molecular connectivities.^{23–25} Until the 1990s, the

* Author to whom correspondence should be addressed. E-mail: M.Hewitt@ljmu.ac.uk.

majority of these approaches utilized linear regression models. However, from the early 1990s, more complex alternatives were being explored, including artificial neural networks (ANNs)^{26,27} and genetic algorithms.²⁸ This has resulted in a plethora of modeling approaches being available today, each of which is capable of utilizing a large number of molecular descriptors, which can be calculated using specific software.

Of the factors that were determined to be significant, molecular size and shape and hydrogen-bonding ability (polarity) are considered among those which are fundamental with regard to predicting aqueous solubility. For a compound to enter solution, a cavity must be formed in the solvent (water). Therefore, it makes sense that the larger the solute is, the larger the cavity required (hence, the relationship with molecular size). Similarly, the larger the cavity required, the more water–water hydrogen bonds which require breaking to create the cavity. More details can be found on the dissolution process and solute–solvent interaction in the literature.^{29,30}

There are now many models available within the literature for predicting aqueous solubility,^{15,29} but these are accompanied by errors of ~ 1 order of magnitude. A primary source of this error is the error associated with the data from which the models are derived. It has been estimated that experimental aqueous solubility data contain an experimental error of up to 1.5 log units.¹⁵ Therefore, this makes developing predictive models inherently difficult, because they are limited by the quality of the experimental data. In addition, the solubility process is not controlled by a single property, but rather is a complex process that is dependent on various molecular attributes in addition to interactions between solute molecules themselves.²⁹

Given the obvious importance of aqueous solubility in pharmaceutical sciences and the plethora of modeling approaches and molecular descriptors available today, a challenge was set recently by Llinàs et al., termed the “Solubility Challenge”.³¹ The challenge was to predict the solubility of 32 druglike compounds, using a high-quality training set of 100 compounds with experimentally determined solubilities. The challenge was laid down with the objective of assessing the current state of the art in solubility prediction and going on to make recommendations as to which of the available methods are best to use when making a prediction. Competitors were invited to submit their predictions to the editors of the *Journal of Chemical Information and Modeling* and a follow-up publication discussed the results.³²

Our laboratory participated in the Solubility Challenge and following the publication of the results, the present study provides an assessment of differing modeling approaches used in the prediction of solubility. Using four contrasting approaches used in our own Solubility Challenge entries, the Solubility Challenge dataset was used to develop predictive models. The results of these were then submitted and, upon publication of the results, a full analysis of predictive performance was made.

METHODS

Datasets. Training Set. A high-quality training set composed of 100 diverse druglike compounds was obtained from Llinàs et al.³¹ These compounds were chosen by the authors

because they represented a wide range of available pharmaceutical chemicals in use today. Each of these compounds had to meet three criteria: (a) contain an ionizable group with a pK_a value of 1–13, (b) be commercially available and relatively affordable, and (c) not pose any significant hazards. The experimentally determined intrinsic solubility of each compound was measured by Llinàs et al.³¹

Four compounds in the training set (chlorprothixene, phthalic acid, sulindac and trichlormethiazide) were identified as existing in two polymorphic forms, effectively increasing the number of compounds to 104 (we called these simply “form I” and “form II”). For the purposes of this analysis, only one of these forms was considered in each case. In addition, five of the training-set compounds (5-fluorouracil, levofloxacin, L-proline, orbifloxacin, and procainamide) were too soluble to be measured and were therefore removed from the analysis. Finally, two more compounds were removed, because they were shown to decompose during analysis, yielding a final training set of 97 compounds. This training set is shown in Table 1, together with the experimental intrinsic solubility and the category to which each compound is later assigned (explained in detail in a later section).

Test Set. Similarly, a test set of 32 druglike molecules was also obtained from Llinàs et al.³¹ These compounds were chosen to reflect the diversity of available drugs and were again required to meet the previously detailed criteria. Intrinsic solubilities were measured experimentally by Llinàs et al.,³¹ but these were released only following completion of the Solubility Challenge. The challenge was, of course, to predict the intrinsic solubilities of these 32 compounds listed in Table 2.

Solubility Measurement. Intrinsic solubility was determined experimentally by Llinàs et al.³¹ and these data were obtained from the Solubility Challenge publication.³¹ The authors utilized a potentiometric method, providing very accurate and rapid solubility determinations. This entailed a “Chasing Equilibrium” technique (CheqSol),³³ producing a final precipitate that is thermodynamically driven, yielding highly reproducible solubility data with an associated error of ~ 0.05 log units. Full experimental details are given in Llinàs et al.³¹ In all cases, intrinsic solubilities were provided in units of $\mu\text{g/mL}$. However, in this study, solubilities were converted into molar units (Mol/L) allowing solubility to be quantified independently of molecular weight.

Structure Generation and Descriptor Calculation. Each dataset (training and test) obtained from the Solubility Challenge publication³¹ contained graphical representations of chemical structure together with the chemical name for each compound. SMILES notations for each structure were obtained from the chemical name, via the use of various online chemical databases, including ChemIDPlus³⁴ and ChemBioFinder.³⁵ As a quality control measure, each SMILES notation was checked against the corresponding structure given in the original article. For each dataset (training and test), the SMILES notations were then converted into an SDF file containing energy optimized three-dimensional structures, using MolConvert, version 5.2.1.0, (a module within the MarvinSketch software developed by ChemAxon).³⁶

Using chemical structures in SDF format, many physico-chemical descriptors were calculated. The EPISUITE software package (version 3.12),³⁷ developed by the Syracuse

Table 1. Training Set Compounds with Log(Intrinsic Solubility) Data and Category Information^a

name	log(intrinsic solubility) (M)	CAS No.	category ^b
1-naphthol ^{Val}	-1.14	90-15-3	4
2-amino-5-bromobenzoic acid	-2.40	5794-88-7	2
4-iodophenol	-1.05	540-38-5	4
5-bromo-2,4-dihydroxybenzoic acid	-1.99	7355-22-8	2
5-hydroxybenzoic acid	-0.60	99-96-7	2
acetaminophen	-0.24	103-90-2	4
acetazolamide	-1.79	59-66-5	3
alprenolol ^{Val}	-2.03	13655-52-2	4
amantadine	-1.03	768-94-5	6
amiodarone	-7.98	1951-25-3	6
amitriptyline	-3.99	50-48-6	6
amodiaquine	-5.34	86-42-0	1
atropine	-1.46	51-55-8	1
azathioprine	-2.65	446-86-6	outside of the test set categories
benzylimidazole	-1.46	4238-71-5	5
bromogranine ^{Val}	-3.46	830-93-3	6
bupivacaine	-2.68	38396-39-3	6
carprofen	-4.14	52263-47-5	2
carvedilol	-3.87	72956-09-3	1
cephalothin	-2.54	153-61-7	2
chlorpheniramine	-2.11	132-22-9	6
chlorpromazine	-4.57	50-53-3	6
chlorpropamide ^{Val}	-2.69	94-20-2	3
chlorprothixene			
form I	-6.25	113-59-7	6
form II	-5.37	113-59-7	polymorph - excluded
cimetidine	-1.09	51481-61-9	outside of the test set categories
ciprofloxacin	-3.12	85721-33-1	2
danofloxacin	-2.45	112398-08-0	2
deprenyl	-1.78	2323-36-6	outside of the test set categories
desipramine	-3.06	50-47-5	6
diazoxide	-2.72	364-98-7	3
diclofenac ^{Val}	-4.93	15307-86-5	2
difloxacin	-3.20	98106-17-3	2
diltiazem	-2.78	42399-41-7	6
diphenhydramine	-2.36	58-73-1	6
diphenylhydantoin	-3.26	57-41-0	outside of the test set categories
enrofloxacin	-2.74	93106-60-6	2
famotidine	-2.18	76824-35-6	outside of the test set categories
fenoprofen	-3.08	31879-05-7	2
flufenamic acid	-4.80	530-78-9	2
flumequine	-3.15	42835-25-6	2
flurbiprofen	-3.54	5104-49-4	2
glipizide	-5.14	29094-61-9	3
guanine	-3.61	73-40-5	6
hexobarbital	-2.04	56-29-1	outside of the test set categories
hydroflumethiazide	-2.49	135-09-1	3
ibuprofen	-2.91	15687-27-1	2
lomefloxacin ^{Val}	-1.88	98079-51-7	2
loperamide	-6.75	53179-11-6	4
maprotiline	-4.13	10262-69-8	6
meclizine ^{Val}	-6.07	569-65-3	outside of the test set categories
mefenamic acid	-6.12	61-68-7	2
metoclopramide	-3.05	364-62-5	6
metronidazole	-0.45	443-48-1	outside of the test set categories
miconazole	-4.69	22916-47-8	5
nalidixic acid ^{Val}	-2.98	389-08-2	2
naloxone	-2.42	465-65-6	1
naproxen	-3.86	22204-53-1	2
niflumic acid ^{Val}	-4.04	4394-00-7	2
nitrofurantoin	-2.62	67-20-9	outside of the test set categories
norfloxacin	-2.26	70458-96-7	2
nortriptyline	-3.44	72-69-5	6
ofloxacin ^{Val}	-0.83	82419-36-1	2
oxytetracycline	-2.75	6153-64-6	outside of the test set categories
papaverine	-3.40	58-74-2	outside of the test set categories
phenanthroline ^{Val}	-0.88	12678-01-2	outside of the test set categories
phenazopyridine	-3.52	94-78-0	outside of the test set categories
phenobarbital	-1.66	50-06-6	outside of the test set categories
phenylbutazone ^{Val}	-3.88	50-33-9	outside of the test set categories
phthalic acid			
form I	-0.83	88-99-3	2

Table 1. Continued

name	log(intrinsic solubility) (M)	CAS No.	category ^b
form II	-0.71	88-99-3	polymorph - excluded
pindolol ^{Val}	-3.19	13523-86-9	1
piroxicam ^{Val}	-4.32	36322-90-4	3
procaine	-1.09	59-46-1	6
propranolol	-2.90	525-66-6	1
quinine	-2.30	130-95-0	outside of the test set categories
ranitidine	-2.00	66357-35-5	outside of the test set categories
sarafloxacin	-2.72	98105-99-8	2
sertraline	-4.32	79617-96-2	6
sparfloxacin	-2.96	110871-86-8	2
sulfacetamide	-0.85	144-80-9	3
sulfamethazine ^{Val}	-2.17	57-68-1	3
sulfasalazine	-5.74	599-79-1	outside of the test set categories
sulfathiazole	-2.10	72-14-0	3
sulindac			
form I	-4.06	38194-50-2	outside of the test set categories
form II	-3.23	38194-50-2	outside of the test set categories
tetracaine ^{Val}	-2.43	94-24-6	6
tetracycline ^{Val}	-2.57	60-54-8	outside of the test set categories
thymol	-1.37	89-83-8	outside of the test set categories
tolmetin	-3.50	26171-23-3	2
trichloromethiazide ^{Val}			
form I	-3.11	133-67-5	3
form II	-2.76	133-67-5	polymorph - excluded
trimethoprim	-2.41	738-70-5	outside of the test set categories
trimipramine	-4.27	739-71-9	6
tryptamine	-2.50	61-54-1	6
verapamil	-3.63	52-53-9	outside of the test set categories
warfarin	-4.27	81-81-2	4

^a Compounds indicated by the superscript "Val" are those used as a subsequent validation set. Details of the category definitions are given later in the Methods section. ^b Categories are defined as follows: (1) compounds with BOTH a hydrogen-bond-accepting N and hydrogen-bond-donating OH, (2) carboxylic acids, (3) sulfonamides (hydrogen-bond acceptors, NOT donors), (4) compounds with OH groups (hydrogen-bond acceptors and donors) but NO hydrogen-bond-accepting N (in some cases N is present but too weakly basic to hydrogen-bond significantly), (5) heterocyclic compounds with N sufficiently basic to be protonated at neutral pH, and (6) hydrogen-bond-accepting N with NO hydrogen-bond donors.

Research Corporation (SRC), was used to calculate the logarithm of the octanol:water partition coefficient ($\log P$), together with the melting and boiling points of each compound. Although the WSKOWWIN module of EPISUTE is able to make predictions of aqueous solubility, previous studies have shown the predictions of this software to be rather poor compared to other *in silico* models.¹⁵ As a result, these predictions were not considered in this study.

The HYBOT-PLUS module contained within the MOLPRO software package (Version 2.1.0.706)³⁸ was also used to generate several hydrogen bonding descriptors. In addition, an array of descriptors was also calculated using the Dragon Professional (Version 5.3) software developed by the Milano Chemometrics and QSAR Research Group.³⁹ Initially, a very large number of descriptors were calculated (>1000). However, following the removal of highly correlated descriptors (correlation of >0.90 (the lowest correlation level available within the descriptor reduction tool)), this was reduced to a final descriptor set of 426.

Model Development and Statistical Analysis. Given that one of the study objectives was to perform an unbiased evaluation of the performance of different modeling approaches to predict aqueous solubility, a range of differing modeling techniques and software tools were utilized. These ranged greatly in their complexity and approach and reflected the scope of approaches used to predict solubility to date. Four modeling approaches were investigated in this study: multiple linear regression (MLR), artificial neural networks

(ANN), category formation and local modeling, and *in silico* predictions using existing models. Each method is discussed in detail below. The full data matrix containing experimental solubility data together with individual model predictions is available as Supporting Information.

Multiple Linear Regression (MLR) Model. The simplest approach used in the study, a traditional regression model, seemed to be an obvious starting point. Used widely in the modeling of solubility, many linear regression models (QSARs) have been proposed for solubility prediction.⁴⁰ To aid in model evaluation and validation, the training set was further subdivided to create a small validation set (20% of compounds). To ensure that a representative validation set was created, the training set compounds were ordered, according to measured solubility, and every fifth compound was taken as a validation compound. Although there are many ways to split the data (e.g., random selection or based on physicochemical properties), this approach has been determined to work as well as others and is the simplest to perform. This led to the development of a 19-compound validation set, as indicated by the superscript "Val" in Table 1.

Descriptor selection was performed using a genetic algorithm contained within the Mobydigs software (Version 1.0).⁴¹ The full set of 426 descriptors was used as input variables, and the experimentally derived intrinsic solubilities of the training set were used as the response variable.

Table 2. Test Set Compounds with Assigned Category Information^a

	name	log intrinsic solubility (μ M)	CAS No.	category ^c
1	2-chloromandelic acid	TSTM ^b	10421-85-9	2
2	acebutolol	-2.71	37517-30-9	1
3	amoxicillin	-2.03	26787-78-0	2
4	bendroflumethiazide	-3.89	73-48-3	3
5	benzocaine	-2.74	94-09-7	6
6	benzthiazide	-4.46	91-33-8	3
7	clozapine	-3.24	5786-21-0	6
8	dibucaine	-4.39	85-79-0	6
9	diethylstilbestrol	-4.43	56-53-1	4
10	diflunisal	-5.94	22494-42-4	2
11	dipyridamole	-5.16	58-32-2	4
12	ephedrine	TSTM ^b	299-42-3	1
13	folic acid	-5.25	59-30-3	2
14	furosemide	-4.23	54-31-9	2
15	hydrochlorothiazide	-2.68	58-93-5	3
16	imipramine	-4.11	50-49-7	6
17	indomethacin	-2.94	53-86-1	2
18	ketoprofen	-3.21	22071-15-4	2
19	lidocaine	-1.87	137-58-6	6
20	marbofloxacin	TSTM ^b	115550-35-1	2
21	meclofenamic acid	-6.27	644-62-2	2
22	naphthoic acid	-3.77	1320-04-3	2
23	probenecid	-4.86	57-66-9	2
24	pseudoephedrine	TSTM ^b	90-82-4	1
25	pyrimethamine	-4.11	58-14-0	6
26	salicylic acid	-1.93	69-72-7	2
27	sulfamerazine	-3.12	127-79-7	3
28	sulfamethizole	-2.78	144-82-1	3
29	terfenadine	-7.74	50679-08-8	1
30	thiabendazole	-3.48	148-79-8	5
31	tolbutamide	-3.46	64-77-7	3
32	trazodone	-3.47	19794-93-5	6

^a Details of the category definitions are given later in the Methods section. ^b TSTM = too soluble to measure. ^c Categories are defined as follows: (1) compounds with BOTH a hydrogen-bond-accepting N and hydrogen-bond-donating OH, (2) carboxylic acids, (3) sulfonamides (hydrogen-bond acceptors NOT donors), (4) compounds with OH groups (hydrogen-bond acceptors and donors) but NO hydrogen-bond-accepting N (in some cases N is present but too weakly basic to hydrogen-bond significantly), (5) heterocyclic compounds with N sufficiently basic to be protonated at neutral pH, and (6) hydrogen-bond-accepting N with NO hydrogen-bond donors.

The default settings for the Mobydigs software were used for the genetic algorithm, with the exception that the maximum number of descriptors allowed in a model was limited to five. This reduces the likelihood of model overfitting and maintains a compound:descriptor ratio greater than 5:1, which aids in the development of a robust model.⁴² The genetic algorithm was allowed to run for 50 000 iterations before the best model (ranked top according to validation set predictivity) was recorded. To assess the optimal number of descriptors, the analysis was repeated setting the maximum number of descriptors allowed in the model at four, and then again at three. The results are given and discussed in detail in the Results section and summary statistics are supplied.

Artificial Neural Network (ANN). In light of the increased use of more complex modeling methods to predict aqueous solubility recently, it was thought prudent to investigate one of the most utilized of these methods, artificial neural networks (ANNs). As mentioned previously, ANNs have been widely utilized in the prediction of solubility, in addition to many other ADME properties. Inspired by the central nervous system, these models train a network of artificial neurons to reproduce the property being modeled

from a set of input variables (i.e., descriptors). The perceived advantage of ANNs over MLR methods is their capability of applying nonlinear pattern recognition functions aiding in the detection of more complex statistical relationships.⁴³

ANN analysis was performed using the Intelligent Problem Solver algorithms contained within the Statistica statistical software⁴⁴ using intrinsic solubility as the dependent variable and the same three descriptors selected by the genetic algorithm as independent variables. The Intelligent Problem Solver algorithm attempts to build the optimum neural network by training and validating several linear and three-layer perceptron neural networks. In this study, the algorithm was allowed to run for a default 500 cycles, with the training set of 97 compounds (with the same 19 compound validation set being utilized as a subset for the estimation of the training error).

Following completion of the 500 cycles, 50 neural networks were retained for further analysis. These networks were selected to be diverse in terms of the type and architecture of the networks tested during the Intelligent Problem Solver routine. It was from these 50 networks that the final network was selected. This selection was based on the balance between training and validation error in conjunction with perceived model complexity, with simpler architectures being preferred. The quality of the final neural network model was then assessed for fit using the coefficient of determination (r^2) and root-mean-square error (RMSE), as with the previous MLR model.

Chemical Category Formation and Local Modeling.

Given that the compounds considered in this study are all pharmaceuticals, they span incredibly diverse structural and mechanistic domains. Therefore, given the possible problems with attempting to model these in a single global model,⁴⁵ a third approach was undertaken, whereby chemicals in the training set were split into a small number of unique chemical categories, in the hope that each category, when taken individually, would result in the development of a more meaningful and predictive local model.

Training and test set compounds were assigned to one of six chemical categories, as listed below and shown in Tables 1 and 2, respectively. Developed in-house, these categories were centered primarily upon hydrogen-bonding ability, a well-known determinant in aqueous solubility.²⁰ When present, highly electronegative atoms attract a hydrogen atom's sole electron toward themselves, leaving the strongly charged hydrogen nucleus exposed. In this state, the exposed positive nucleus can exert considerable attraction on the electrons in other molecules, as well as other atoms in the same molecule, forming a protonic bridge that is substantially stronger than most other types of dipole interactions.⁴⁶ Understandably, the strength of these hydrogen bonds can play a significant role in aqueous solubility. The following six categories were used in this study:

- (1) Compounds with BOTH a hydrogen-bond-accepting N and hydrogen-bond-donating OH,
- (2) Carboxylic acids,
- (3) Sulfonamides (hydrogen-bond acceptors NOT donors),
- (4) Compounds with OH groups (hydrogen-bond acceptors and donors) but NO hydrogen-bond-accepting N (in some cases, N is present but too weakly basic to hydrogen-bond significantly),

(5) Heterocyclic compounds with N sufficiently basic to be protonated at neutral pH, and

(6) Hydrogen-bond accepting N with NO hydrogen-bond donors

Following category formation, each subset of compounds was then considered individually and, based on the size of the subset, QSAR or read-across methods were used to generate solubility predictions. The models derived within each category and their predictive performance are detailed in the Results section.

In Silico Prediction. In this final approach, many existing commercially available models were used to predict aqueous solubility. This approach both highlighted the availability of solubility models and, at the same time, evaluated their performance. In this study, existing models from four commercial sources were used to predict solubility: the ChemSilico CSLogWS Model, the Optibrium StarDrop Model, the Pharma Algorithms Solubility Predictor, and SPARC. These are listed below with a short summary of each model.

ChemSilico CSLogWS Model. Predictions of intrinsic solubility were kindly generated and provided by ChemSilico using their CSLogWS model (http://www.chemsilico.com/CS_prWS/WShome.html). This neural network model was developed using a training set of over 5964 diverse compounds and utilized 519 topological and E-state descriptors.

Optibrium (a Spin-out Company of BioFocus DPI) StarDrop Model. Intrinsic aqueous solubility was predicted using a model that is part of the StarDrop interactive software platform and was developed by Optibrium (<http://www.optibrium.com/stardrop.php>). No information on model type or training chemicals could be found. Again, these predictions were kindly provided by Optibrium.

Pharma Algorithms Solubility Predictor. Predictions of aqueous solubility were made online using the freely accessible ADME/TOX WEB application, which is available from the Pharma Algorithms website (<http://www.pharma-algorithms.com>). Using the ADME Boxes v4.0 platform, this web application was used to make predictions of water solubility (log *S*), using the Solubility Predictor. Using a training set of over 6800 compounds, the log solubility in pure water at 25 °C was predicted.

SPARC. Again, post-Challenge, intrinsic aqueous solubility predictions were made using the freely accessible model contained within the web-based tool SPARC (version 4.2) (<http://ibmlc2.chem.uga.edu/sparc/>). The solubility model within SPARC does not calculate solubility from first principles, but utilizes an activity coefficient model. Unlike other models, SPARC bases its predictions on solute–solvent interactions, combining this knowledge with the Flory–Huggins excess entropy of mixing contribution for placing a solute molecule in solution. More information on the methodology behind the SPARC model can be found in Hilal and Karickhoff.⁴⁷

Simulations Plus Solubility Models. Intrinsic aqueous solubility predictions were generated by Simulations Plus model developers using eight solubility models, developed specifically for use in the Solubility Challenge. It must be stressed that, currently, none of these models are available commercially, but they do employ many complex modeling approaches often utilized in commercial models and, therefore, are of interest. Furthermore, these Simulations Plus

models are derived using the same training and test data supplied within the challenge so these models can be directly compared with those developed in this study. The details of the eight models are given below:

SHARK: A consensus (arithmetic mean) prediction from an ensemble of 32 ANN models, each built using only two-dimensional (2-D) descriptors in their development. This model takes the 94 compound dataset from the Solubility Challenge and splits it into a 80%–90% training set and 10%–20% test set for validating model predictivity. The splitting of data into training and test data is performed on the basis of compound clustering with the aid of self-organizing Kohonen maps.

YINAN: A consensus prediction was taken from 25 2-D ANN ensembles. The same training and test procedure as that in the SHARK model was used.

UIQBB: A consensus prediction from an ensemble of 32 ANN models, each built using only three-dimensional (3-D) descriptors. The same training and test procedure as that in the SHARK model was used.

LGGAV: A consensus prediction was taken from 23 3-D ANN ensembles. The same training and test procedure as that in the SHARK model was used.

A69EM: A consensus prediction from an ensemble of 32 ANN models, each built using only 2-D descriptors, as in the SHARK model. However, this model was trained using the full 94 compound dataset taken from the Solubility Challenge. Thirty-eight (38) additional compounds, with intrinsic solubility data, were collected from the literature and used as an external test set.

NSLIC: A consensus prediction from an ensemble of 32 ANN models, each built using only 3-D descriptors, as in the SHARK model. However, this model was trained using the full 94 compound dataset taken from the Solubility Challenge. 38 additional compounds, with intrinsic solubility data, were collected from the literature and used as an external test set.

AM108: A consensus prediction from an ensemble of 32 ANN models, each built using only 2-D descriptors, as in the SHARK model. However, this model used a larger 132 compound data set (94 Solubility Challenge chemicals plus 38 compounds obtained from the literature), which was split into training and test datasets as done previously.

OLASM: The final Simulations Plus model yielded a consensus prediction from 15 2-D ANN ensembles, using the same datasets as those used in the AM108 model.

Predictions from all models were kindly supplied by Simulations Plus.

Consensus Prediction. To ascertain any benefits offered by consensus modeling, a consensus prediction was also generated using predictions from four of the models detailed above (the ChemSilico, Optibrium, Pharma Algorithms, and Simulations Plus (SHARK) models). At the time of modeling, predictions were available only for these four models. Subsequently, Simulations Plus released predictions using seven additional models, which are discussed later for the sake of completeness. A simple mean value was taken from the four predictions in logarithmic form for each compound, and this was taken to be the consensus prediction.

Statistical Analysis. To assess the performance of each model and to allow for model comparison, it was necessary to calculate many summary statistics. Where possible, the

correlation coefficients (R^2), standard errors (s), and the Fisher statistic (F) were calculated using Minitab for Windows.⁴⁸

In addition, the root-mean-square error (RMSE) was also calculated also providing a measure of model performance. RMSE values were calculated using eq 1, where "Pred" refers to the predicted solubility values and "Obs" represents the experimentally observed values:

$$\text{RMSE} = \sqrt{\frac{\sum (\text{Pred} - \text{Obs})^2}{N - 1}} \quad (1)$$

Applicability Domain Assessment. In an attempt to assess the relationship between applicability domain coverage and predictive performance, the freely available Toxmatch software tool was used.⁴⁹ Toxmatch is able to quantify the similarity of two chemicals based on chemical structure. Although the software allows for both chemical structure and structural similarity to be coded for and quantified differently using different methods, this study utilized a structural fingerprint approach, whereby the similarity was quantified using the Tanimoto distance. Structural similarity is then quantified as a numerical value scaled between zero and one, with zero being no similarity and one being very high similarity/identical.

Using such an approach, a similarity cutoff is required to quantify when compounds share enough structural components to be classed as similar. A previous study by the authors which utilized the Toxmatch software found a cutoff of 0.6 to be sufficient to form valid structural categories.⁵⁰ Given that the current dataset has similar structural diversity, it seemed prudent to use this same cutoff. Using a lower cutoff (e.g., 0.4) would flag compounds as being similar when they share only subtle structural features and may lead to a fuzzy and ill-defined applicability domain assessment. Further information of applicability domain definition can be found in Shacham et al.⁵¹

RESULTS

Utilizing the previously mentioned four distinctly different modeling approaches, predictions of intrinsic solubility for the 32 test set compounds were made. The results of each approach, including models, summary statistics and predicted solubility values are discussed in detail below.

Multiple Linear Regression (MLR) Model. Following model building (as detailed in the Methods section), an optimal three-descriptor regression model was developed. The inclusion of additional descriptors was accompanied by only a small increase in model performance, not enough to justify the added model complexity. The optimal regression model, together with summary statistics, is given below (see eq 2). When possible in this study, each model was evaluated using three correlation coefficients, these being calculated for the training set compounds ($R^2_{(\text{Train})}$), the validation set compounds ($R^2_{(\text{Val})}$), and the test set compounds ($R^2_{(\text{Test})}$). Although $R^2_{(\text{Train})}$ describes the model's fit with the training data, the real test for the model comes when making predictions for external test set compounds. This measure of predictivity is given by the correlation coefficients for the validation ($R^2_{(\text{Val})}$) and test ($R^2_{(\text{Test})}$) sets, with the prediction of the test set solubilities being the ultimate goal in this study.

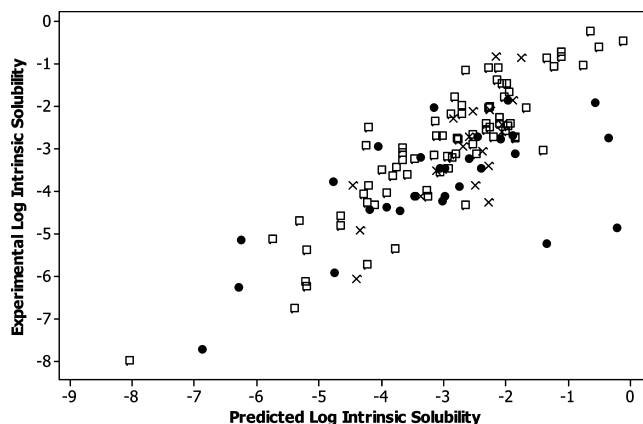


Figure 1. Plot of experimental versus predicted log(intrinsic solubility) using the MLR model (eq 2) for (□) the training set, (×) the validation set, and (●) the test set.

The standard error (s) and the Fisher statistic (F) are also given.

$$\begin{aligned} \log(\text{intrinsic solubility}) &= 7.706 - 0.613 \log \\ &\quad P - 0.007\text{BPt} - 5.404\text{R2e} + \\ R^2_{(\text{Train})} &= 0.74 \quad R^2_{(\text{Val})} = 0.67 \quad R^2_{(\text{Test})} = 0.51 \quad (2) \\ \text{RMSE}_{(\text{Test})} &= 0.95 \\ s &= 0.77 \quad F = 216.53 \end{aligned}$$

where $\log P$ is the logarithm of the octanol:water partition coefficient (calculated using the KOWWIN module of the United States Environmental Protection Agency's (USEPA's) EPISUITE software³⁷), BPt is the boiling point (given in Celsius) (calculated using the MPBPWIN module of the USEPA's EPISUITE software³⁷), R2e+ is the R maximal autocorrelation of lag 2/weighted by atomic Sanderson electronegativities (calculated using DRAGON Professional (Version 5.3) software³⁹).

In agreement with the historical modeling of aqueous solubility, the inverse relationship with hydrophobicity (as described by $\log P$) is again shown to be highly correlated with solubility. Similarly, a correlation with the boiling point of a compound (BPt) was also observed. Because an increasing boiling point (and melting point) is often associated with an increase in intermolecular interactions (including hydrogen bonding), the importance of this descriptor is not surprising. However, despite the molecular weight of the training set compounds varying from 138.12 (5-hydroxybenzoic acid) to 645.32 (amiodarone), no strong correlation with molecular size is present, with the differences in boiling point being due to other factors such as hydrogen bonding or polarity.⁵² The final descriptor (R2e+), a descriptor belonging to a set known as the GETAWAY descriptors,⁵³ was also used in the model. Because GETAWAY is an acronym for Geometry, Topology, and Atom-Weights Assembly, this descriptor clearly is coding for a molecular property related to molecular size and connectivity. Although not easy to interpret from its title, this descriptor has also been used in other modeling efforts to model the melting point.⁵¹

As can be seen from the model statistics and also from Figure 1, the statistical fit of the training and validation data was reasonably good (R^2 values of 0.74 and 0.67, respectively). However, there is a reduction in performance when

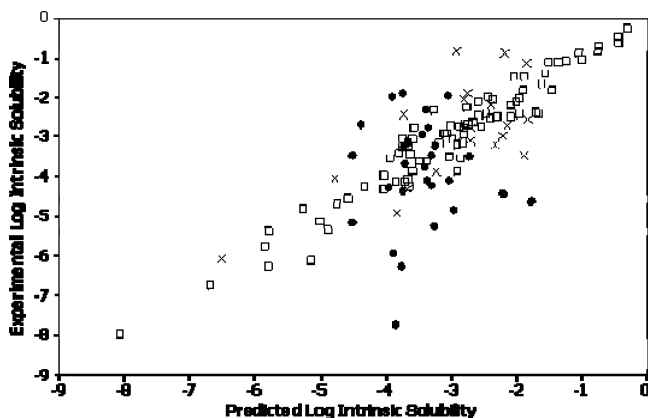


Figure 2. Plot of experimental versus predicted log(intrinsic solubility) using the ANN model for (□) the training set, (×) the validation set, and (●) the test set.

looking at the predictions made for the test set, with an R^2 value of 0.51. This can also be seen in the root mean square error (RMSE) value for the test set of 0.95.

Artificial Neural Network (ANN). The Intelligent Problem Solver routine within the Statistica software⁴⁴ generated 10 neural network models of differing architecture. As stated previously, the same descriptors as those selected by the genetic algorithm and used in the MLR model were used as input variables. Of the models generated, the model with the best predictive performance was found to be a two-descriptor multilayer perceptron neural network. Interestingly, in this case, the inclusion of the third descriptor (boiling point) did not offer any increase in predictive performance. The following descriptors were used as inputs by the neural network:

(1) $\log P$ (calculated using the KOWWIN module of the USEPA's EPISUITE software³⁷).

(2) $R2e+ = R$ maximal autocorrelation of lag 2/weighted by atomic Sanderson electronegativities (calculated using DRAGON Professional (Version 5.3)³⁹).

As with the MLR model, both hydrophobicity and molecular size/connectivity are again used as descriptors, highlighting their importance. Of course, the arrangement of atoms with a molecular structure may be related to many properties, including polarization, molecular shape and size, etc. The significance of these descriptors seems obvious, when one considers the processes involved in solubility, where a compound's size, shape, polarity, etc. are crucial.

The summary statistics of the neural network model are given below, together with a plot of experimental versus predicted solubility (see Figure 2).

$$R^2_{(\text{Train})} = 0.79 \quad R^2_{(\text{Val})} = 0.56 \quad R^2_{(\text{Test})} = 0.40 \quad \text{RMSE}_{(\text{Test})} = 1.51$$

As evident from Figure 2 and the accompanying statistics, the statistical fit of the training set was high and comparable to that observed in the MLR model ($R^2_{(\text{Train})} = 0.79$). Unfortunately, although a reduction in performance was expected when utilizing external data, the predictivity of the validation compounds was significantly reduced ($R^2_{(\text{Val})} = 0.56$). Following completion of the study and the release of the test set solubilities, it is now apparent that the neural network model is likely to have been overtrained (overfitted to the training data). This is especially clear when the test

set predictions are considered, which yielded a low R^2 value of 0.40 and a large RMSE value of 1.51, showing that the model has very poor predictive capacity outside its training data. It is interesting as the training set compounds are predicted within a tight band, much narrower than would be expected, as seen in Figure 2. This could indicate that the test set compounds are outside the domain of the training set, meaning that they may be structurally or mechanistically unrepresented within the training set compounds. However, a simple examination of the datasets reveals that this is not true for the majority of the compounds and therefore points heavily to overtraining of the network.

Chemical Category Formation and Local Modeling.

Following the splitting of all training and test compounds into one of the six derived categories, each subset of compounds was then modeled individually. The modeling approach utilized was dependent on the number of training compounds present within each category. Where possible, QSAR models were developed. Where data were sparse, read-across predictions were made. Modeling efforts for each of the six categories are summarized in Table 3 and discussed individually below. In addition, plots of observed solubility versus predicted solubility are shown in Figure 3.

Category 1: Compounds with BOTH a Hydrogen-Bond-Accepting N and Hydrogen-Bond-Donating OH. As the name of the category suggests, this category contains compounds that possess both a N atom capable of acting as a hydrogen-bond acceptor and a hydroxyl (OH) substituent that acts as a hydrogen-bond donor. Initial modeling revealed that a single parameter regression model could be generated with excellent statistical fit. In contrast to the previous regression modeling, no validation set was utilized, given the limited number of compounds in the training datasets.

The single descriptor used by this model was heat of formation (Δ_{HF}), calculated using TSAR for Windows.⁵⁴ Δ_{HF} represents the enthalpy change during the formation of 1 mol of a substance⁵⁵ (i.e., covalent bond formation). Although this is a measure of relative thermal stability, it is difficult to relate this to solubility. It is evident from the negative correlation coefficient (see Table 3) that increasing Δ_{HF} leads to reduced solubility. Given the small number of compounds present within this category, it is also possible that this is a chance correlation. However, for the purposes of this study, and as no other descriptor offered the same level of predictive performance, the Δ_{HF} value was used in the final QSAR model.

Unfortunately, two of the four test compounds (ephedrine and pseudoephedrine) were reported as being too soluble to measure and, therefore, assessment of the test set predictions made by this model could not be made using traditional correlation coefficients because of the lack of data points. However, upon examination of the predictions and from visual inspection of Figure 3, the predictive performance of this model seems poor. Acebutolol has an experimental log(intrinsic solubility) of -2.71 but was predicted to be -1.15 , a difference of over 1.5 log units. Similarly, terfenadine has an experimental log(intrinsic solubility) of -7.74 but was predicted to be -4.02 , a difference of 3.7 log units. This compound can be easily identified as an outlying point in Figure 3. Although the model gave good training statistics, because of the small number of training chemicals, predictions from this model were very poor.

Table 3. Summary of Each Model, Together with Statistics and Descriptor Definitions, Developed within Each of the Six Defined Categories

category	$N_{(\text{Train})}$	$N_{(\text{Test})}$	type	model ^a	descriptors
(1) Compounds with BOTH a hydrogen-bond-accepting N and hydrogen-bond-donating OH	6	4 (2)	QSAR	LogIS = $0.384 - 0.0248\Delta_{\text{HF}}$ $R^2_{(\text{Train})} = 0.91$, $R^2_{(\text{Test})} = 0.39$, $F = 50.39$, $\text{RMSE}_{(\text{Test})} = 4.04$	Δ_{HF} = heat of formation
(2) Carboxylic acids	26	12 (10)	QSAR	LogIS = $5.318 - 0.705 \log P - 0.0669\text{Alpha}$ $R^2_{(\text{Train})} = 0.78$, $R^2_{(\text{Test})} = 0.03$, $s = 0.59$, $F = 91.07$, $\text{RMSE}_{(\text{Test})} = 1.65$	$\log P$ Alpha = Molecular polarizability
(3) Sulfonamides (hydrogen-bond acceptors, NOT donors)	10	6	QSAR	LogIS = $2.92 - 0.733 \log P$ $R^2_{(\text{Train})} = 0.84$, $R^2_{(\text{Test})} = 0.10$, $s = 0.41$, $F = 49.04$, $\text{RMSE}_{(\text{Test})} = 0.50$	$\log P$
(4) Compounds with OH groups (hydrogen-bond acceptors and donors) but NO hydrogen-bond-accepting N	6	2	QSAR	LogIS = $5.029 - 0.747 {}^5\chi_{(\text{path})}$ $R^2_{(\text{Train})} = 0.98$, $R^2_{(\text{Test})} = 0.30$, $F = 281.59$, $\text{RMSE}_{(\text{Test})} = 2.75$	${}^5\chi_{(\text{path})}$ = Kier Chi5 (path) index
(5) Heterocyclic compounds with N sufficiently basic to be protonated at neutral pH	2	1	Quantitative Read-Across	Read-across prediction based on a single descriptor — $\log P$ experimental \log intrinsic solubility of test compound = 1.82, predicted \log intrinsic solubility = 3.13	$\log P$
(6) Hydrogen-bond-accepting N with NO hydrogen-bond donors	20	7	QSAR	LogIS = $-1.518 - 0.00872\text{IM1}(\text{Size}) - 16.694\text{Max}(\text{Q}-)R^2_{(\text{Train})} = 0.71$, $R^2_{(\text{Test})} = 0.37$, $s = 0.82$, $F = 47.57$, $\text{RMSE}_{(\text{Test})} = 0.92$	IM1(Size) = inertia moment 1 size Max(Q-) = most negative partial atomic charge

^a LogIS = log(intrinsic solubility). Numbers shown in brackets relate to the number of compounds following the removal of those too soluble to measure. ^b Correlation coefficient value for test set compounds ($R^2_{(\text{Test})}$) could not be calculated, because of an insufficient number of compounds.

Category 2: Carboxylic Acids. This was the most populated category, containing 26 training compounds and 12 test compounds. As the name suggests, the category is populated solely with carboxylic acids. Given the number of members, a traditional QSAR model was developed, culminating in a two-descriptor model, based on hydrophobicity ($\log P$) again calculated using KOWWIN³⁷ and hydrogen-bonding (Alpha) calculated using HYBOT, which is a module contained within the MOLPRO software.³⁸ As seen previously in both the MLR and ANN models, hydrophobicity, as described by the octanol:water partition coefficient ($\log P$) is a key, albeit inversely related, determinant of water solubility. As expected, the descriptor has a negative coefficient, meaning that a higher $\log P$ value (elevated lipophilicity) results in reduced aqueous solubility. The second descriptor (Alpha) is a descriptor of molecular polarizability that is strongly correlated with molecular size. As with $\log P$, Alpha shows a negative coefficient, showing that solubility decreases with increasing values of Alpha.

Unfortunately, despite promising training statistics, when the experimental solubilities were released, the predictivity of the model was shown to be surprisingly poor ($R^2_{\text{Test}} = 0.03$). However, upon closer examination, it appears that two test compounds (indomethacin and folic acid) were poorly predicted statistical outliers, both with very high leverage, as can be seen in Figure 3. Removal of these two compounds increased the test set predictivity, strikingly, to ($R^2_{\text{Test}} = 0.70$).

Category 3: Sulfonamides (Hydrogen-Bond Acceptors, NOT Donors). Ten sulfonamide derivatives were identified within the training set and six more in the test set. From these, a simple regression model was derived, based solely on hydrophobicity ($\log P$), again calculated using KOWWIN.³⁷ As observed previously, the same inverse relationship was observed between hydrophobicity and solubility, providing a model with good training statistics ($R^2_{(\text{Train})} = 0.84$).

Again, the predictive performance of the model was poor, giving an $R^2_{(\text{Test})}$ value of 0.10. Two of the six test compounds were predicted particularly poorly (hydrochlorothiazide and tolbutamide), but no obvious explanation for this could be found.

Category 4: Compounds with OH Groups (Hydrogen-Bond Acceptors and Donors) but NO Hydrogen-Bond-Accepting N. Six training compounds entered this category including naphthol, acetaminophen, and warfarin. Two test compounds (diethylstilbestrol and dipyrindamole) were also classified as belonging to this category. Again, a simple single parameter model was obtained with excellent training statistics ($R^2_{(\text{Train})} = 0.98$). The model utilized a Kier topological index (the fifth-order path index) as its sole descriptor (calculated using TSAR for Windows⁵⁴), coding for molecular size, degree of branching, flexibility and overall shape.⁵³

As with the first category, predictivity of the test set could not be determined using a correlation coefficient as only two test compounds were present. However, from considering the two predictions and visual inspection of the plot in Figure 3, predictive performance seems mixed. One test compound, diethylstilbestrol, was predicted reasonably well, with an experimental intrinsic solubility of -4.43 and a prediction of -3.66 . However, the second compound, dipyrindamole, was predicted very poorly with an experimental intrinsic solubility of -5.16 and a prediction of -7.80 . Closer inspection of the chemical structure reveals this to be a structurally distinct chemical within this category, making it outside of the applicability domain of the model. Therefore, a poor prediction is not unexpected.

Category 5: Heterocyclic Compounds with N Sufficiently Basic to Be Protonated at Neutral pH. For this category, only two training compounds (benzylimidazole and miconazole) and one test compound (thiabendazole) were identified, obviously making a traditional QSAR model impossible, because of the scarcity of data. Therefore, it was necessary to perform a quantitative read-across prediction for the single test compound. In order to make such a prediction, a quantitative property of the compound must be utilized as a scaling factor for solubility. As such, this property must be pivotal in determining the degree of a compound's solubility. Therefore, it was decided to use a property that was repeatedly being highlighted as being a fundamental determinant of solubility (namely, hydrophobicity). The $\log P$

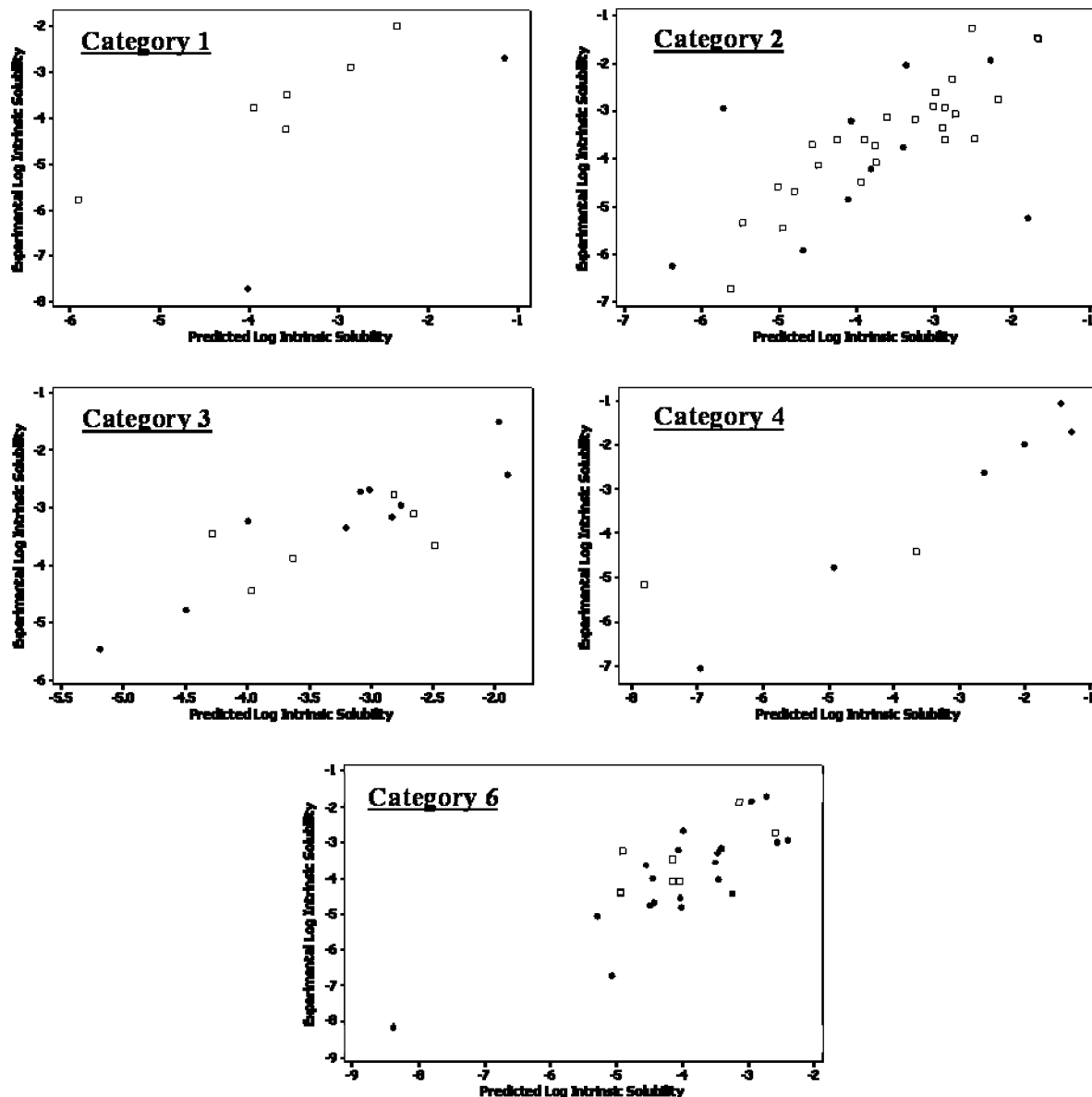


Figure 3. Plots of experimental versus predicted log(intrinsic solubility) for the five categories where QSAR models were developed (see Table 3). Each plot contains the predictions for (□) the training set and (●) the test set. Because a read-across prediction was made for a single compound within category 5, no such plot can be displayed.

values of the two training compounds were plotted against experimental solubility and a linear regression line was derived. The equation of this line was used to then predict the solubility of the test compound, by substituting in the test compound's log *P* value.

Note that this is far from an ideal read-across scenario. Preferably, two (or possibly more) training compounds would be carefully chosen from a larger dataset. These compounds would usually be those most closely resembling the test compound, in terms of both chemical structure and physicochemical properties. Ideally, one would want structurally similar training compounds which, in this example, would also have log *P* values as close as possible to that of the test compound. Furthermore, these should enclose the value of the test compound, meaning that predictions are made within the current knowledge (i.e., via interpolation, not extrapolation). However, because only two training chemicals are available in this example, no choice was available. Moreover, the range of log *P* values of the two training chemicals (2.31 and 6.25, respectively) did not cover that of the test

compound (2.00). Nevertheless, given the lack of data, such a read-across prediction was deemed necessary and would be accompanied with a cautionary warning, with regard to its accuracy.

Presubmission, there was no way to evaluate the prediction made by the model, but after the solubilities were published, it was clear that the prediction was poor. The experimental log(intrinsic solubility) of thiabendazole was determined to be -3.48 , whereas the read-across prediction was -1.17 .

Category 6: Hydrogen-Bond-Accepting N with NO Hydrogen-Bond Donors. This final category contained 20 training compounds and 7 test compounds; therefore, as with the other well-populated categories, a QSAR model was derived. A two-descriptor model was developed utilizing descriptors relating to hydrogen-bonding potential ($\text{Max}(\text{Q}-)$, calculated using HYBOT³⁸) and molecular volume/size ($\text{IM1}(\text{Size})$, calculated using Dragon Professional,³⁹ both of which are well-established factors in determining aqueous solubility. Although not as high as some of the previous models, the statistical fit was reasonable with an $R^2_{(\text{Train})}$ value

Table 4. Summary Statistics for the Four Commercial Solubility Models and Consensus Prediction, Together with Subsequently Available Simulations Plus Predictions and New Consensus Model

model	$R^2_{(\text{Test})}$	s	F	RMSE
ChemSilico (CSLogWS)	0.35	1.10	15.56	1.24
Optibrium (StarDrop)	0.28	1.15	11.69	1.38
Pharma Algorithms	0.19	1.22	7.48	1.45
Simulations Plus—SHARK	0.57	0.90	36.59	0.90
<i>in silico</i> consensus	0.38	1.07	17.54	1.14
Simulations Plus—YINAN	0.55	0.92	33.72	0.91
Simulations Plus—UIQBB	0.65	0.81	50.90	0.82
Simulations Plus—LGGAV	0.53	0.93	31.70	0.93
Simulations Plus—A69EM	0.57	0.90	36.04	0.90
Simulations Plus—NSLIC	0.54	0.93	32.60	0.93
Simulations Plus—AM108	0.61	0.85	43.66	0.87
Simulations Plus—OLASM	0.55	0.91	34.51	0.92
SPARC Solubility Model	0.50	0.96	28.17	1.56
New <i>in silico</i> consensus	0.60	0.68	41.39	0.90

of 0.66. Also, in keeping with the previous models, the predictivity of the model was relatively poor ($R^2_{(\text{Test})} = 0.37$).

In Silico Prediction. This final approach involved the test set compounds being run through several commercially available solubility models. In addition, a consensus prediction was also generated based on four predictions (those available at the time of submission). In contrast to the other modeling approaches used in this study, no training data or model development was required, as the models were already developed. As such, presubmission, the performance of each predictive model was unknown. In addition to the four models available at the time of submission, seven additional predictions were subsequently released by Simulations Plus, plus predictions generated by the freely accessible SPARC model. To further explore consensus modeling, a second consensus prediction was made using all available *in silico* models.

The results of each *in silico* model, together with summary statistics, are given in Table 4. Plots of experimental versus predicted solubility are given in Figure 4, for the four models used in the submitted consensus approach, together with the resultant consensus prediction. For completeness, although not submitted as part of the Solubility Challenge, a second consensus model was derived considering predictions from all 12 models, termed “New *in silico* consensus”.

It is evident from the statistics and plots (Table 4 and Figure 4) that the performance of the commercial models was not consistent and varied significantly between models. For the test chemicals used in this study, the Simulations Plus SHARK model gave the best predictions ($R^2_{(\text{Test})} = 0.57$), followed by the ChemSilico model (0.35). Interestingly, the Pharma Algorithms model performed very poorly in this study, with an $R^2_{(\text{Test})}$ value of just 0.19. As expected, the consensus prediction, which is a simple mean value, falls midway between the extremes ($R^2_{(\text{Test})} = 0.38$) and in this example, is not able to yield an improved prediction over any one individual model. Interestingly, the second consensus model, which was developed after the Challenge deadline, showed improved performance. The new consensus model achieved an $R^2_{(\text{Test})}$ value of 0.60. It is thought that the increase in model performance that is observed in the new consensus model is a result of the inclusion of the significantly more predictive UIQBB model and possibly the increased knowledge coded for in the SPARC model methodology. Unlike the other models, the SPARCs predic-

tions are based upon solute–solvent interactions and entropy calculations. Table 4 clearly indicates that, although the SPARC model alone does not perform as well, its use within a consensus framework, combining this knowledge with existing QSARs, results in improved solubility prediction.

DISCUSSION

With the release of the Solubility Challenge and the provision of a high-quality dataset came an opportunity to investigate our current ability to predict aqueous solubility. Because of the blind nature of the test predictions being made, this study also presented itself as a bias-free case study for the evaluation of differing modeling approaches suitable for modeling solubility. Given the need to predict solubility accurately, for a vast array of applications, any insights into ways of improving its prediction are of significant interest. Utilizing the four contrasting modeling techniques in this study (MLR, ANN, category formation, and commercial *in silico* models), the benefits of complex models (if any) and the nature of descriptors used in the prediction of solubility were evaluated.

As detailed previously, prior modeling activities have resulted in the development of numerous predictive models for aqueous solubility, utilizing various physicochemical descriptors as predictors. Of these, the major and recurring predictors include molecular size (often correlated, at least partially, with hydrophobicity, melting/boiling point and connectivity indices), polarity (hydrogen bonding), and various group/atom contributions. Although many of the models considered in this study utilize physicochemical descriptors well-known to be important in the process of aqueous solubility, the performance of the differing models varies significantly. Many of the models possess reasonable training statistics; however, when applied to external compounds, all models showed greatly reduced performance.

Previous reviews have considered numerous models available within the open literature for the prediction of aqueous solubility, together with many commercially available models.^{15,56} However, the performance of these models is often difficult to assess, because of the differing validation methods used (cross-validation or external test data) and the statistics used. One surprising finding was the variability in the performance of commercially available models. Taskinen and Norinder⁵⁶ compiled a review whereby the correlation of model predictions and experimental values was observed to vary over a range of 0.25–0.63 for a given test set. However, as with all evaluations of model performance, the makeup of the test set is of paramount importance, because it may be outside of the applicability domain of the given model. In such situations, it is not surprising that poor predictions are obtained. Unfortunately, with many of the available models, definitions of the model applicability domain are not available, because training set details are not released.

The most significant discovery from this study was the performance of the simplest modeling approach, a traditional MLR model. This simple three-descriptor model outperformed all other employed methods and, much to our surprise, was placed in the top 10 of the 99 entries in the Solubility Challenge.³¹ Unfortunately, none of the other modeling approaches used in this study was able to improve

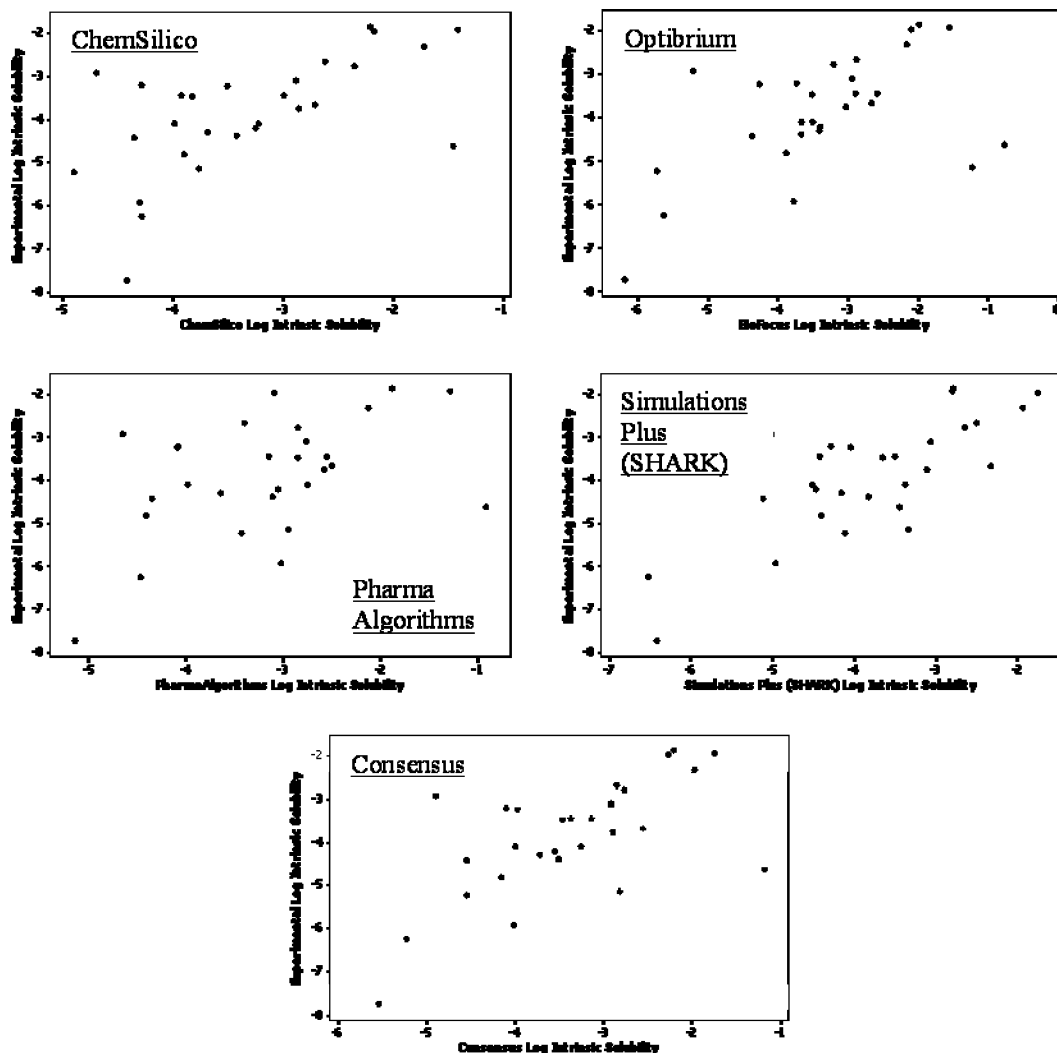


Figure 4. Plot of experimental versus predicted log(intrinsic solubility) using each of the four commercial models plus a consensus prediction.

upon the predictions made by the MLR model. Although the ANN model demonstrated promising training statistics, the model performed poorly for the validation chemicals and very poorly for the test chemicals. This pronounced loss of predictive power suggests a degree of overtraining in the model. Unfortunately, the only real solution to this problem is to provide the ANN with significantly more data points, because ANNs are data-hungry and require significantly larger datasets.⁴⁰

With regard to the category approach, the results generally were unexpectedly poor, with the models showing very low predictivity. It was hoped that splitting the chemicals into defined categories would aid in model development. However, despite successful category formation, the approach failed to generate good predictive models. Similar to the ANN approach, the primary cause of this failure was probably the lack of data. As a result of splitting the training set, each category was sparsely populated, with three categories containing <10 compounds, making the development of robust QSAR models unlikely.

In addition, the definition of suitable categories also makes this approach more challenging. The aim of forming chemical categories is to segregate a large diverse dataset into several smaller chemically or structurally restricted categories. It was hoped that the modeling of these individual categories would

reveal more predictive models. However, in addition to the paucity of data to populate these categories, the categories themselves may be too broad and not able adequately to segregate the dataset into suitable categories suitable for local modeling. If additional data were available and suitable categories were in place to split the data into series of related compounds, this approach might have been more successful. However, with the data available in this study, the category approach clearly is not able to generate meaningful models. Similarly, the read-across prediction made for category 5 is far from ideal as a result of insufficient data being available.

It must be stated here that category definition and formation can be addressed in many ways, with only one such method being investigated in this study. Therefore, it may be a weakness in the category formation method not yielding suitable categories, rather than a weakness in the general approach, which is evident in this study. Categories based on other factors (including structural similarity, mechanism of action, or physicochemical properties) may yield improved results. For example, grouping chemicals that have similar packing characteristics may further increase predictive performance. Unfortunately, investigation into this topic is outside the remit of the current study, but should be considered as a possible explanation for the observed failings of this approach.

In the final approach, the prediction of solubility using commercially available models, the results obtained were again surprisingly poor. Although previous reviews have shown the performance of commercial software to vary greatly,⁵⁶ the results obtained were below even these expectations. The results shown in Table 4 reveal that, for the test set used in this study, the performance of the commercial models varies significantly. The Simulations Plus UIQBB model was the most predictive, yielding an $R^2_{(\text{Test})}$ value of 0.65. In contrast, the Pharma Algorithms model performed very poorly, with an $R^2_{(\text{Test})}$ value of just 0.19. As one might expect, the submitted consensus prediction lies somewhere between these extremes with a $R^2_{(\text{Test})}$ value of 0.38, with the subsequently consensus model showing slightly improved predictivity. Using this approach, the weakness of one particular model is compensated by another model, therefore yielding a more robust consensus model. However, in this example, given the additional performance of the Simulations Plus UIQBB model, the consensus models were not able to improve upon the predictions of the individual model. Also, Table 4 clearly indicates that the additional Simulation Plus models show comparative performance, despite their different methodologies.

There are numerous reasons for the observed poor performance of the commercial models. The first, and most probable, is that many of the test compounds may have fallen outside of the applicability domain of the commercial models. It is highly likely that the training sets of the commercial models contained mostly non-druglike compounds, whereas the test set compounds are all drugs. Because the training data for each of these models are not provided, it is not possible to assess the applicability of each of these models. It may be that only certain compounds are within the applicability domains of each model, accounting for vast differences in modeling performance seen between compounds. Furthermore, the differing training sets used in the development of these models most probably account for the observed performance differences. As stated in the Methods section, detailed information on the training data for each model is not available. When available, this information is restricted simply to the number of compounds and, in some cases, whether they are diverse in structure or druglike, etc. As a result, this makes determining the applicability domains of these models impossible. Therefore, if any of the test set compounds are outside of the models' applicability domain, it is not possible to identify and comment on these in this present study.

In addition to differing training data, it is also possible that differences in model design and architecture could be responsible for differing performance. Although this study has shown that more complex methods offer little, if any, benefit in the prediction of solubility, because of their reduced transparency, more complex methods are also more open to misuse (e.g., the observed overtraining of the ANN model).

With the exception of many of the commercial solubility models, each model was developed using the 97-compound training set provided by the Solubility Challenge.³¹ Following the release of the experimental solubility measurements for the test set compounds, it is possible to consider the issue of applicability domain coverage in detail. If a prediction is made for a compound falling outside of the applicability domain of a model (as dictated by the training data), any

prediction will be via extrapolation outside of known data, potentially making that prediction highly dubious. To investigate whether the applicability domain is an issue, six compounds were studied in detail. Three of these compounds (bendroflumethiazide, imiprimine, and tolbutamide) were chosen because they were predicted well; the remaining three compounds (folic acid, indomethacin, and probenecid) were predicted poorly by all models. The latter two compounds were highlighted in the findings of the Solubility Challenge as being the two least-well-predicted compounds, when considering all entries with correct predictions (within 0.5 log units) in only 4% and 2% of cases, respectively.³²

Using structural similarity as a basis for applicability domain definition, the position of these compounds within the applicability domain of the training set was investigated. Utilizing the freely available structural similarity tool Toxmatch (Version 1.05),⁴⁹ it became evident that the applicability domain of the model is a crucial determinant of prediction accuracy. The Toxmatch software is able to assess structural similarity using multiple methodologies including structural fragment and atom environment approaches. Full details of the methods used in this study are given in the Methods section. For the three compounds considered that were well-predicted, these clearly lay within the applicability domain, by the presence of structurally similar compounds in the training set. For example, the test chemical imiprimine has two analogous compounds in the training set (desipramine and trimipramine), both of which are identified with similarity coefficients of 1. Similarly, bendroflumethiazide and tolbutamide are also structurally similar to many training set compounds. In contrast, those chemicals with poor predictions (folic acid, indomethacin, and probenecid) are, upon closer examination, not represented by the training compounds. Therefore, this means that these compounds are outside of the predictive domain of the models derived using this training set; hence, it is of little surprise to find that these compounds are poorly predicted. Using folic acid as an example, the experimental log(intrinsic solubility) was determined to be -5.23 ; however, the closest prediction was that by the Simulation Plus (SHARK) model (-4.12 , which is a difference of over one log unit). The relationship between position within predictive space (i.e., in or out of the applicability domain of the model) and the model performance was striking. The importance of considering the applicability domain of a model when making predictions is clear, because it can be the cause of many erroneous predictions.

This study set out to investigate differing methods of predicting aqueous solubility, in an effort to make recommendations on which modeling approach(es) are better-suited to the prediction of solubility. Given that solubility is a global molecular property and is reliant upon, and dictated by, multiple physicochemical properties, it was expected that more complex and nonlinear modeling methods (such as ANN models) may have proved superior. However, the findings of this study reveal no such benefits, showing instead that simple modeling methods (such as MLR) perform as well as, if not outperform, more complex methods. A possible reason for this is the absence of descriptors that code for a particular physicochemical property important in solubility. Although a vast array of descriptors was used, spanning a plethora of properties (including hydrophobicity, hydrogen

bonding, charge, and molecular shape and size), it is possible that other property descriptors, not covered by the current descriptor set, are lacking. Similarly, it may be that the optimal modeling approach has not yet been applied. Although this study tried to cover a wide range of modeling approaches, concentrating on those most often utilized in the literature, other approaches may prove more successful. For example, a more complex consensus modeling approach using weighted model contributions could be employed. Possibly, a consensus model utilizing predictions derived from a greater number of commercial models, covering a greater range of modeling methods and physicochemical properties, could possibly capture the strengths offered by differing modeling approaches to yield a superior consensus prediction.

Other things being equal, a simple model offering greater transparency, both in terms of model architecture and physicochemical relevance, is sought. Such a model also has the benefit of having increased portability—and, therefore, enhanced applicability—making it more accessible to potential users and more acceptable to regulators.

Finally, the design of the datasets utilized in the solubility challenge should be considered. The predictions made from all models clearly indicate that solubility is very difficult to predict accurately. Even the best model obtained in this study resulted in multiple erroneous predictions (see Figure 1). Indeed, this is a recurrent theme in the QSAR area, and it is generally accepted that the accuracy of past and present solubility models is limited by the accuracy of the experimental solubility data.⁵⁴ Therefore, the datasets provided in the Solubility Challenge (both training and test) possibly were not ideally suited to the task. Although the solubility measurements for the datasets were described as being “highly accurate”,³¹ the datasets did pose several problems to the QSAR modeler. First, and most importantly, the quality of the data (the experimental intrinsic solubility measurements), may not have been consistent. Although the Chasing Equilibrium technique used does ensure that thermodynamic equilibrium is reached, maximizing the accuracy of the results, it was clearly stated that, where literature data were available, these were not always in close agreement with these measured values, despite the associated error being stated as ~ 0.05 log units.³¹ Although solubility is a difficult property to measure accurately, especially for those compounds with very high or low solubilities, the presence of contradictory data in the literature questions the accuracy of some solubility determinations. For example, the test compound indomethacin, which is one of the worst predicted compounds in the Solubility Challenge, was said to have an intrinsic solubility of $410 \mu\text{g/mL}$. However, indomethacin is considered to be practically insoluble in water⁵⁷ and has reported solubilities in the range of $2.0\text{--}7 \mu\text{g/mL}$.^{58,59} This inconsistency brings into question the quality of the solubility measurements, or, perhaps more specifically, the place of compounds with conflicting solubility profiles in the test set.

Several compounds within both the training and test sets were also determined to be too soluble to measure (TSTM). In our opinion, these compounds should not be used within the test set, because these compounds were simply dismissed in later analysis. Finally, four compounds in the training set and three compounds in the test set were shown to exhibit polymorphism. In this study, as with many QSAR studies,

this was not addressed. The effect of polymorphs on solubility and the determination for which polymorph solubility is being measured was not considered. From the follow-up publication of the Solubility Challenge, it was stated that polymorphism was not considered by any contestant and no predictions were made as a function of polymorphic state.³² A final consideration is that of dataset size. Because many models, including several of the *in silico* models used in this study, are derived using a training set that consists of thousands of chemicals, the datasets made available in this study are rather restricted in comparison. In particular, the ANN approaches would benefit greatly from additional data, possibly allowing for the detection of more complex relationships between solubility and molecular properties. As is so often the case in SAR/(Q)SAR studies, the search for, and collection of, suitable data is paramount. Therefore, it is possible that the lack of data is limiting model performance.

CONCLUSIONS

The final conclusion of this study must be that predicting aqueous solubility is indeed still a formidable challenge! The results obtained in this study clearly indicate that no one model was able to predict solubility accurately. What is clear is that there are many approaches one can take in the development of a predictive solubility model, with varying levels of complexity.

This study, and the Solubility Challenge itself, have acted as a good review of the available approaches and some of the problems encountered when performing such modeling. The results of this study suggest that, for the prediction of aqueous solubility, more complex modeling methods do not necessarily yield better predictions and the added complexity involved is not warranted. Simple MLR models, based on interpretable physicochemical descriptors, resulted in equivalent levels of predictability.

Therefore, the findings of this study illustrate that the limitations in the current ability to predict aqueous solubility are probably not a result of inadequate modeling methodologies, but are more probably a result of an insufficient appreciation for the complexity of the solubility phenomenon. Given the limited success to date and the limitations posed by the reliability (reproducibility) of experimental data, the prediction of aqueous solubility is likely to remain a difficult area. Nonetheless, it is an issue that is still of great importance, with aqueous solubility being at the heart of pharmaceutical design, ecotoxicology and mammalian toxicological evaluation, and risk assessment. Given the dependence of *in silico* methods on high-quality data, more emphasis should be placed on obtaining more high-quality data from which to develop predictive models.

RECOMMENDATIONS

- The first step in solubility prediction should be to consider a simple modeling approach. It is suggested by the current study that the progression to more complex modeling methods may not be warranted and may yield predictions no better than those obtained from simple modeling approaches.

- The consideration of data quality is paramount. Even with the “high-quality” dataset provided in the Solubility

Challenge, questions regarding the quality of the data are raised. The data used in this study are representative of real-life datasets in that they have their limitations. Knowledge regarding the limitations of one's data is vital in acknowledging the limitations of any subsequent model. Therefore, one should always strive to consider the accuracy and/or error observed within one's data, because this limits model quality.

- The predictive domain of the model (applicability domain) is a critical consideration when assessing the reliability of a prediction. Although predictions made outside of a model's applicability domain are not necessarily "wrong", such predictions are considerably less reliable and should be treated with extreme caution. Therefore, knowledge of a model's applicability domain is fundamental when making a solubility prediction. Given the current limitations when making solubility predictions within a model's applicability domain, it would be preferable to make no prediction rather than a grossly inaccurate one.

- Although far from ideal, currently available solubility models are useful for the initial screening of chemicals, aiming at identifying those with undesirable solubility characteristics. If a more accurate prediction is required, assessing the agreement of predictions from multiple models may allow the user to gauge the reliability of the prediction with more confidence.

- Finally, this study has shown that, despite current predictive models being of use in screening applications, we are far from obtaining accurate high-quality predictions using current methods. With this in mind, alternative nonstatistical approaches should be fully explored. These approaches may include the study and inclusion of mechanistic reasoning (mechanism of action and mechanisms of solvation), solute packing characteristics, etc. Even if such methods were unable to improve upon current levels of predictivity, the additional information would be vital in increasing prediction confidence and maximizing the utilization of predictive solubility models.

Note Added after ASAP Publication. There were minor text errors in the version published ASAP November 2, 2009; the corrected version was published ASAP November 4 2009.

Supporting Information Available: Full data matrix containing experimental solubility data, together with individual model predictions. (Excel file.) This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- (1) Musha, A. Body water in man. I. Total body water in normal subjects and edematous patients. *Tohoku J. Exp. Med.* **1956**, *63*, 309–317.
- (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 Delivery Rev.* **1997**, *23*, 3–25.
- (3) Di, L.; Kerns, E. H. Biological assay challenges from compound solubility: strategies for bioassay optimization. *Drug Discovery Today* **2006**, *11*, 446–451.
- (4) Barker, S. A.; Khossravi, D. Drug delivery strategies for the new millennium. *Drug Discovery Today* **2001**, *6*, 75–77.
- (5) Borchardt, R. T.; Kerns, E. H.; Hageman, M.; Thakker, D. R.; Stevens, J. L., Eds. *Optimizing the "Drug-Like" Properties of Leads in Drug Discovery*; Springer: New York, 2006.
- (6) Ekins, S.; Rose, J. In silico ADME/Tox: The state of the art. *J. Mol. Graph. Model.* **2002**, *20*, 305–309.
- (7) Stegemann, S.; Leveiller, F.; Franchi, D.; de Jong, H.; Lindén, H. When poor solubility becomes an issue: From early stage to proof of concept. *Eur. J. Pharm. Sci.* **2007**, *31*, 249–261.
- (8) Kennedy, T. Managing the drug discovery/development interface. *Drug Discovery Today* **1997**, *2*, 436–444.
- (9) Palmer, A. M. New horizons in drug metabolism, pharmacokinetics and drug discovery. *Drug News Perspect.* **2003**, *16*, 57–62.
- (10) Kerns, E. H.; Li, D. *Drug-like Properties: Concepts, Structure Design and Methods: from ADME to Toxicity Optimization*; Academic Press: Boston, 2008.
- (11) Bhattachar, S. N.; Deschenes, L. A.; Wesley, J. A. Solubility: It's not just for physical chemists. *Drug Discovery Today* **2006**, *11*, 1012–1018.
- (12) Dai, W.; Pollock-Dove, C.; Dong, L. C.; Li, S. Advanced screening assays to rapidly identify solubility-enhancing formulations: High-throughput, miniaturization and automation. *Adv. Drug Delivery Rev.* **2008**, *60*, 657–672.
- (13) Alsenz, J.; Kansy, M. High throughput solubility measurement in drug discovery and development. *Adv. Drug Delivery Rev.* **2007**, *59*, 546–567.
- (14) Fühner, H. The aqueous solubility of homologous series. *Ber. Dtsch. Chem. Ges.* **1924**, *57B*, 510–515.
- (15) Dearden, J. C. In silico prediction of aqueous solubility. *Expert. Opin. Drug Discovery* **2006**, *1*, 31–52.
- (16) Reynolds, J. A.; Gilbert, D. B.; Tanford, C. Empirical correlation between hydrophobic free energy and aqueous cavity surface area. *Proc. Natl. Acad. Sci.* **1974**, *71*, 2925–2927.
- (17) Hansch, C.; Quinlan, J. E.; Lawrence, G. L. The linear free-energy relationship between partition coefficients and aqueous solubility of organic liquids. *J. Org. Chem.* **1968**, *33*, 347–350.
- (18) Faller, B.; Ertl, P. Computational approaches to determine drug solubility. *Adv. Drug Delivery Rev.* **2007**, *59*, 533–545.
- (19) Yalkowsky, S. H.; Valvani, S. C. Solubility and partitioning I: solubility of nonelectrolytes in water. *J. Pharm. Sci.* **1980**, *69*, 912–922.
- (20) Kamlet, M. J.; Doherty, R. M.; Abboud, J.-L. M.; Abraham, M. H.; Taft, R. W. Linear solvation energy relationships: 36. Molecular properties governing solubilities of organic nonelectrolytes in water. *J. Pharm. Sci.* **1986**, *75*, 338–348.
- (21) Irmann, F. A simple correlation of water solubility and structure of hydrocarbons and hydrocarbon halides. *Chem.-Ing.-Tech.* **1965**, *37*, 789–798.
- (22) Wakita, K.; Yoshimoto, M.; Miyamoto, S.; Watanabe, H. A method for calculation of the aqueous solubility of organic compounds by using new fragment solubility constants. *Chem. Pharm. Bull.* **1986**, *34*, 4663–4681.
- (23) Randic, M. Characterization of molecular branching. *J. Am. Chem. Soc.* **1975**, *97*, 6609–6615.
- (24) Kier, L. B.; Hall, L. H. *Molecular Connectivity in Chemistry and Drug Design*; Academic Press: New York, 1976.
- (25) Kier, L. B.; Hall, L. H. *Molecular Connectivity in Structure–Activity Analysis*; John Wiley & Sons: New York, 1986.
- (26) 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.
- (27) Liu, R.; So, S.-S. Development of quantitative structure-property relationship models for early ADME evaluation in drug discovery. 1. Aqueous solubility. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1633–1639.
- (28) Wegner, J. K.; Zell, A. Prediction of aqueous solubility and partition coefficient optimized by a genetic algorithm based descriptor selection method. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1077–1084.
- (29) Delaney, J. S. Predicting aqueous solubility from structure. *Drug Discovery Today* **2005**, *10*, 289–295.
- (30) Costa, M. F.; Pádua, A. A. H. In *Developments and Applications in Solubility*; Letcher, T. M., Ed.; Royal Society of Chemistry: London, 2007; Chapter 10, pp 153–170.
- (31) Llinàs, A.; Glen, R. C.; Goodman, J. M. Solubility challenge: Can you predict solubilities of 32 molecules using a database of 100 reliable measurements. *J. Chem. Inf. Model.* **2008**, *48*, 1289–1303.
- (32) Hopfinger, A. J.; Esposito, E. X.; Llinàs, A.; Glen, R. C.; Goodman, J. M. Finding of the Challenge to predict aqueous solubility. *J. Chem. Inf. Model.* **2009**, *49*, 1–5.
- (33) Stuart, M.; Box, K. Chasing equilibrium: Measuring the intrinsic solubility of weak acids and bases. *Anal. Chem.* **2005**, *77*, 983–990.
- (34) ChemIDPlus Advanced. <http://chem.sis.nlm.nih.gov/chemidplus/> (accessed April 22, 2009).
- (35) ChemBioFinder. <http://chembiofinder.cambridgesoft.com/chembiofinder/SimpleSearch.aspx> (accessed April 22, 2009).
- (36) ChemAxon. <http://www.chemaxon.com/> (accessed April 22, 2009).
- (37) US EPA Estimation Program Interface (EPI) Suite. <http://www.epa.gov/oppt/exposure/pubs/episuite.htm> (accessed April 22, 2009).
- (38) Raevsky, O. A.; Grigor'ev, V. Ju. Trepalin, S. V. HYBOT (HYdrogen Bond Thermodynamics) Program Package (Version 2.1.0.706). Reg-

- istration by Russian State Patent Agency N 990090 of 26.02.99 (Raevsky, O. A., Skvortsov, V. S., Grigor'ev, V. Ju. Trepalin, S. V.).
- (39) Dragon Professional Software Package, Version 5.3 for Windows; Milano Chemometrics and QSAR Research Group; Milano, Italy, 2009.
- (40) Jorgensen, W. L.; Duffy, E. M. Prediction of drug solubility from structure. *Adv. Drug Delivery Rev.* **2002**, 355–366.
- (41) Mobydigs—Software for Multilinear Regression Analysis and Variable Selection by Genetic Algorithm, Version 1.0 for Windows; Milano Chemometrics and QSAR Research Group, Milano, Italy, 2006.
- (42) Topliss, J. G.; Edwards, R. P. Chance factors in studies of quantitative structure–activity relationships. *J. Med. Chem.* **1979**, 22, 1238–1244.
- (43) Butina, D.; Segall, M. D.; Frankcombe, K. Predicting ADME properties *in silico*: Methods and models. *Drug Discovery Today* **2002**, 7, S83–S88.
- (44) Statistica Statistical Software for Windows (Version 6.1); StatSoft, Inc., Tulsa, OK, 2004.
- (45) Enoch, S. J.; Cronin, M. T. D.; Schultz, T. W.; Madden, J. C. An evaluation of global QSAR models for the prediction of the toxicity of phenols to *Tetrahymena pyriformis*. *Chemosphere* **2008**, 71, 1225–1232.
- (46) Ghasemi, J.; Saaidpour, S. Quantitative structure–property relationship study of *n*-octanol–water partition coefficients of some of diverse drugs using multiple linear regression. *Anal. Chim. Acta* **2007**, 604, 99–106.
- (47) Hilal, S. H. Karickhoff, S. W. Verification and validation of the SPARC model. U.S. Environmental Protection Agency Report EPA/600/R-03/033. (Available via the Internet at http://www.epa.gov/athens/publications/reports/EPA_600_R_03_033.pdf, accessed September 23, 2009.)
- (48) Minitab for Windows Statistical Software, Version 15; Minitab, Inc., State College, PA, 2007.
- (49) Toxmatch, Version 1.06; IdeaConsult, Ltd., Sofia, Bulgaria, 2008.
- (50) Enoch, S. J.; Cronin, M. T. D.; Madden, J. C.; Hewitt, M. Formation of structural categories to allow for read-across for teratogenicity. *QSAR Comb. Sci.* **2009**, 28, 696–708.
- (51) Shacham, M.; Brauner, N.; Cholakov, G.; Stateva, R. P. Identifying applicability domains for quantitative structure property relationships. Report from the 17th European Symposium on Computer Aided Process Engineering, 2007. (Available via the Internet at ftp://ftp.bgu.ac.il/shacham/publ_papers/Escape17_327_07.pdf, accessed April 22, 2009.)
- (52) Massart, D. L. *Chemometrics: A Textbook*, 5th Edition; Elsevier: Amsterdam, 1988.
- (53) Consonni, V.; Todeschini, R.; Pavan, M. Structure/response correlations and similarity/diversity analysis by GETAWAY descriptors. 1. Theory of the novel 3D molecular descriptors. *J. Chem. Inf. Comput. Sci.* **2002**, 42, 682–692.
- (54) *TSAR for Windows Version 3.3*; Accelrys, Inc.: Oxford, England, 2009.
- (55) Ebbing, D. D.; Gammon, S. D. *General Chemistry*, 9th Edition; Houghton Mifflin: Boston, 2009; Chapter 6, pp 223–262.
- (56) Taylor, J. B. Trigg, D. J. In *ADME-Tox Approaches*; Testa, B. van de Waterbeemd, H., Eds.; Comprehensive Medicinal Chemistry II, Vol. 5; Elsevier: Oxford, 2007; pp 627–648.
- (57) O'Neil, M. J. *The Merck Index*, 13th Edition; Merck & Co., Inc.: Whitehouse Station, NJ, 2001.
- (58) Avdeef, A. Physicochemical profiling (solubility, permeability and charge state). *Curr. Top. Med. Chem.* **2001**, 1, 277–351.
- (59) Nokhodchi, A. The effect of type and concentration of vehicles on the dissolution of a poorly soluble drug (indomethacin) from liquisolid compacts. *J. Pharm. Pharmaceut. Sci.* **2005**, 8, 18–25.

CI900286S