

Combinatorial QSAR Modeling of Chemical Toxicants Tested against *Tetrahymena pyriformis*

Hao Zhu,[†] Alexander Tropsha,^{*,†} Denis Fourches,[‡] Alexandre Varnek,[‡] Ester Papa,[§] Paola Gramatica,[§] Tomas Öberg,^{||} Phuong Dao,[⊥] Artem Cherkasov,[⊥] and Igor V. Tetko^{#,∇}

Laboratory for Molecular Modeling, Division of Medicinal Chemistry and Natural Products, and Carolina Exploratory Center for Cheminformatics Research, School of Pharmacy, CB 7360, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, Laboratories of Chemoinformatics, Institute of Chemistry, Louis Pasteur University, Strasbourg, France, QSAR Research Unit in Environmental Chemistry and Ecotoxicology, Department of Structural and Functional Biology, University of Insubria, Varese, Italy, School of Pure and Applied Natural Sciences, University of Kalmar, SE-391 82 Kalmar, Sweden, Division of Infectious Diseases, Faculty of Medicine, University of British Columbia, 2733 Heather Street, Vancouver, British Columbia, V5Z 3J5, Canada, Helmholtz Center Munich—German Research Center for Environmental Health, Institute for Bioinformatics, Neuherberg, D-85764, Germany, and Institute of Bioorganic & Petrochemistry, Murmanskaya 1, Kyiv-94, 02660, Ukraine

Received November 28, 2007

Selecting most rigorous quantitative structure–activity relationship (QSAR) approaches is of great importance in the development of robust and predictive models of chemical toxicity. To address this issue in a systematic way, we have formed an international virtual collaboratory consisting of six independent groups with shared interests in computational chemical toxicology. We have compiled an aqueous toxicity data set containing 983 unique compounds tested in the same laboratory over a decade against *Tetrahymena pyriformis*. A modeling set including 644 compounds was selected randomly from the original set and distributed to all groups that used their own QSAR tools for model development. The remaining 339 compounds in the original set (external set I) as well as 110 additional compounds (external set II) published recently by the same laboratory (after this computational study was already in progress) were used as two independent validation sets to assess the external predictive power of individual models. In total, our virtual collaboratory has developed 15 different types of QSAR models of aquatic toxicity for the training set. The internal prediction accuracy for the modeling set ranged from 0.76 to 0.93 as measured by the leave-one-out cross-validation correlation coefficient (Q_{abs}^2). The prediction accuracy for the external validation sets I and II ranged from 0.71 to 0.85 (linear regression coefficient R_{absI}^2) and from 0.38 to 0.83 (linear regression coefficient R_{absII}^2), respectively. The use of an applicability domain threshold implemented in most models generally improved the external prediction accuracy but at the same time led to a decrease in chemical space coverage. Finally, several consensus models were developed by averaging the predicted aquatic toxicity for every compound using all 15 models, with or without taking into account their respective applicability domains. We find that consensus models afford higher prediction accuracy for the external validation data sets with the highest space coverage as compared to individual constituent models. Our studies prove the power of a collaborative and consensual approach to QSAR model development. The best validated models of aquatic toxicity developed by our collaboratory (both individual and consensus) can be used as reliable computational predictors of aquatic toxicity and are available from any of the participating laboratories.

1. INTRODUCTION

Chemical toxicity can be associated with many hazardous biological effects such as gene damage, carcinogenicity, or the induction of lethal rodent or human diseases. It is important to evaluate the toxicity of all commercial chemicals, especially the high production volume compounds as

well as drugs or drug candidates, before releasing them into the market. To address this need, standard experimental protocols have been established by the chemical industry, pharmaceutical companies, and government agencies to test chemicals for their toxic potential. For example, a so-called “Standard Battery for Genotoxicity Test” was established by the International Conference on Harmonization, U.S. Environmental Protection Administration (EPA), U.S. Food and Drug Administration, and other regulatory agencies. This test includes one bacterial reverse mutation assay (e.g., *Salmonella typhimurium* mutation test), one mammalian cell gene mutation assay (e.g., mouse lymphoma cell mutation test), and one *in vivo* micronucleus test. The test battery varies slightly for pharmaceutical compounds, industrial com-

* Corresponding author. Phone: +1 (919) 966-2955. Fax: +1 (919) 966-0204. E-mail: alex_tropsha@unc.edu.

[†] University of North Carolina at Chapel Hill.

[‡] Louis Pasteur University.

[§] University of Insubria.

^{||} University of Kalmar.

[⊥] University of British Columbia.

[#] German Research Center for Environmental Health, Institute for Bioinformatics.

[∇] Institute of Bioorganic & Petrochemistry.

pounds, and pesticides. The current strategies and guidelines for toxicity testing have been described.¹

Although the experimental protocols for toxicity testing have been developed for many years and the cost of compound testing has reduced significantly, computational chemical toxicology continues to be a viable approach to reduce both the number of efforts and the cost of experimental toxicity assessment.² Significant savings could be achieved if accurate predictions of potential toxicity could be used to prioritize compound selection for experimental testing. Many quantitative structure–activity relationship (QSAR) studies have been conducted for different toxicity end points to address this challenge.^{3–6}

The most critical limitation of many traditional QSAR studies is their low external predictive power, that is, their ability to predict accurately the underlying end point toxicity for compounds that were not used for model development. The low external prediction accuracy of QSAR models in spite of the high accuracy of the training set models is a well-known phenomenon^{3,7,8} frequently referred to as the Kubinyi paradox.^{9,10} There could be many reasons for this discrepancy between the internal (fitness) and external predictive power of QSAR models. The most common is that training set models are based on data interpolation, and therefore they inherently have limited applicability in the chemistry space, whereas any external prediction implies inherent and, frequently, excessive extrapolation of the training set models. Poor external predictive power of QSAR models could be due to the lack of or incorrect use of external validation during the modeling process. Furthermore, each statistical method used in QSAR studies has its specific advantages, weaknesses, and practical constraints, so it is important to select the most suitable QSAR methodology for a specific toxicity end point. We have addressed some of these problems in our earlier publications.^{8,11}

In this paper, we report on the results of combinatorial QSAR modeling of a diverse series of organic compounds tested for aquatic toxicity in *Tetrahymena pyriformis* in the same laboratory over nearly a decade.^{12–18} This computational study was conducted in collaboration between six academic groups specializing in cheminformatics and computational toxicology. The common goals for our virtual collaboratory were to explore the relative strengths of various QSAR approaches in their ability to develop robust and externally predictive models of this particular toxicity end point. We have endeavored to develop the most statistically robust, validated, and externally predictive QSAR models of aquatic toxicity. The members of our collaboratory included scientists from the University of North Carolina at Chapel Hill in the United States (UNC), University of Louis Pasteur (ULP) in France, University of Insubria (UI) in Italy, University of Kalmar (UK) in Sweden, Virtual Computational Chemistry Laboratory (VCCLAB) in Germany, and the University of British Columbia (UBC) in Canada. Each group relied on its own QSAR modeling approaches to develop toxicity models using the same modeling set, and we agreed to evaluate the realistic model performance using the same external validation set(s). Thus, this study presents an example of a fruitful international collaboration between researchers that use different techniques and approaches but share general principles of QSAR model development and validation. Significantly, we did not make any assumptions

about the purported mechanisms of aquatic toxicity yet were able to develop statistically significant models for all experimentally tested compounds. In this regard, it is relevant to cite an opinion expressed in an earlier publication by Schultz that “models that accurately predict acute toxicity without first identifying toxic mechanisms are highly desirable”.¹³

We were excited to observe that the consensus model integrating all validated individual models was found to be the most externally predictive. Our results indicate that consensus models could be used as reliable predictors of aquatic toxicity for chemical compounds. In addition to the scientific merits of our investigations, we believe that this study presents a model of collaboration that integrates the expertise of participating laboratories toward establishing best practices and reliable solutions for difficult problems in chemical toxicology.

2. METHODS

2.1. Data Sets. The growth inhibition of the ciliated protozoan *T. pyriformis* is a commonly accepted toxicity screening tool that has been under development and implementation by Schultz and co-workers for many years.¹² In the past 10 years, this group has published the results from the standard *T. pyriformis* toxicity test protocol for more than 1000 different compounds providing a unique data set for modeling aquatic toxicity.

The *T. pyriformis* toxicity data set used in this study was compiled from several publications of the Schultz group^{12,14–17} as well as from data available at the Tetratox database Web site (<http://www.vet.utk.edu/TETRATOX/>). The data were collected from publicly available sources and may not include all test results from the Schultz laboratory. We will make every attempt to enrich this data collection as additional experimental data become available and use the new data as an external validation set in future studies. After deleting duplicates as well as several compounds with conflicting test results and correcting several chemical structures in the original data sources, our final data set included 983 unique compounds (the structural information is included in the Supporting Information). The *T. pyriformis* toxicity of each compound was expressed as the logarithm of 50% growth inhibitory concentration (pIGC50) values. For the purposes of this study, the data set was randomly divided into two parts: (1) the modeling set of 644 compounds and (2) the validation set including 339 compounds. The former set was used for model development by each participating group, and the latter set was used to estimate the external prediction power of each model as a universal metric of model performance.

When this project was already well underway, a new data set had become available from the most recent publication by the Schultz group.¹⁸ It provided us with an additional external set to evaluate the predictive power and reliability of all QSAR models. Among compounds reported in ref 18, 110 were unique, that is, not present among the original set of 983 compounds; thus, these 110 compounds formed the second independent validation set for our study. Figure 1 shows the activity distributions of compounds in both the training and the two validation sets. Obviously, all three data sets consist of similar fractions of compounds with low,

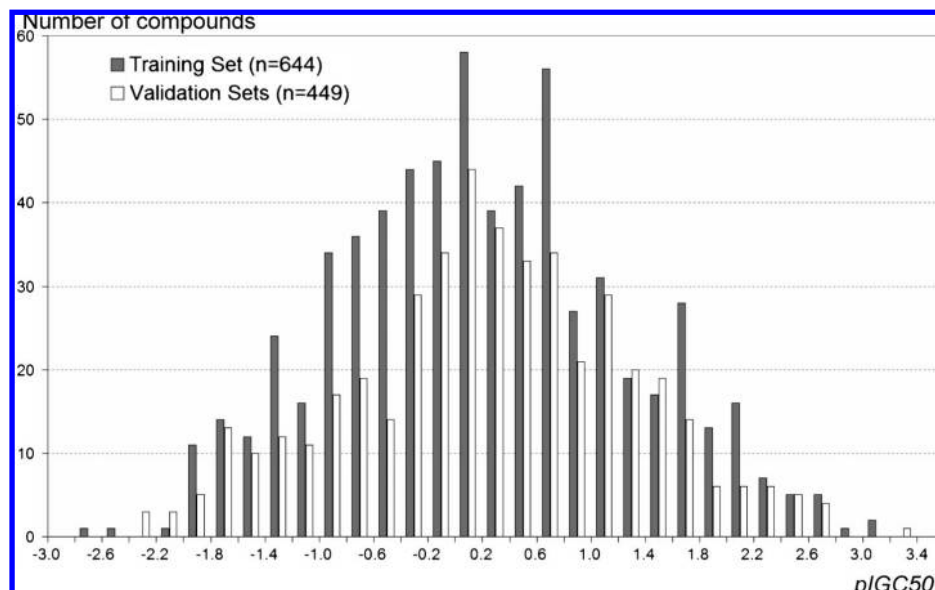


Figure 1. Experimental pIC_{50} for the training set of 644 molecules and for the combined validation set of 449 molecules.

Table 1. Overview of QSAR Modeling Approaches Employed by Six Cheminformatic Groups Involved in This Study

group ID	modeling techniques	descriptor type	applicability domain definition
UNC	kNN, SVM	MolconnZ, Dragon	Euclidean distance threshold between a test compound and compounds in the modeling set
ULP	MLR, SVM, kNN	fragments (ISIDA), molecular (CODESSA-Pro)	Euclidean distance threshold between a compound and compounds in the modeling set; bounding box
UI	MLR/OLS	Dragon	leverage approach
UK	PLS	Dragon	residual standard deviation and leverage within the PLSR model
VCCLAB	ASNN	E-state indices	maximal correlation coefficient of the test molecule to the training set molecules in the space of models
UBC	MLR, ANN, SVM, PLS	IND_I	undefined

intermediate, and high toxicity values (expressed as pIC_{50}). A complete list of the compounds in all three data sets is provided in Supporting Information Table 1.

2.2. Universal Statistical Figures of Merit for All Models. Different groups have employed different techniques and (sometimes) different statistical parameters to evaluate the performance of models developed independently for the modeling set (described below). To harmonize the results of this study, the same standard parameters were chosen to describe each model's performance as applied to the modeling and external test set predictions. Thus, we have employed Q_{abs}^2 (the squared leave-one-out cross-validation correlation coefficient) for the modeling set, R_{abs}^2 (frequently described as the coefficient of determination) for the external validations sets, and MAE (mean absolute error) for the linear correlation between predicted (Y_{pred}) and experimental (Y_{exp}) data (here, $Y = pIC_{50}$); these parameters are defined as follows:

$$Q_{abs}^2 = 1 - \sum_Y (Y_{exp} - Y_{LOO})^2 / \sum_Y (Y_{exp} - \langle Y \rangle_{exp})^2 \quad (1)$$

$$R_{abs}^2 = 1 - \sum_Y (Y_{exp} - Y_{pred})^2 / \sum_Y (Y_{exp} - \langle Y \rangle_{exp})^2 \quad (2)$$

$$MAE = \sum_Y |Y - Y_{pred}| / n \quad (3)$$

Many other statistical characteristics can be used to evaluate model performance; however, we restricted ourselves to these

three parameters that provide minimal but sufficient information concerning any model's ability to reproduce both the trends in experimental data for the test sets and the mean accuracy of predicting all experimental values. The models were considered acceptable if R_{abs}^2 exceeded 0.5.

2.3. QSAR Approaches. Each participating group has developed previously its own QSAR approaches including descriptor generation and statistical data modeling protocols. In addition, each group (with one exception) has developed and/or implemented the model-specific applicability domains (ADs) of the resulting QSAR models. The brief summary of QSAR techniques used in participating groups is given in Table 1, and major details of the techniques used in this study are described in the remaining parts of this section where references to individual methods for additional in-depth description are given.

2.3.1. Descriptors. A brief summary of descriptors used by each group is given below. In some cases, different groups used descriptors of the same type but generated with different software packages (e.g., both MolConnZ (MZ) and Dragon generate topological and electrotopological descriptors but differ in some other descriptors).

UNC. The MZ software available from Edusoft was used.¹⁹ It affords computation of a wide range of topological indices of molecular structure (e.g., molecular connectivity indices, k molecular shape indices, topological and electrotopological state indices, differential connectivity indices, etc.), but several descriptors depend upon the arbitrary

numbering of atoms in a molecule and are introduced solely for book-keeping purposes.^{20–23} The latter descriptors as well as those with zero variance across the modeling set were not used in model generation. Furthermore, due to different absolute ranges of descriptor values, range scaling was applied to all descriptors. The total number of MZ descriptors generated for the 644 compounds in the modeling set was 336.

ULP. Two types of descriptors were used: substructural molecular fragments calculated with the ISIDA program^{24–26} and molecular descriptors calculated with the CODESSA-Pro program.²⁷ All of these descriptors were derived solely from 2D chemical structure and did not require any experimental data or expensive theoretical calculations.

Two subclasses of fragment descriptors available in ISIDA were used: “sequences” and “augmented atoms”. The sequences may contain connected atoms and bonds, atoms only, or bonds only. For each type of sequences, the minimal ($n_{\min} \geq 2$) and maximal ($n_{\max} \leq 15$) number of constituent atoms is defined. An “augmented atom” represents a particular atom with its environment including either neighboring atoms and bonds, or atoms only, or bonds only. Hybridization of atoms can be taken into account.

The CODESSA-Pro program calculates several hundred molecular descriptors belonging to the following classes: constitutional, geometrical, topological, electrostatic, quantum chemical, and thermodynamic. Unlike fragment descriptors, the calculation of some molecular descriptors requires non-negligible CPU time since it involves semiempirical quantum mechanical calculations. Fragment descriptors calculated by ISIDA can also be used in CODESSA-Pro as external descriptors.

UI. A set of 929 theoretical molecular descriptors was computed using the software Dragon v.5.4.²⁸ Only simple structural descriptors, directly derived from the SMILES notation for each studied compound, were calculated; the three-dimensional (3D) descriptors were not computed. The typology of the included descriptors is as follows: 0D constitutional (atom and group counts), 1D functional groups, 1D atom-centered fragments, 2D topological descriptors, 2D walk and path counts, 2D autocorrelations, 2D connectivity indices, 2D information indices, 2D topological charge indices, 2D eigenvalue-based indices, 2D topological descriptors, 2D edge-adjacency indices, 2D Burden eigenvalues, and molecular properties. Constant and near-constant variables (178 total) were deleted. One of the pairwise more-than-98%-correlated variables (271 total) was deleted as well; thus, a final set included 480 descriptors, which were used for QSAR modeling. The procedures to calculate these descriptors and relevant references were reported previously.²⁹

UK. Three-dimensional molecular structures were calculated from the SMILES notations using the CORINA software.³⁰ The 3D molecular structures in MDL SD file format were subsequently used as input for the generation of 1664 descriptors using the Dragon v.5.4 software.²⁸ The generated molecular descriptors include those described above by the UI group as well as an additional 14 charge and 721 3D descriptors (Randic molecular profiles, geometrical, RDF, 3D-MoRSE, WHIM, and GETAWAY descriptors). These additional 735 descriptors are discussed elsewhere²⁹

VCCLAB. The electrotopological state (E-state) indices introduced by Hall and Kier^{31,32} combine both electronic and topological characteristics of the analyzed molecules. For each atom type in a molecule, the E-state index values are summed and used in a group contribution manner. In this study, we have used an extended set of atom-type E-state indices which was developed to improve the coverage of functional groups and the neighborhood of nitrogen and oxygen atoms.^{33,34} The atom-type E-state indices were calculated using the in-house program.^{35,36} Similarly to our previous studies,³⁴ molecular weight (MW) and the number of non-hydrogen atoms were used as additional descriptors. In addition, considering the importance of lipophilicity, we also included topological polar surface area, the number of hydrogen-bond-acceptor atoms, and the number of hydrogen-bond-donor atoms as three additional parameters.

UBC. The “inductive” descriptors IND_I were developed in a series of papers by Cherkasov and coauthors.^{37–39} These molecular parameters are based on the models of inductive and steric effects, inductive electronegativity, and molecular capacitance and are accessed from fundamental parameters of bound atoms, such as absolute electronegativities (χ), covalent radii (R), and intramolecular distances (r).

This approach allows one to compute as many as 50 “inductive” QSAR descriptors (cf. refs 37–39 for additional details). It should be mentioned that when all 50 “inductive” descriptors are computed for conventional chemical data sets there is typically a cross-correlation between the parameters, as some of the descriptors reflect closely related properties. Thus, highly correlated “inductive” descriptors should be typically eliminated prior to creating QSAR models.

In the present study, all 50 IND_I parameters were computed for all compounds from the modeling and the two validation sets. A separate cross-correlation analysis was not conducted since the descriptor-generating software include scripts for multiple linear regression (MLR), artificial neural network (ANN), support vector machines (SVM), k -nearest neighbor (k NN), and linear discriminant analysis modeling that take care of this problem inherently. The descriptors were computed from 3D structures of molecules optimized with the MMFF molecular force field (as implemented within the MOE program⁴⁰). The SVL scripts developed in-house to compute the IND_I parameters were employed (the scripts can be freely downloaded through the SVL exchange).

2.3.2. Modeling Approaches. This section presents an overview of computational data analytical approaches used by each participating group to develop QSAR models.

UNC. Training set models were built using variable-selection k NN and SVM approaches that were developed and implemented in this group. The k NN QSAR method⁴¹ employs the k NN classification principle and the variable selection procedure. Briefly, a subset of $nvar$ (number of selected variables) descriptors is selected randomly at the onset of the calculations. The $nvar$ is set to different values, and the training set models are developed with leave-one-out cross-validation, where each compound is eliminated from the training set and its biological activity is predicted as the average activity of the k most similar molecules, where the value of k is optimized as well ($k = 1–5$). The similarity is characterized by the Euclidean distance between compounds in multidimensional descriptor space. A method of simulated annealing with the Metropolis-like acceptance

criteria is used to optimize the selection of variables. The objective of this method is to obtain the best leave-one-out cross-validated Q_{abs}^2 possible by optimizing n_{var} and k . The additional details of the method can be found elsewhere.⁴¹

SVM was developed by Vapnik⁴² as a general data modeling methodology where both the training set error and the model complexity are incorporated into a special loss function that is minimized during model development. The methodology allows one to regulate the importance of the training set error versus the model complexity to develop the optimal model that best predicts a test set. Later, SVM was extended to afford the development of SVM regression models for data sets with noninteger activities, such as QSAR.

In our studies, the linear SVM was used. The performance of SVM depends on the selection of several internal parameters of the algorithm (C and ϵ). To find models with the highest accuracy for both training and test sets, the calculations were carried out for all combinations of C and ϵ , with the C value varying from 0.1 to 100 with a step of 10 and ϵ varying from 0.0 to 0.5 with a step of 0.34. For example, if the total number of training/test sets generated for one type of descriptors was 36, $36 \times 10 \times 2 = 720$ models were constructed. Further details of the k NN and SVM method implementation are given elsewhere.^{43,44}

As emphasized in our earlier reports,^{8,11} training-set-only modeling is insufficient to achieve models with validated predictive power. For this reason, the 644-compound modeling set was divided into multiple training/test sets using the sphere-exclusion algorithm.⁴⁵ For each collection of descriptors, the modeling set was divided into 36–50 training/test sets of different relative sizes. Both k NN and SVM QSAR toxicity models were developed using training set data only, which was part of the modeling set, and the resulting models were validated by predicting the toxicity of compounds in the corresponding test sets. Therefore, the statistical significance of the k NN and SVM QSAR toxicity models was characterized not only with the cross-validated Q_{abs}^2 for the training sets but also with a linear fit R_{abs}^2 for the test sets. The model acceptability thresholds in this study were $Q_{\text{abs}}^2/R_{\text{abs}}^2 > 0.75/0.75$; that is, only models that met these criteria were kept and used for the consensus prediction of new compounds. The importance of this procedure was discussed previously.⁴⁵

ULP. The pool of fragment and molecular descriptors is much larger than the number of compounds in the training set; therefore, a variable-selection technique should be applied to build statistically significant multilinear regressions. In CODESSA-Pro, the forward stepwise procedure ("best multilinear regression"⁴⁶) is applied to select a limited number of descriptors. A more sophisticated technique is implemented in ISIDA-MLR, that is, the forward stepwise procedure, which selects a user-defined number of descriptors (usually 60–80% from the size of the training set) followed by t test backward stepwise selection.⁴⁷ The optimized descriptor subset was used by either ISIDA-MLR or CODESSA-Pro to build a multilinear correlation equation in the form of $\text{pIGC}_{50} = a_0 + \sum a_i \times X_i$, where X_i is the value of i th descriptor, a_i is its contribution, and a_0 is a descriptor-independent term. For fragment descriptors, X_i is the occurrence of the i th fragment.

To obtain ISIDA/SVM models, we used descriptors selected as described for ISIDA/MLR. ISIDA implements the open-source LibSVM package to build ISIDA-SVM models for the training set.

Similar to the k NN approach used by UNC, ISIDA- k NN assumes that similar compounds have similar properties: the target property of a compound is calculated as a distance-weighted mean of property values for its k nearest neighbors in the chemical space. However, ISIDA- k NN implements a different approach for variable selection: it uses an original stepwise algorithm, which iteratively selects pools of descriptors leading to reliable k NN models. The number of variables in pools is increased step by step according to the LOO cross validation coefficient (Q_{abs}^2) of corresponding models and a Metropolis criterion, avoiding the convergence to local solutions. Then, models are sorted according to their statistical parameters, that is, Q_{abs}^2 for the training set and R_{abs}^2 for an internal test set. Finally, selected k NN models are used to screen compounds in the external set. For each compound, the program computes the property as an arithmetic mean of values obtained with these selected k NN models; predictions that appeared as outliers within the distribution of predicted values for each compound were excluded according to the Grubbs's statistics.⁴⁸

UI. Models were built with MLR by the ordinary least-squares method (OLS) and variable selection by the genetic algorithm using the MOBY DIGS package.⁴⁹ The aim of this approach is to develop the simplest model based on the minimum number of individual molecular descriptors following the parsimony principle.

Dragon calculates a large number of descriptors in order to capture all possible diverse structural information for the underlying data set, making it practically impossible to employ a MLR approach without variable selection. Thus, genetic algorithm—variable subset selection (GA-VSS)⁵⁰ was applied to the input set of 480 descriptors to select the most relevant subsets that afford models with the highest predictive power in modeling the studied end point. The outcome of the GA-VSS procedure is a population of 100 regression models, ordered according to their decreasing internal predictive performance as estimated by the leave-one-out cross-validated correlation coefficient Q_{abs}^2 . As the first step, all of the models with one or two variables were developed by the all-subset-method procedure in order to explore all low-dimensional QSAR models. The number of descriptors was subsequently increased one by one, by GA selection, and new models were formed. The GA optimization was terminated when increasing the model size did not increase the Q_{abs}^2 value to any significant degree. In this study, the best tradeoff between complexity and predictive power was obtained for models including only six individual molecular descriptors. Particular attention was paid to the collinearity of the selected molecular descriptors: in fact, to avoid multicollinearity, regression was calculated only for variable subsets with an acceptable multivariate correlation with a response, by applying the Q under influence of K (QUIK) rule.⁵¹

According to this rule, only those models with a global correlation of the $[X + y]$ block (K_{XY}) greater than the global correlation of the X block (K_{XX}) variable (X being the molecular descriptors and y the response variable) were considered acceptable.

Moreover, the bootstrapping approach,⁵² repeated 5000 times for each validated model, was applied to avoid an overestimation of model predictive power and to verify its robustness and internal predictivity (Q_{BOOT}^2). Finally, the models were checked for reliability by Y scrambling to verify the absence of a chance correlation.¹¹

UK. The data analysis and multivariate calibrations were carried out with the Unscrambler software.⁵³ Partial least-squares (PLS) regression was used for data analysis and modeling. PLS regression is based on a linear transformation of the original descriptors to a limited number of orthogonal factors, attempting to maximize the covariance between the descriptors and the response variable. The term “latent variable” is used to denote the PLS factors, since they can be interpreted as describing the inherent chemical properties. Multivariate calibration with PLS is reviewed by Martens and Næs⁵⁴ and Wold et al.⁵⁵ Nonsignificant descriptor variables were assigned zero weight; these variables were identified using a jackknife method for significance testing of the PLS model parameters during cross-validation.

All descriptor variables were preprocessed by autoscaling to zero mean and unit variance. Cross-validation was used to establish the rank of the calibration model (number of latent variables), and an external validation set was used to estimate the prediction error. The calibration model was characterized by the standard deviations of the prediction residuals for the calibration objects and the external validation sets respectively: root-mean-square error of calibration and root-mean-square error of prediction. The explained variances are defined as sums of squares due to regression divided by sums of squares about the mean: R^2 (square of the multiple correlation coefficients for the calibration objects) and Q^2 (square of the multiple correlation coefficients for the external test set).

VCCLAB. Associative neural network (ASNN) represents a combination of an ensemble of feed-forward neural networks and k NN. This method uses the correlation between ensemble responses (each molecule is represented in a space of neural network models as a vector of model predictions) as a measure of distance amid the analyzed cases for the nearest-neighbor technique. Thus, ASNN performs k NN in a space of ensemble residuals. This provides an improved prediction by the bias correction of the neural network ensemble.^{56,57} The neural networks ensemble of 50 networks with one hidden layer was used. After several preliminary runs, we fixed three hidden neurons for all data sets. The efficient partition algorithm was used to train the neural network ensemble.⁵⁸ The calculations were performed using the program available at <http://www.vcclab.org/lab/asnn>. The leave-one-out cross validation correlation coefficients Q_{abs}^2 calculated for neural networks as described elsewhere⁵⁹ were reported as the model accuracy for the training set.

UBC. The applicability of various statistical and machine-learning approaches for creating QSAR models was explored including MLR, ANN, PLS, and SVM. In all calculations, the “inductive” QSAR descriptors were used as independent variables and experimental log IG_{50} parameters as dependent properties. The Weka software (version 3.5.6)⁶⁰ was used; it includes the following modules: Linear Regression for MLR, MultilayerPerceptron for ANN, PLSClassifier for PLS, and SVMreg for SVM. Similar data-mining approaches have been used by other collaborating groups; additional details

about the implementation of these approaches in the Weka modules are given elsewhere.⁶⁰

All four types of QSAR models (MLR-, ANN-, PLS-, and SVM-based) have been investigated using a 90%/10% division of 644 compounds of the modeling set as well as by using the LOO cross-validation. We have used default settings for MLR and PLS modules. For ANN, we reset settings with 10 hidden nodes for only one hidden layer, weight decay, learning rate = 0.8, and momentum = 0.1. For SVM, several models with different types of kernels (linear kernel; polynomial kernel of degrees 1, 2, and 3) and values of complexity parameter C (1, 10, 50, 100, and 1000) have been built for the training set. The best model, that is, the one that results in the highest Q_{abs}^2 for the training set, was chosen for further analyses using two external validation sets. The results reported for all SVM models in this paper were obtained by setting the polynomial kernel degree to 2 and $C = 1$. The performance of all QSAR models was assessed by the standard statistical properties that included Pearson squared correlation coefficient r^2 , SDE, and coverage.

In addition, both LOO and 10-fold cross-validation analyses were conducted. The model performance was assessed by cross-validation parameters Q_{abs}^2 computed for LOO and the $Q^2(\text{c10})$ parameter for 10-fold cross-validation.

2.4. Model Applicability Domains and Chemical Space Coverage. Defining model applicability domains is an active area of modern QSAR research.^{61,62} Every QSAR model can formally predict the relevant target property for any compound for which chemical descriptors can be calculated. However, since each model is developed using compounds in the training set only (that cover only a small fraction of the entire chemistry (i.e., descriptor) space), the special applicability domain for each model should always be defined. This restriction prevents making predictions for compounds that differ substantially from those in the modeling set. Generally, there is no universal method of defining the AD in the descriptor space, especially when using variable selection techniques. Global applicability domains are defined in the complete chemistry space, that is, using all descriptors, whereas local domains are defined in the context of specific variable selection models using only selected (optimized) variables. Each participating group (with one exception) has adopted its own definition of the AD in the context of the respective QSAR methods. Another closely related parameter is chemistry space coverage. Thus, as a consequence of defining the AD, only a certain fraction of compounds in any external data set is expected to fall within such a domain. This fraction is therefore referred to as the data set coverage. The definitions of AD used by each group are described below.

UNC. The AD is calculated from the distribution of similarities between each compound and its k nearest neighbors in the training set (similarities are computed as Euclidean distances between compounds represented by their multiple chemical descriptors). Based on the previous studies, the standard cutoff value to define the applicability domain for a QSAR model places its boundary at one-half of the standard deviation calculated for the distribution of distances between each compound in the training set and its k nearest neighbors in the same set (assuming a Boltzmann-like distribution of these distances). Thus, if the distance of the test compound from any of its k nearest neighbors in the

training set exceeds the threshold, the prediction is considered unreliable. The detailed description of the algorithm to define the AD is given elsewhere.^{11,45}

ULP. Applicability domains in ISIDA-MLR and ISIDA-*k*NN were calculated with an approach similar to that described by UNC above. Additionally, the "Bounding Box" AD has been used for ISIDA-MLR calculations. Thus, for each fragment descriptor involved in the model, its minimal and maximal occurrences within compounds in the training set were retrieved and defined as an allowed range for this fragment. For a given validation set compound, the model was considered unreliable for the prediction if the occurrence of one of its fragment descriptors was outside the corresponding range defined for the training set. For ISIDA-SVM and CODESSA-MLR calculations, the ADs of selected models were not defined.

UI. Hat values from the leverage matrix, representing the compound "distance" from the model structural space, were used to check structurally influential chemicals—*X* outliers (with high leverage values: $h > h^*$, the critical value being $h^* = 3p'/n$, where p' is the number of model variables plus one and n the number of the objects used to calculate the model).⁶³ Moreover, the presence of outliers for the response (*Y* outliers) was also verified, and such problematic compounds in the modeling set were identified as those with standardized residuals greater than 2.5.

UK. A PLS calibration model can determine the valid domain for the descriptor variables. New validation and prediction objects are assessed by comparing the residual standard deviation (the Euclidean distance to the PLS model) and the leverage (the Mahalanobis distance within the PLS model space) to that of the calibration objects. These two distance measures were used to decide whether or not a new object was within the AD of the training set model. Here, the 5% significance level was chosen as the limit for the residual standard deviation, and the limit for the leverage was set to 3 times the average leverage for the calibration objects. The leverage is directly proportional to Hotelling's *T*² diagnostic (a multivariate generalization of the standard *t* test).⁶⁴

VCLLAB. The ensemble of $N = 100$ models was used to calculate the ultimate training set ASNN model. Thus, for any molecule, a vector with 100 predictions is always calculated. This vector corresponds to a new representation of a molecule in so-called model space. For each analyzed molecule from the test set, we determined a molecule in the training set that had a maximum correlation with the analyzed molecule in the model space.⁶² A cut-off value of $r^2 = 0.7$ was used to define the AD of the ASNN model. Thus, if the analyzed molecule had $r^2 > 0.7$ at least to one of training set molecules in the space of models, it was considered inside of the AD of the model. Otherwise, it was considered outside of the AD of the model.

UBC. The range of a descriptor values is defined as an interval $[0.85\text{MIN} - 1.15\text{MAX}]$ where MIN and MAX are the minimum and maximum values appearing in the training set for a given descriptor (i.e., 15% deviation from the range of descriptor values present in the training set was allowed). The test set compound is considered to fit the AD if all its descriptor values are within the described range. For the case of the studied data sets, only one entry in validation set I did not fit the AD; all compounds in the second validation

test have been covered by the AD. It was found that exclusion of that single AD outlier did not change the prediction statistics, and therefore, the AD described above afforded 100% coverage of both validation sets.

3. RESULTS

The statistical parameters of the predictions obtained from all QSAR models for the modeling set and the two external validation sets are shown in Table 2. The results indicate that most of the models were successful in reproducing the experimental data for the 644-compound modeling set. A total of 9 out of all 15 models afforded a Q_{abs}^2 higher than 0.80, and only one model had MAE greater than 0.4 for this self-validation test.

It is of interest to notice that on average the results for validation set II were not as good as those for validation set I for almost all models. The most likely reason for this observation is the greater general dissimilarity of the compounds in validation set II to the compounds in the modeling set. This conclusion can be illustrated by considering the model AD as implemented by the UNC group. About 50% of the compounds in validation set II were identified as outside the AD, which was calculated using all descriptors. In contrast, for validation set I, only ca. 20% of the compounds were found to be outside the AD.

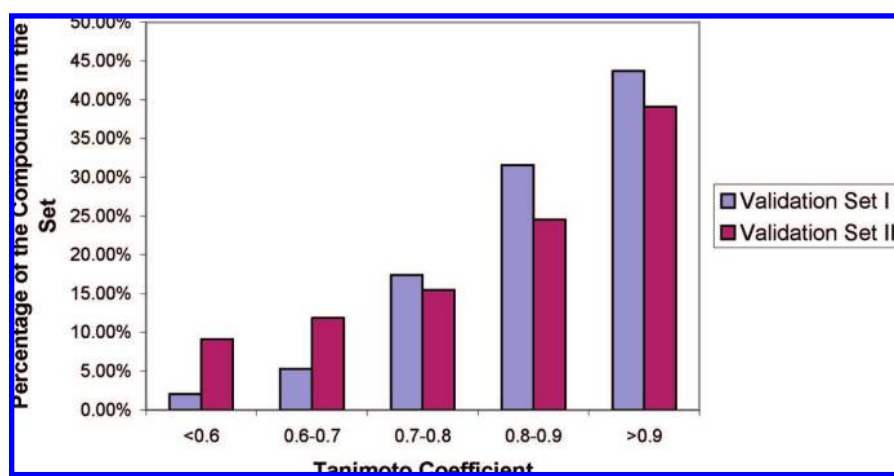
To investigate the level of (dis)similarity between the modeling set and two validation sets, the MACCS structural keys⁶⁵ were calculated for all compounds in the three sets using the MOE software.⁴⁰ The Tanimoto similarity coefficients⁶⁶ between all pairs of compounds in the two validation sets versus the modeling set were then calculated. Figure 2 shows the distribution of Tanimoto similarity between compounds in both validation sets versus that in the modeling set. Clearly, the compounds in validation set I are more "similar" to the compounds in the modeling set than the compounds in validation set II. This result provides a clear demonstration that even a validation set obtained from the same data pool as the modeling set may not serve as a real substitute for a truly external validation set. For this reason, using a totally independent data set (such as validation set II) could lead to more realistic estimates of the true external power of QSAR models. We shall now describe the results from individual groups.

UNC. Both Dragon and MolconnZ descriptors were used for *k*NN and SVM QSAR modeling. The $0.75/0.75 Q_{\text{abs}}^2/R_{\text{abs}}^2$ cutoff for the training and test sets, respectively, generated from the modeling sets was used to select the acceptable models. The total numbers of models that satisfied these cutoff criteria were 542, 192, 60, and 114 for *k*NN-Dragon, *k*NN-MolconnZ, SVM-Dragon, and SVM-MolconnZ QSAR, respectively. The toxicity for each validation set compound was predicted by averaging the predictions obtained with all training set models that had this compound within their respective AD. As mentioned above, all compounds in the validation sets that were out of the global AD (i.e., defined using all descriptors) were excluded. Because the global AD was used, the AD of the resulting SVM models was the same as that of the *k*NN models. Therefore, the coverage of the test sets obtained from these two approaches was identical.

Table 2. Statistical Results Obtained with All QSAR Models for the Modeling and External Validation Sets

model	group ID	modeling set ($n = 644$)			validation set I ($n = 339$)			validation set II ($n = 110$)		
		Q_{abs}^2	MAE	coverage (%)	R_{abs}^2	MAE	coverage (%)	R_{abs}^2	MAE	coverage (%)
kNN-Dragon	UNC	0.92	0.22	100	0.85	0.27	80.2	0.72	0.33	52.7
kNN-MolconnZ	UNC	0.91	0.23	99.8	0.84	0.30	84.3	0.44	0.39	53.6
SVM-Dragon	UNC	0.93	0.21	100	0.81	0.31	80.2	0.83	0.27	52.7
SVM-MolconnZ	UNC	0.89	0.25	100	0.83	0.30	84.3	0.55	0.37	53.6
ISIDA-kNN	ULP	0.77	0.37	100	0.73	0.36	78.5	0.63	0.37	42.7
ISIDA-SVM	ULP	0.95	0.15	100	0.76	0.32	100	0.38	0.50	100
ISIDA-MLR	ULP	0.94	0.20	100	0.81	0.31	95.9	0.65	0.41	51.8
CODESSA-MLR	ULP	0.72	0.42	100	0.71	0.44	100	0.58	0.47	100
OLS	UI	0.86	0.30	92.1	0.77	0.35	97.0	0.59	0.43	98.2
PLS	UK	0.88	0.28	97.7	0.81	0.34	96.1	0.59	0.40	95.5
ASNN	VCCLAB	0.83	0.31	83.9	0.87	0.28	87.4	0.75	0.32	71.8
PLS-IND_I	UBC	0.76	0.39	100	0.74	0.39	99.7	0.45	0.54	100
MLR-IND_I	UBC	0.77	0.39	100	0.75	0.40	99.7	0.46	0.53	100
ANN-IND_I	UBC	0.77	0.39	100	0.76	0.39	99.7	0.46	0.53	100
SVM-IND_I	UBC	0.79	0.31	100	0.79	0.35	99.7	0.53	0.46	100
consensus model I^a		0.92	0.23	100	0.85	0.29	100	0.67	0.39	100
consensus model II^b		0.92	0.22	100	0.87	0.27	100	0.70	0.34	100
consensus model IIB^c		0.92	0.22	100	0.87	0.27	100	0.70	0.36	100
consensus model III^d		0.92	0.22	100	0.86	0.28	99.7	0.70	0.34	98.2

^a Consensus model I: average of the 15 selected models without considering their individual applicability domains. ^b Consensus model II: average of the nine models (kNN-Dragon, kNN-MolconnZ, SVM-Dragon, SVM-MolconnZ, ISIDA-kNN, ISIDA-MLR, OLS, PLS, and ASNN) **using** their individual applicability domains. ^c Consensus model IIB: average of the nine models (kNN-Dragon, kNN-MolconnZ, SVM-Dragon, SVM-MolconnZ, ISIDA-kNN, ISIDA-MLR, OLS, PLS, and ASNN) **without using** their individual applicability domains. ^d Consensus model III: average of predictions with a minimal number of one model among the nine having an individual AD, excluding predictions according to the Grubbs's statistics.

**Figure 2.** The comparison between the similarities of the modeling set vs validation set I and the modeling set vs validation set II.

Validation Set I. The accuracies of prediction for this external validation set were lower than those for the modeling set but still very high ($R_{\text{abs}}^2 > 0.8$, $\text{MAE} < 0.32$) for all four models. Since the statistical parameters of all four models were similar (Table 2), it is difficult to identify the best-performing model in terms of combination of kNN or SVM with either MolConnZ or Dragon descriptors.

Validation Set II. The R_{abs}^2 of validation set II ranged from 0.44 to 0.83, and the corresponding MAE ranged from 0.27 to 0.39, with kNN-MolConnZ and SVM-MolConnZ models having relatively lower predictive power (Table 2). Therefore, the MolConnZ descriptors proved to be less successful for modeling of this external set. On the other hand, the mean atomic polarizability and Moriguchi octanol–water partition coefficient Dragon descriptors were selected as the top two most significant descriptors for the final kNN-Dragon model. Thus, the underperformance of MolConnZ descriptors may be because such property descriptors are not included in the

software. In previous reports, property descriptors such as lipophilicity were also established as critical descriptors for aquatic toxicity models.^{67,68}

ULP. At the training stage, the best ISIDA-MLR model used 109 fragment descriptors. Only 26 fragments were selected by the ISIDA-kNN variable selection procedure. The CODESSA-MLR model involves six molecular descriptors: *average atom weight*, *molecular surface area*, *FPSA2* and *FPSA-3 fractional positive surface areas*, *WNSA1 weighted negative surface area*, and the *relative number of S atoms*. Both ISIDA and CODESSA-Pro models led to reasonable statistical parameters, $Q_{\text{abs}}^2 = 0.72\text{--}0.95$ and $\text{MAE} = 0.15\text{--}0.42$, with a complete coverage of the training set. It should be noted that the linear ISIDA-MLR and nonlinear ISIDA-SVM models gave similar results for this set. Then, the selected models were applied to the two external validation sets.

Table 3. Statistical Results Obtained with All QSAR Models for External Validation Sets with Full Coverage (100% - No AD)

model	group ID	validation set I (n = 339)		validation set II (n = 110)	
		R_{abs}^2	MAE	R_{abs}^2	MAE
kNN-Dragon	UNC	0.84	0.29	0.59	0.43
kNN-MolconnZ	UNC	0.83	0.31	0.49	0.49
SVM-Dragon	UNC	0.70	0.37	0.53	0.42
SVM-MolconnZ	UNC	0.77	0.33	0.58	0.44
ISIDA-kNN	ULP	0.71	0.39	0.37	0.54
ISIDA-SVM	ULP	0.76	0.32	0.38	0.50
ISIDA-MLR	ULP	0.49	0.38	0.43	0.49
		0.71 ^a	0.35 ^a		
CODESSA-MLR	ULP	0.71	0.44	0.58	0.47
OLS	UI	0.77	0.36	0.59	0.42
PLS	UK	0.81	0.34	0.59	0.41
ASNN	VCCLAB	0.85	0.30	0.66	0.38
PLS-IND_I	UBC	0.74	0.39	0.45	0.54
MLR-IND_I	UBC	0.75	0.40	0.46	0.53
ANN-IND_I	UBC	0.76	0.39	0.46	0.53
SVM-IND_I	UBC	0.79	0.35	0.53	0.46
consensus model I^b		0.85	0.29	0.67	0.39

^a Without one outlier (see text). ^b Consensus model without considering the applicability domain.

Validation Set I. All ISIDA models led to reasonable predictions: $R_{\text{abs}}^2 = 0.73\text{--}0.81$ and $\text{MAE} = 0.31\text{--}0.36$, with a coverage ranging from 78.5 to 100% (see Table 2). Compared to other methods, CODESSA-MLR calculations displayed a good correspondence between predicted and experimental pIGC_{50} with $R_{\text{abs}}^2 = 0.71$, but rather large prediction error ($\text{MAE} = 0.44$).

The performance of ISIDA-kNN, ISIDA-SVM, and CODESSA-MLR models was also reasonable ($R_{\text{abs}}^2 = 0.71\text{--}0.76$ and $\text{MAE} = 0.32\text{--}0.44$) even when the applicability domain was not applied (Table 3). Poor statistical parameters ($R_{\text{abs}}^2 = 0.49$ and $\text{MAE} = 0.38$) of the ISIDA-MLR model could be related to only one outlier, 2,2,2-tribromoethanol, for which experimental ($\text{pIGC}_{50} = 0.11$) and predicted (-9.03) values were very different. Without this outlier, the ISIDA-MLR model was much better: $R_{\text{abs}}^2 = 0.71$ and $\text{MAE} = 0.35$. In fact, 2,2,2-tribromoethanol contains three Br-C-Br and C-Br fragments which are very poorly represented in the training set. Thus, the observed outlying value could be explained by bad statistics related to the aforementioned fragments. It should be also noted that this particular compound was found as an outlier with SVM/Dragon, SVM/MZ, Codessa Pro, and PLS, and it falls outside the AD of six models including ISIDA/MLR.

Validation Set II. If the applicability domain was not applied, the performance of all ISIDA models was relatively poor: $R_{\text{abs}}^2 < 0.5$ and $\text{MAE} = 0.49\text{--}0.54$ (see Table 3). The CODESSA-MLR calculations led to somewhat better results ($R_{\text{abs}}^2 = 0.58$ and $\text{MAE} = 0.47$).

When AD mode was indeed activated (see Table 2), ISIDA-kNN and ISIDA-MLR models afforded fairly reasonable values of $R_{\text{abs}}^2 = 0.63\text{--}0.65$ and $\text{MAE} = 0.37\text{--}0.41$; however, this improvement was also associated with a relatively low coverage of 42.7–51.8% of the data set.

UI. After several attempts to model all 644 chemicals using OLS regression, 26 compounds were found to be out of the global AD for a collection of different models generated by GA selection. These compounds that were strongly affecting the performance of models using different molecular descriptors for the complete training set were

excluded as outliers; thus, the final modeling data set consisted of 618 chemicals.

The best predictive model, based on six variables, was selected from a population of 100 models of different descriptor typology (where the number of variables used in the models varied between 1 and 6 as described in the Methods). When considering the population of the 80 best six-dimensional models, the range of Q_{loo}^2 was from 0.82 to 0.84. The best model ($\text{MAE} = 0.30$) was finally chosen from those included in the population according to the QUIK rule ($\Delta K_{\text{xy}} = 7\%$) and also evaluated for its robustness ($Q_{\text{boot}}^2 = 0.83$, $R_{\text{Y-scrambling}}^2 = 0.01$).

The variables included in this model, in order of importance as defined in the model by their standardized coefficients, are AMR (Ghose–Crippen molar refractivity), Me (mean atomic Sanderson electronegativity), nHAcc (number of H bonds atoms acceptors), O-056 (fragment: alcohol), H-046 (H attached on C (sp³) without heteroatoms on the adjacent C), and O-058 (fragment: =O). It is important to note that the descriptor AlogP was selected as an important variable among the population of the GA-developed models. This observation (even though this variable was not included in the proposed best OLS model because it was substituted by other descriptors) highlights the well-known importance of lipophilicity in modeling fish aquatic toxicity.^{67,68}

The evaluation of the AD of the proposed OLS model on the training set of 618 compounds revealed the presence of 18 compounds out of the X-structural domain and 13 Y outliers out of the response domain (domain coverage 92.1%).

Validation Set I. The parameters of the external predictive power of the best OLS model as applied to validation set I were high ($R_{\text{abs}}^2 = 0.77$, $\text{MAE} = 0.35$) and comparable to those obtained on the training set. The model was found to cover 97% of the domain for validation set I (10 structural outliers). The exclusion from validation set I of these compounds did not give any significant increase in the model performance (R_{abs}^2 remained at 0.77 irrespective of whether compounds out of the AD were included or not).

Validation Set II. The performance of the OLS model as applied to validation set II was relatively low ($R_{\text{abs}}^2 = 0.59$, $\text{MAE} = 0.43$). The model was found to cover about 98% of the domain for validation set II (two structural outliers). Also in this case, the exclusion from the validation set of these compounds did not give any significant increase in the model performance (R^2 was 0.59 irrespective of the use of AD). Apparently, the number of compounds outside the structural AD of the two validation sets was too small (10 in validation set I and only two in validation set II) to perturb the prediction accuracy in any significant way.

UK. A total of 31 outliers were identified and excluded from the modeling set on the basis of an examination of projections of PLS factors versus the response variable. In the final PLS calibration model, 515 descriptor variables were selected for inclusion on the basis of significance tests using jackknifing for the calibration set in preliminary runs. The number of latent variables to retain in the PLS model was estimated at five using cross-validation with 20 randomly assigned validation segments of equal size. These five latent variables capture 64.4% of the variance in the descriptor variables, thus demonstrating that the information contained in the descriptors is effectively used in the calibration model. The explained calibration variance (r_{Cal}^2) for the dependent variable (the logarithm of the 50% growth inhibitory concentration) was 87.9%, and the root-mean-square error of calibration was 0.36 log units. The explained prediction variance for the cross-validation compounds (q_{CV}^2) was 85.7%, and the root-mean-square error of prediction was 0.39 log units.

The PLS model also defines a valid domain for the descriptor variables. A total of 23 compounds in the external validation set were substantially different from the calibration compounds and fell outside the 5% confidence bound for the residuals and the leverage limit of 0.03. These compounds were thus excluded from further use, and the model was subsequently validated with the remaining 426 compounds from the external validation set, where the explained variance (q_{Ext}^2) for the dependent variable was 78.8% and the root-mean-square error of prediction was 0.47 log units. If all compounds in the external test set had been retained, the explained variance and root-mean-square error of prediction would still remain almost the same (79.1% and 0.48 log units, respectively).

Validation Set I. After removal of 18 compounds outside the model AD, the model was validated with the 321 remaining compounds. The explained variance (q_{Ext}^2) for the dependent variable was 81.5%, the root-mean-square error of prediction was 0.44, and the mean absolute error was 0.34 log units. If all objects in the external test set had been retained, the explained variance, the root-mean-square error of prediction, and the mean absolute error would still remain almost the same (81.8%, 0.45, and 0.34 log units, respectively).

Validation Set II. After removal of five compounds outside the model AD, the model was validated with the 105 remaining compounds. The explained variance (q_{Ext}^2) for the dependent variable was 60.5%, the root-mean-square error of prediction was 0.55, and the mean absolute error was 0.41 log units. If all compounds in the external test set had been retained, the explained variance, the root-mean-square error

of prediction, and the mean absolute error would also remain almost the same (62.7%, 0.57, and 0.41 log units, respectively).

VCCLAB. All 644 molecules from the training set were used to build the ASNN model. The model involved a total of 58 descriptors, and it calculated leave-one-out $\text{MAE} = 0.29$ for the training set. To better validate the prediction ability of the method, we also applied a 5-fold cross-validation procedure. The calculated results, $\text{MAE} = 0.32$, were similar to those calculated using all molecules, thus indicating stability of the model. The cross-validation studies identified a set of 30 outlying molecules (or 4.6%) which had $\text{MAE} > 1$ log unit. The statistical results for the validation sets are reported in Table 3. In order to better explore the data, we also applied several other machine learning methods using the same settings and scripts developed in our previous study.⁶⁹ SVM with a radial basic function kernel calculated similar performance with $\text{MAE} = 0.31$ and 0.38 for both validation sets. The singular value decomposition calculated a lower prediction ability, $\text{MAE} = 0.35$ and 0.39, while the $k\text{NN}$ method failed to model these data with $\text{MAE} = 0.38$ and 0.63. Notice that all these methods were used with all descriptors and default settings, while some methods may require description selection. Thus, these results should be only considered as an exploratory analysis of data. We also applied variable selection pruning methods⁷⁰ to detect the set of most important descriptors. The minimal set of 14 descriptors ($\text{MAE} = 0.31$ and 0.40 for the validation sets) included molecular weight, the number of non-hydrogen atoms, and the numbers of donors and acceptors. The first two descriptors represent the bulk effect and correlate with the lipophilicity of molecules, which is one of the most important descriptors in the models of Schultz and colleagues.^{12,14–18} The other two descriptors are directly related to electrophilic properties, which is also an important parameter in these models. It is interesting that this set also included three types of E-state indices, SsOH(alc) , SsOH(phen) , and SsOH(acid) , corresponding to the oxygen atom in the hydroxy group at different binding environments. It presumably allowed further quantification of both lipophilic and electrophilic properties of the molecules. Indeed, these indices were proposed as an extension of the basic set of E-state indices^{31,32} for the ALOGP model.^{33,34}

We should notice that the model based on a minimal set of 14 descriptors had similar performance to the model built with all descriptors. However, in the final report, we decided to use the model that was built with all descriptors because this model would be less sensitive to the missing descriptors problem for future prediction of chemical scaffolds that were not covered by the training set. We found that the E-state indices that provided a complete representation of a molecule were not redundant or duplicative.

Validation Set I. The results for this validation set, $\text{MAE} = 0.30$, were in good agreement with the accuracy of the model for the training set. This set had 10 outlying molecules (3%) identified according to $\text{MAE} > 1$ log unit criteria.

Validation Set II. The prediction accuracy for this set, $\text{MAE} = 0.38$, was lower compared to that for the training and the first validation sets. This set had six outlying molecules, thus contributing the highest percent, 5.5%, of these molecules amid all sets.

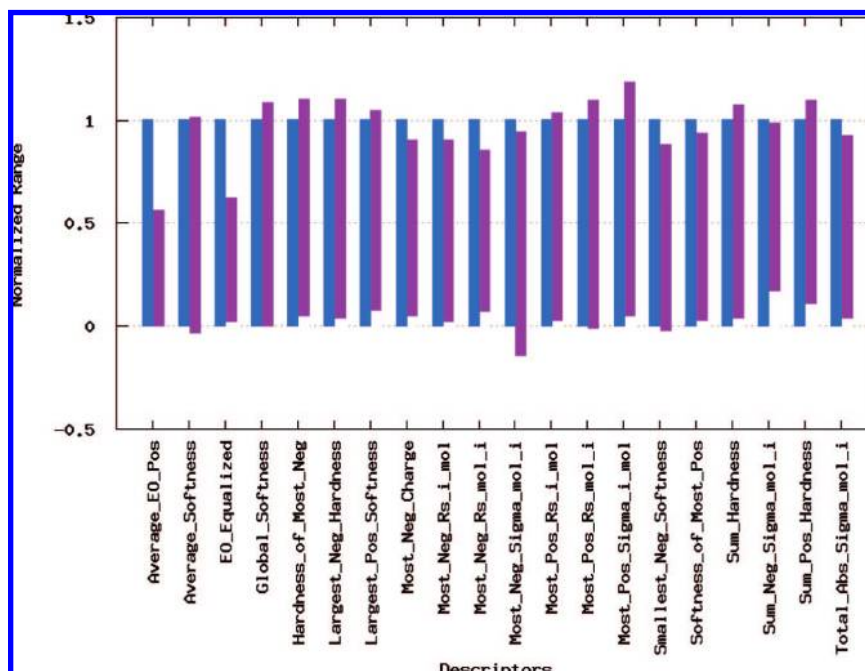


Figure 3. Normalized inductive descriptor ranges of the training set (blue) compared to the validation set I (purple).

A total of 18 outlying molecules (nine for each validation set) were outside the AD of the model. There was about the same percentage of molecules, 84% and 87%, within the AD for the training and first validation sets, respectively (Table 2). The accuracy of prediction for molecules outside the AD was lower as measured by MAE = 0.40 and 0.45 for the training and first validation sets, respectively. A lower percent of molecules, 72%, was within the domain for the second validation set (Table 2). The prediction accuracy for molecules outside the domain for this set was the lowest, as in terms of MAE = 0.52.

UBC. To eliminate possible cross-correlation among 50 descriptors, the AutoQSAR SVL script was used, which is based on the PLS method. This procedure resulted in the selection of 20 descriptors. This approach enables automated QSAR modeling “on the fly” and is available through the SVL exchange.

All of these 20 selected descriptors were used in all four models (MLR-, PLS-, ANN-, and SVM-based). The resulting models allowed very accurate training with the 10-fold cross-validation parameters within a narrow range of 0.76–0.79. Similar to other reported approaches, the predicted toxicity values for each validation set compound were calculated by averaging the predictions from all 10 training set models resulting from the 10-fold cross validation analysis.

Validation Set I. The ranges of descriptors for the training set were calculated, and only 1 out of 339 compounds in the first validating set was found to be outside the AD. Excluding the outlier compound, neither the R_{abs}^2 nor the MAE values changed, which overall were on par with the training set results.

Validation Set II. The statistical parameters for predicting the second validation set were significantly worse than those for the first validation set. For instance, the highest value of R_{abs}^2 was 0.53 for validation set II compared to the highest R_{abs}^2 of 0.79 for the first validation set.

To investigate possible reasons for the difference in the observed prediction accuracies for sets I and II, we have considered the ranges and distributions of descriptor values in the training and validation sets. For each descriptor type used in the modeling, we have normalized all of its values using the minimum and maximum values of each descriptor in the training set for range scaling. Consequently, the descriptor values have been transformed to be within the range of [0,1] for the training set. Most of the descriptor values for validation sets I and II were found within the ranges of corresponding descriptors for the training set, but several were outside these ranges. Figures 3 and 4 show the histograms for range-scaled descriptor value distribution in the training set versus that in validation sets I and II, where the lengths of the histograms correspond to the normalized ranges of the descriptor values and their mutual positioning is defined by medians of descriptors values within the sets (the histogram centers have been placed at the median values). Thus, the extent of histogram overlap is not only determined by the range of descriptors values but also by their distribution within the range.

As can be seen from the graphs, neither set I nor set II contained extreme values of descriptors that would significantly extend beyond the training set ranges (as our AD analysis already illustrated). At the same time, in the case of set II, the distribution of descriptors values for the training set is clearly much more unbalanced versus the distribution for the test set. This, perhaps, may be considered as one possible reason for the less accurate prediction of toxicity values for compounds from set II. We also note that for both external validation sets the space coverage was 100%. This observation may imply that the AD used in these studies may be too generous, especially as applied to validation set II. We suggest that optimal prediction by the QSAR model can be achieved in those cases, when the values of descriptors in training and external sets have both similar ranges and similar distributions.

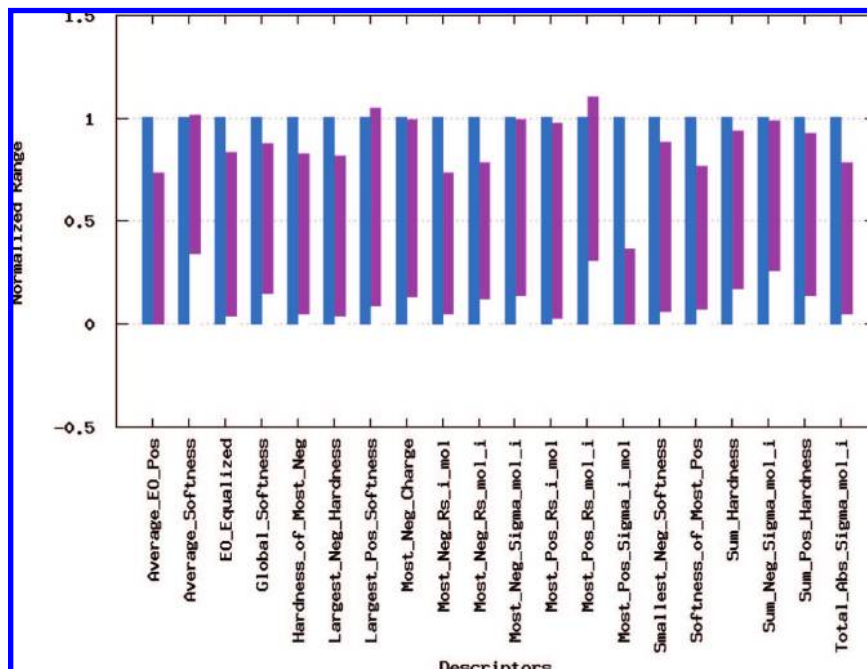


Figure 4. Normalized inductive descriptor ranges of the training set (blue) compared to the validation set II (purple).

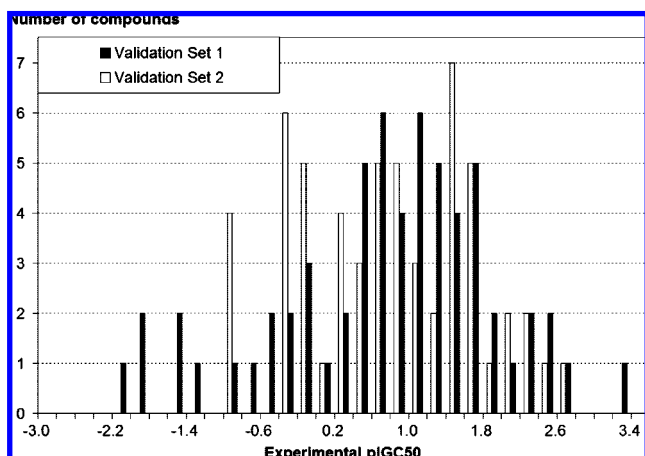


Figure 5. The pIC_{50} toxicity distribution of compounds that have been marked as outside of the applicability domain of at least three models.

4. DISCUSSION

Comparison between Methods and Models. The objective of this prospective methodological study was to explore the suitability of different QSAR modeling tools for the analysis of a data set with an important toxicological end point. Typically, such data sets are analyzed with one (or several) modeling technique, with a great emphasis on the (high value of) statistical parameters of the training set models. Such an approach is exemplified by the studies of Schultz and co-workers, who generated the experimental data used in our analysis.^{12,14–18} In a series of publications that included both experimental results and QSAR models based on those results, the authors typically used one modeling method (e.g., linear regression analysis) and reported the single best model in each individual publication for the respective data sets. The largest data set used in earlier publications by Schultz and co-workers included only 467 compounds.

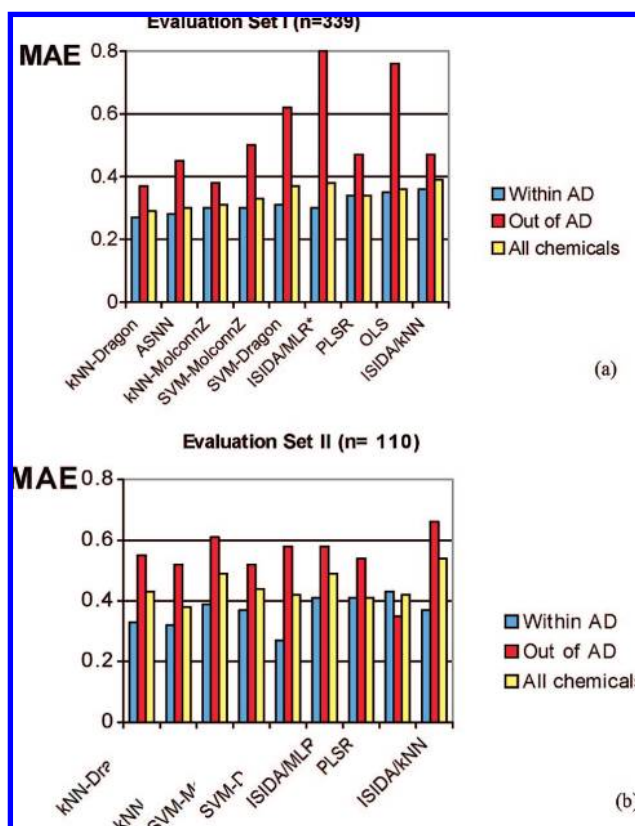


Figure 6. The MAEs of nine QSAR models for the first validation set (a) and the second validation set (b). *The histogram corresponding to compounds outside of AD for the ISIDA/MLP method (MAE = 1.19) was truncated.

In this paper, we went well beyond the modeling studies reported in the original publications^{12,14–18} in several respects. First, we have compiled all reported data on chemical toxicity against *T. pyriformis* in a single large data set and attempted to develop global QSAR models for the entire set. Second, we have employed multiple QSAR modeling techniques thanks to the engagement of six

Table 4. Statistical Parameters Obtained with Consensus Model III for Training and Validation Sets: Coverage and Accuracy of Consensus Model III vs Its Minimal Number of Incorporated Single Models (Among the Nine Models That Have Implemented AD)

minimal number of models	consensus model III								
	training set ($n = 644$)			validation set I ($n = 339$)			validation set II ($n = 110$)		
	R_{abs}^2	MAE	coverage (%)	R_{abs}^2	MAE	coverage (%)	R_{abs}^2	MAE	coverage (%)
1^a	0.92	0.22	100	0.86	0.28	99.7	0.70	0.34	98.2
2	0.92	0.22	100	0.86	0.28	99.1	0.68	0.35	95.5
3	0.92	0.22	100	0.87	0.27	97.9	0.69	0.32	87.3
4^b	0.92	0.22	100	0.86	0.27	96.2	0.69	0.32	70.9
5	0.92	0.22	100	0.87	0.26	90.6	0.76	0.29	61.8
6	0.92	0.22	99.8	0.87	0.26	88.5	0.77	0.29	55.5
7	0.92	0.22	99.5	0.87	0.26	81.7	0.78	0.29	48.2
8	0.93	0.21	96.4	0.87	0.25	70.8	0.81	0.26	30.9
9	0.94	0.20	77.6	0.88	0.23	56.6	0.77	0.29	20.9

^a Consensus model III involving a minimal number of five models has the best balance between accuracy (reasonable R_{abs}^2 and MAE) and coverage. ^b Example: if the minimal number of models is equal to four, it means that the toxicity of a given compound is predicted only if it is found as inside the AD of at least four models (among the nine ones having an AD: *k*NN-Dragon, *k*NN-MolconnZ, SVM-Dragon, SVM-MolconnZ, ISIDA-*k*NN, ISIDA-MLR, OLS, PLS, and ASNN).

collaborating groups. Third, we have focused on defining model performance criteria not only using training set data but most importantly using external validation sets that were not used in model development in any way (unlike any common cross-validation procedure).⁷¹ This focus afforded us the opportunity to evaluate and compare all models using simple and objective universal criteria of external predictive accuracy, which in our opinion is the most important single figure of merit for a QSAR model that is of practical significance for experimental toxicologists. Fourth, we have explored the significance of applicability domains and the power of consensus modeling in maximizing the accuracy of external predictivity of our models.

We believe that results of our analysis lend strong support for our strategy. Indeed, all models performed quite well for the training set (Table 2) with even the lowest Q_{abs}^2 among them as high as 0.72. However, there was much greater variation between these models when looking at their (universal and objective) performance criteria as applied to validation sets I and II, both with (Table 2) and without (Table 3) the applicability domain.

It is of a particular interest to explore and compare the performances of all models without the applicability domain (see Table 3) since, in this case, the comparison can be made for both validation sets including all compounds (full coverage). For validation set I, all models demonstrated similar performance with an average MAE of 0.36 ± 0.04 . Fisher's test indicates that there is no significant statistical difference between MAE values equal to 0.29 and 0.33 (that means, for example, that the results generated with *k*NN-Dragon, *k*NN-MolConnZ, ASNN, and ISIDA-SVM are equivalent for the first validation set). For the second validation set, the MAE average for all models is 0.48 ± 0.06 (significantly higher compared to the first validation set). The ASNN method afforded the lowest MAE of 0.38, significantly lower than 0.48, according to the Fisher's test.

The activity distribution for compounds, which were found to lie outside the AD by at least three individual models, is shown in Figure 5. Apparently, this distribution is similar to that of all compounds (cf. Figure 1); that is, there are similar fractions of low, intermediate, and highly toxic compounds irrespective of whether they are found within or

without the applicability domains. This result indicates that, in modeling complex end points such as aquatic toxicity, when multiple mechanisms of action could be involved, there is no simple relationship between compounds' chemical similarity and their end point toxicity.

Role of the Applicability Domain for Individual Models. Of 15 QSAR approaches used in this paper, nine implemented method-specific applicability domains. Models that did not define the AD showed a reduced predictive accuracy for validation set II even though they yielded reasonable results for validation set I. Only CODESSA-MLR (which did not employ any AD) approached in accuracy the lower bound of the models using the AD as measured by $R_{\text{abs}}^2 = 0.58$ but still had one of the highest MAEs of 0.47 (Table 2). On the other hand, among models employing the AD, only *k*NN-MolconnZ had a relatively low accuracy of prediction for validation set II, with R_{abs}^2 below 0.5. For all other models, R_{abs}^2 ranged between 0.55 and 0.83.

On average, the use of applicability domains improved the performance of individual models, although the improvement came at the expense of lower chemistry space coverage (cf. Tables 2 and 3). The direct comparison between individual models appears difficult due to different definitions of AD and different interplay between coverage and accuracy for different models.

The choice of descriptors played a more important role than the choice of modeling techniques. This observation could only be made in a few cases when different approaches utilized exactly the same descriptor sets. For instance, the results of UNC studies (Table 2) clearly indicated that Dragon descriptors afforded significantly better models, both with SVM and *k*NN, than MolconnZ descriptors. Dragon and MolconnZ share many descriptors, but the most significant difference between the two methods is that Dragon has additional physical chemical descriptors that apparently play an important role in defining aquatic toxicity. Similarly, ISIDA *k*NN and ISIDA MLR afforded relatively similar results when applicability domains were used. Finally, the last four individual models reported in Table 2 also produced similar results; that is, changing modeling techniques could not help increase the model accuracy in the absence of the AD.

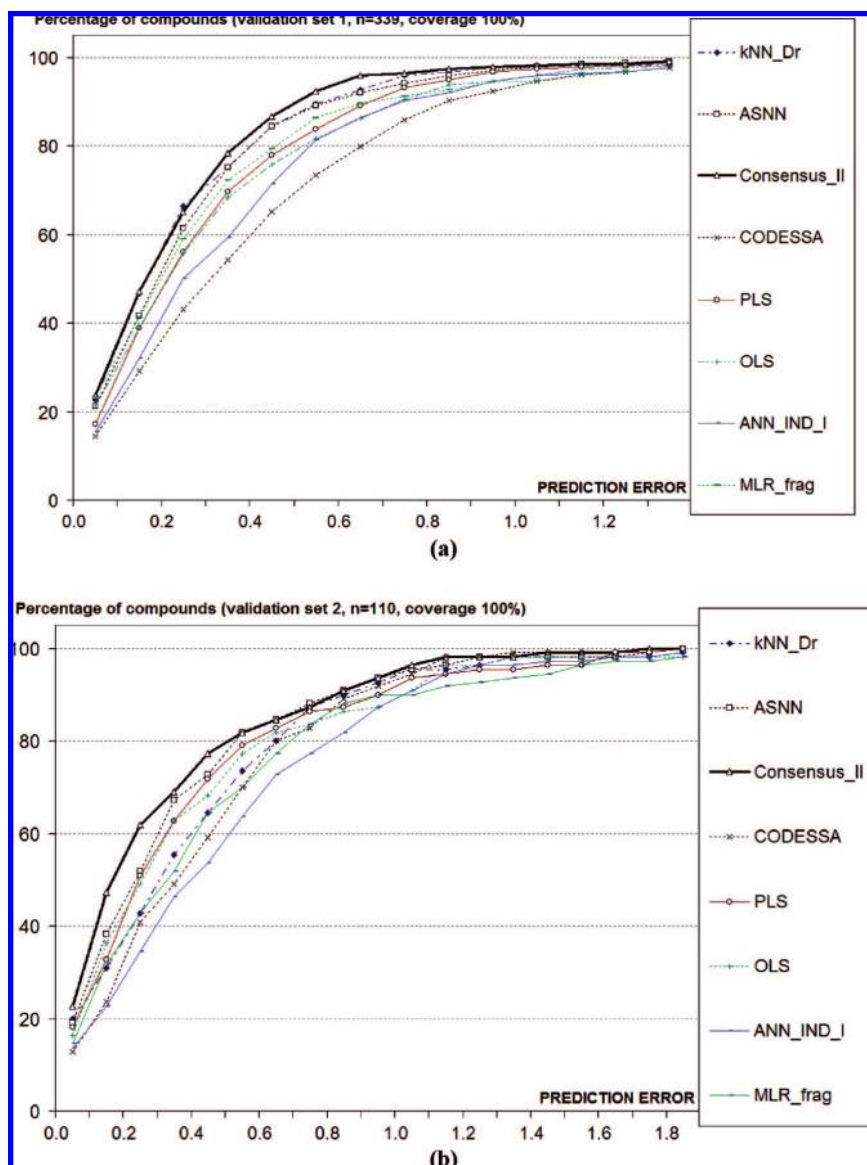


Figure 7. Percentage of compounds for (a) validation set I and (b) validation set II with full coverage (100%) vs the prediction errors obtained with individual models and the consensus model II.

Figure 6 shows the distribution of MAE values for the prediction of both validation sets I and II for nine models that used the AD for three compound sets: located within the AD of each model, outside the AD, and for all compounds. Seven (*kNN*-Dragon, *kNN*-MolconnZ, SVM-Dragon, SVM-MolconnZ, ISIDA-MLR, ISIDA-*kNN*, and ASNN) out of nine QSAR models that used the AD showed improvement in the prediction accuracy for both validation sets as a result of excluding those compounds outside the AD. The results of OLS and PLS practically did not change after applying the AD criteria. This is not surprising given that there were only very few compounds that were outside of the structural AD in these two models.

Overall, we conclude that the use of the AD generally ensures a higher accuracy of prediction for the external sets. However, we should note that the higher accuracy of prediction comes at the expense of reducing the chemical space coverage by the models. It may appear as a deficiency of the modeling with AD. However, one should remember that by default any QSAR model development is restricted to interpolation within the training set data, whereas any

external prediction is by default a model extrapolation attempt. Thus, the AD should be a natural attribute of every training set model irrespective of the descriptor types and optimization methods used. The scientific question that should continue to be explored is how flexible the definition of the AD should be, taking into account the specific distribution of the training set data in the descriptor space and the type of model optimization techniques. All of our groups are actively investigating this important issue.

Consensus Modeling. So far, we have explored and compared the performance of models implemented within individual groups that have collaborated on this project. We have demonstrated that, for the most part, all models succeeded in achieving reasonable accuracy of external prediction, especially when using the AD. It then appeared natural to bring all of the models together to explore the power of consensus prediction, which could be done in several ways. The simplest one is to average all 15 individual predictions for each external compound without considering the applicability domains. The results (see Tables 3 and 4) show very clearly that in all instances, that is, for the training

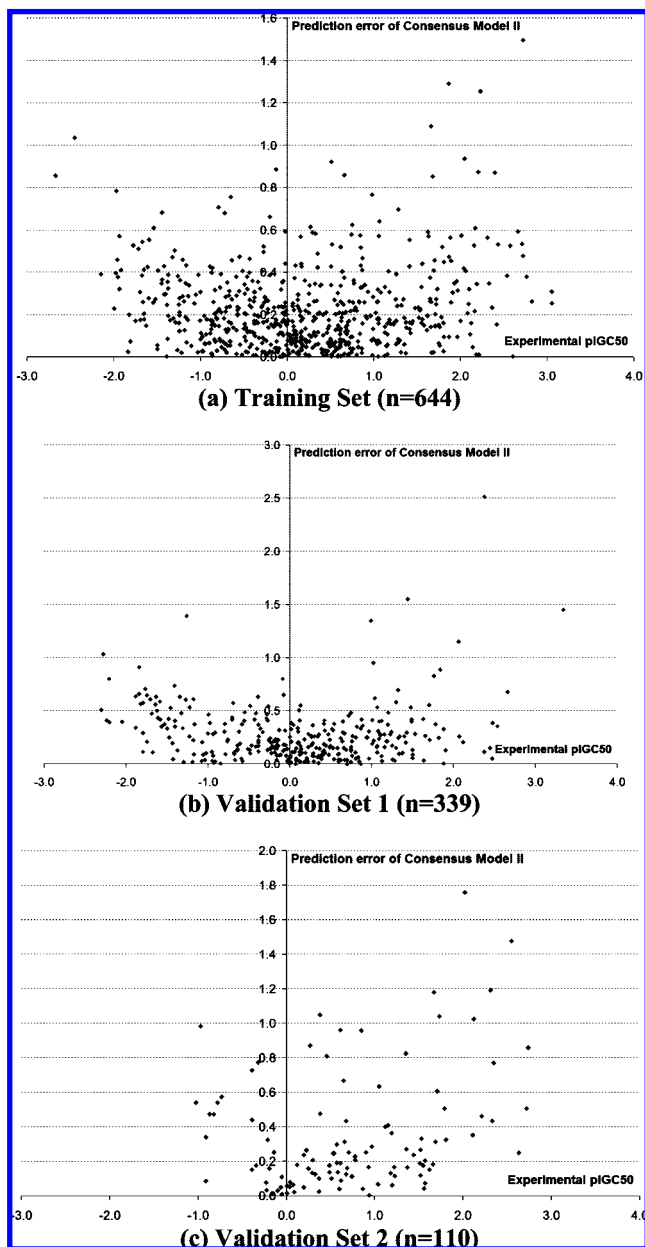


Figure 8. Prediction errors (calculated by consensus model II) versus experimental pIC_{50} for the training (a) and the two validation sets (b and c).

set ($Q_{\text{abs}}^2 = 0.90$ and $\text{MAE} = 0.25$) and both the first ($R_{\text{abs}}^2 = 0.85$ and $\text{MAE} = 0.29$) and the second ($R_{\text{abs}}^2 = 0.67$ and $\text{MAE} = 0.39$) validation sets, **consensus model I** was generally superior to any individual constituent model, except for being *on par* with the ASNN (but the latter had a lower coverage). It should be emphasized that we used both validation sets I and II to evaluate the performance of the models developed with the modeling set of 644 compounds (that in some cases was additionally subdivided into training and test sets) but not to choose the best-performing models for future use since, in real life, the models are only expected to be used in a prospective fashion. As stated above, all 15 models have demonstrated a respectable performance for validation set I, but some of them were less accurate in predicting validation set II. It was quite revealing to observe the impressive stability of consensus model I, which was not perturbed even by models with relatively low prediction accuracies for set II. These results prove that combinatorial

QSAR modeling and consensus prediction afford the most accurate prediction of the external data sets.

While we were satisfied with the results of consensus model I, we have explored additional schemes for consensus prediction. **Consensus model II** was constructed by averaging all available predicted values taking into account the applicability domain of each individual model. Thus, in this case, we used only 9 of 15 models that had the AD defined. Since each model had its unique way of defining the AD, each external compound could be found within the AD of anywhere between one and nine models; so for averaging, we only used models covering the compound. The advantage of this data treatment is that the overall coverage of the prediction is still high because it was rare to have an external compound outside the ADs of all available models. The results (**consensus model II** in Table 2) showed that the prediction accuracy for both the modeling set ($\text{MAE} = 0.22$) and validation sets I and II (0.27 and 0.34, respectively) was again the best compared to any individual model. The same observation could be made for the correlation coefficient R_{abs}^2 . The coverage of consensus model II was 100% for all three data sets. As a corollary, we also examined **consensus model IIB**, which was the same as consensus model II but without using the AD of the nine constituent models. We found that all results were practically the same (see Table 2). It was interesting to observe that, according to a standard statistical Fisher test, there was no significant difference between the statistical parameters of consensus models I and II. Again, this observation suggests that consensus models afford both high space coverage and a high accuracy of prediction.

Figure 7 presents another way of comparing the prediction accuracies of individual models versus that of the consensus model. We plotted the percentages of compounds for validation set I (Figure 7a) and validation set II (Figure 7b) versus the prediction errors obtained with individual models or consensus model II. These plots show that, for any given error threshold, the consensus model consistently predicts the largest number of compounds within this threshold versus that of any of the individual models.

To get a deeper insight into model performance, we have examined the plot displaying the prediction errors (i.e., absolute value of the difference between the predicted and the experimental toxicities) calculated with consensus model II versus the experimental pIC_{50} values (Figure 8). For the modeling set (Figure 8a) and validation set I (Figure 8b), the compounds with extreme values of pIC_{50} (i.e., less than -1.5 and higher than $+1.5$) were often associated with a large prediction error. On the other hand, for validation set II (Figure 8c), there was no obvious correlation between the prediction errors and the experimental pIC_{50} . This observation is likely due to greater dissimilarity between validation set II and both the modeling set and validation set I, which was illustrated in Figure 2.

Finally, **consensus model III** was constructed to examine whether the most conservative approach to selecting models for consensus prediction could prove the most accurate. Under the consensus model II scenario, we made a prediction for an external compound if it was found within the AD of at least one model. Here, we have looked at a progressively smaller number of compounds that would be found within the AD of at least one model (most permissive), two models,

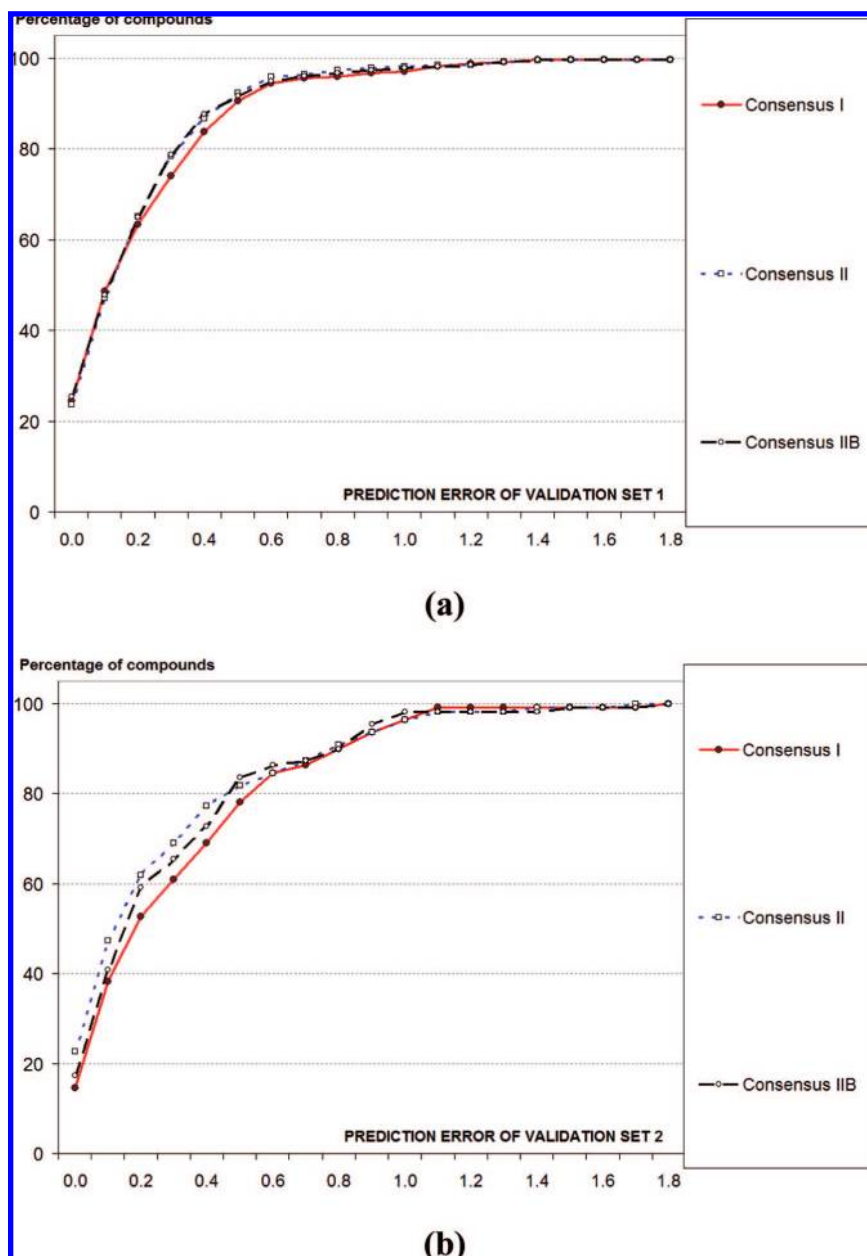


Figure 9. Percentage of compounds for (a) validation set 1 and (b) validation set 2 with full coverage (100%) vs the prediction errors obtained with **consensus models I, II, and IIB.**

and so forth, and up to all nine models (most conservative). In addition, we also refined the predicted values in consensus model III by excluding results that had a large deviation compared to the average values of all available predictions (according to the Grubbs's statistics). The results are shown in Table 4.

As one would expect, for the first validation set, both the correlation coefficient and prediction accuracy are consistent irrespective of the number of models used for consensus prediction. However, the coverage decreases progressively, reaching about 50% for the most conservative model. This result is consistent with the fairly similar and high prediction power of all individual models. On the other hand, for the second validation set, the predictive power (i.e., both R_{abs}^2 and MAE) improves to some extent, but the coverage decreases dramatically, reaching only slightly above 20% for the most conservative model. This sharp decrease indicates again that the second validation set contains

a large fraction of compounds that are more dissimilar to the modeling set than those in validation set I. It also highlights dissimilarity between constituent models that capture different trends within experimental data and have different applicability domain definitions as well.

In summary, we observe that all consensus models afford consistently high prediction accuracy for both the modeling set and validation sets I and II (cf. Tables 2 and 3). The use of AD for consensus models does not seem to have a strong effect on their prediction power but does decrease the space coverage if used conservatively (Table 4). We also observe that, for this data set, ASNN affords predictive power comparable with that of consensus models (given the same space coverage, Table 4). However, the possible advantage of the ASNN model becomes obvious only when examining the results of predicting the validation set II; other individual methods have demonstrated similar performance as applied to both the training set and validation set I. We plan to

examine whether any single approach will emerge as the most reliable as applied to other data sets that our collaborative plans to examine. However, at the moment, we could confidently conclude that the use of combinatorial QSAR and consensus prediction appeals as the methodology of choice in modeling complex toxicity data sets.

Is the Use of AD Necessary? For any individual model, using an AD is definitely critical, even if the better accuracy is often balanced with the low coverage of the external validation set. Consensus models clearly lead to superior prediction results as compared to any individual constituent model. However, consensus models I (average of the 15 models, no AD at all), II (average of the nine models implementing AD), and IIB (average of the nine models without taking into account their AD) yield very similar results for the modeling set and the two validation sets (cf. Table 2). We have applied the statistical Fisher test and have come to the conclusion that all three consensus models are not significantly different. Figure 9 confirms this conclusion by showing that percent compounds in validation sets I and II predicted within a certain error is practically the same for all three consensus models, I, II, and IIB. One may still argue that consensus model II is somewhat better than model I, especially for validation set II: the difference is not significant but seems quite noticeable. If 0.5 log units is considered as a cutoff for the prediction error, the percentage of compounds that could be predicted within this error by consensus model II is 5–10% larger than that using consensus model I (see Figure 9). We conclude that, whereas the use of AD seems imperative for developing individual models—at least for this aquatic toxicity data set—consensus modeling seems to make the use of AD less important. This conclusion is somewhat surprising, and it should be tested on additional data sets. However, if universally true, it will certainly simplify consensus model development.

5. CONCLUSIONS

Several QSAR approaches practiced by six contributing laboratories have been used to develop toxicity models of a large set of diverse organic compounds tested in *T. pyriformis*. The resulting models, most of which have incorporated specific applicability domains, were validated by predicting the toxicity of two relatively large external sets. We found that all models were consistently accurate for the training set and showed somewhat different but comparable performance for validation set I, which was selected from the original large experimental set. However, the models diverged in their performance as applied to validation set II, which included compounds chemically different from the training set. Here, the use of the applicability domain improved the prediction accuracy using individual models; however, the use of AD also decreased the coverage of validation set II (to a different degree for different models), making it difficult to compare individual model performance. Formally, the highest accuracies were achieved by SVM-Dragon and ASNN approaches (0.83 and 0.75, respectively), but this required a decrease in space coverage (to ca. 53% and ca. 72%, respectively); thus, arguably, ASNN had a better balance between the space coverage and accuracy. However, the most significant single result of our studies is the demonstrated superior performance of the consensus

modeling approach when all models are used concurrently and predictions from individual models are averaged. We have shown that both the predictive accuracy and coverage of the final consensus QSAR models were superior as compared to these parameters for individual models. The consensus models appeared robust in terms of being insensitive to both incorporating individual models with low prediction accuracy and the inclusion or exclusion of the AD. Another important result of this study is the power of addressing complex problems in computational toxicology by forming a virtual collaboratory of independent research groups leading to the formulation and empirical testing of best practices in predictive toxicology. This latter endeavor is especially critical in light of the growing interest of regulatory agencies to developing most reliable and predictive models for environmental risk assessment⁷² and placing such models in the public domain. We will make all of our models available to interested scientists upon request and will collaborate toward establishing a publicly available Web server for predicting aquatic toxicity.

ACKNOWLEDGMENT

Although all experimental data were collected from public sources, the authors are thankful to Prof. T. Schultz for many years of research on chemical toxicity in *T. pyriformis* that made this study possible. The UNC coauthors wish to thank Dr. Alexander Golbraikh for his help with programming and interpretation of results. The ULP coauthors thank Prof. Alan R. Katritzky for providing them with the CODESSA-Pro program. This study was supported in part by the NIH RoadMap grant GM076059 and by EPA STAR grant RD832720 (UNC) and partially supported by the Go-Bio BMBF grant 0313883 (AZ-31P4556) (VCCLAB).

Supporting Information Available: Details of the 1093 compounds used in this study. This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- (1) Putman, D. L.; Clarke, J. J.; Escobar, P.; Gudi, R.; Kršmanovic, L. S.; Pant, K.; Wagner, V. O., III; San, R. H. C.; Jacobson-Kram, D. Genetic Toxicity. In *Toxicological Testing Handbook*; Jacobson-Kram, D., Keller, K. A., Eds.; Informa Healthcare: New York, 2006; Vol. 6, pp. 185–248.
- (2) Hengstler, J. G.; Foth, H.; Kahl, R.; Kramer, P. J.; Lilienblum, W.; Schulz, T.; Schweinfurth, H. The REACH concept and its impact on toxicological sciences. *Toxicology* **2006**, *220*, 232–239.
- (3) Richard, A. M. Future of toxicology—predictive toxicology: An expanded view of “chemical toxicity”. *Chem. Res. Toxicol.* **2006**, *19*, 1257–1262.
- (4) Klopman, G.; Zhu, H.; Fuller, M. A.; Saiakhov, R. D. Searching for an enhanced predictive tool for mutagenicity. *SAR QSAR Environ. Res.* **2004**, *15*, 251–263.
- (5) Richard, A. M.; Benigni, R. AI and SAR approaches for predicting chemical carcinogenicity: Survey and status report. *SAR QSAR Environ. Res.* **2002**, *13*, 1–19.
- (6) Yang, C.; Benz, R. D.; Cheeseman, M. A. Landscape of current toxicity databases and database standards. *Curr. Opin. Drug Discovery Dev.* **2006**, *9*, 124–133.
- (7) Stouch, T. R.; Kenyon, J. R.; Johnson, S. R.; Chen, X. Q.; Doweyko, A.; Li, Y. In silico ADME/Tox: why models fail. *J. Comput.-Aided Mol. Des.* **2003**, *17*, 83–92.
- (8) Golbraikh, A.; Tropsha, A. Beware of q²! *J. Mol. Graphics Modell.* **2002**, *20*, 269–276.
- (9) Kubinyi, H.; Hamprecht, F. A.; Mietzner, T. Three-dimensional quantitative similarity-activity relationships (3D QSiAR) from SEAL similarity matrices. *J. Med. Chem.* **1998**, *41*, 2553–2564.
- (10) Van Drie, J. H. Pharmacophore discovery—lessons learned. *Curr. Pharm. Des.* **2003**, *9*, 1649–1664.

- (11) Tropsha, A.; Gramatica, P.; Gombar, V. K. The Importance of Being Earnest: Validation is the Absolute Essential for Successful Application and Interpretation of QSPR Models. *Quant. Struct.-Act. Relat. Comb. Sci.* **2003**, *22*, 69–77.
- (12) Schultz, T. W.; Netzeva, T. I. Development and evaluation of QSARs for ecotoxic endpoints: The benzene response-surface model for *Tetrahymena* toxicity. In *Modeling Environmental Fate and Toxicity*; Cronin, M. T. D., Livingstone, D. J., Eds.; CRC Press: Boca Raton, FL, 2004; Vol. 4, Chapter 12, pp. 265–284.
- (13) Schultz, T. W. Structure-toxicity relationships for benzenes evaluated with *Tetrahymena pyriformis*. *Chem. Res. Toxicol.* **1999**, *12*, 1262–1267.
- (14) Netzeva, T. I.; Schultz, T. W. QSARs for the aquatic toxicity of aromatic aldehydes from *Tetrahymena* data. *Chemosphere* **2005**, *61*, 1632–1643.
- (15) Schultz, T. W.; Sinks, G. D.; Miller, L. A. Population growth impairment of sulfur-containing compounds to *Tetrahymena pyriformis*. *Environ. Toxicol.* **2001**, *16*, 543–549.
- (16) Schultz, T. W.; Yarbrough, J. W.; Woldemeskel, M. Toxicity to *Tetrahymena* and abiotic thiol reactivity of aromatic isothiocyanates. *Cell Biol. Toxicol.* **2005**, *21*, 181–189.
- (17) Aptula, A. O.; Roberts, D. W.; Cronin, M. T. D.; Schultz, T. W. Chemistry-toxicity relationships for the effects of Di- and trihydroxy-benzenes to *Tetrahymena pyriformis*. *Chem. Res. Toxicol.* **2005**, *18*, 844–854.
- (18) Schultz, T. W.; Hewitt, M.; Netzeva, T. I.; Cronin, M. T. D. Assessing applicability domains of toxicological QSARs: Definition, confidence in predicted values, and the role of mechanisms of action. *QSAR Comb. Sci.* **2007**, *26*, 238–254.
- (19) MolConnZ, version 4.05; eduSoft LC: Ashland, VA, 2003.
- (20) Kier, L. B. Inclusion of Symmetry As A Shape Attribute in Kappa-Index Analysis. *Quant. Struct.-Act. Relat.* **1987**, *6*, 8–12.
- (21) Kier, L. B.; Hall, L. H. A Differential Molecular Connectivity Index. *Quant. Struct.-Act. Relat.* **1991**, *10*, 134–140.
- (22) Hall, L. H.; Mohny, B.; Kier, L. B. The Electrotological State – An Atom Index for Qsar. *Quant. Struct.-Act. Relat.* **1991**, *10*, 43–51.
- (23) Hall, L. H.; Kier, L. B. Determination of Topological Equivalence in Molecular Graphs from the Topological State. *Quant. Struct.-Act. Relat.* **1990**, *9*, 115–131.
- (24) Varnek, A.; Fourches, D.; Solov'ev, V. P.; Baulin, V. E.; Turanov, A. N.; Karandashev, V. K.; Fara, D.; Katritzky, A. R. "In silico" design of new uranyl extractants based on phosphoryl-containing podands: QSPR studies, generation and screening of virtual combinatorial library, and experimental tests. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 1365–1382.
- (25) Varnek, A.; Solov'ev, V. P. "In silico" design of potential anti-HIV actives using fragment descriptors. *Comb. Chem. High Throughput. Screening* **2005**, *8*, 403–416.
- (26) Varnek, A.; Fourches, D.; Hoonakker, F.; Solov'ev, V. P. Substructural fragments: an universal language to encode reactions, molecular and supramolecular structures. *J. Comput.-Aided Mol. Des.* **2005**, *19*, 693–703.
- (27) CODESSA-PRO, version 2007; CompuDrug International, Inc.: Sedona, AZ, 2007.
- (28) DRAGON for Windows (Software for Molecular Descriptor Calculations), version 5.4; Talet s.r.l.: Milan, Italy, 2006.
- (29) Todeschini, R.; Consonni, V. *Handbook of Molecular Descriptors*; Wiley: Weinheim, Germany, 2000.
- (30) CORINA, version 3.2; Molecular Networks GmbH: Erlangen, Germany, 2002.
- (31) Kier, L.; Hall, L. *Molecular Structure Description: The Electrotological State*; London, 1999.
- (32) Hall, L. H.; Kier, L. B. Electrotological state indices for atom types – a novel combination of electronic, topological, and valence state information. *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 1039–1045.
- (33) Huuskonen, J. J.; Livingstone, D. J.; Tetko, I. V. Neural network modeling for estimation of partition coefficient based on atom-type electrotological state indices. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 947–955.
- (34) Tetko, I. V.; Tanchuk, V. Y.; Villa, A. E. Prediction of n-octanol/water partition coefficients from PHYSPROP database using artificial neural networks and E-state indices. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1407–1421.
- (35) Virtual Computational Chemistry Laboratory. <http://www.vcclab.org/lab/pclient/> (accessed: Oct 1, 2007).
- (36) Tetko, I. V.; Gasteiger, J.; Todeschini, R.; Mauri, A.; Livingstone, D.; Ertl, P.; Palyulin, V. A.; Radchenko, E. V.; Zefirov, N. S.; Makarenko, A. S.; Tanchuk, V. Y.; Prokopenko, V. V. Virtual computational chemistry laboratory--design and description. *J. Comput.-Aided Mol. Des.* **2005**, *19*, 453–463.
- (37) Cherkasov, A.; Ban, F. Q.; Li, Y.; Fallahi, M.; Hammond, G. L. Progressive docking: A hybrid QSAR/docking approach for accelerating in silico high throughput screening. *J. Med. Chem.* **2006**, *49*, 7466–7478.
- (38) Karakoc, E.; Sahinalp, S. C.; Cherkasov, A. Comparative QSAR- and fragments distribution analysis of drugs, druglikes, metabolic substances, and antimicrobial compounds. *J. Chem. Inf. Model.* **2006**, *46*, 2167–2182.
- (39) Cherkasov, A. Can 'bacterial-metabolite-likeness' model improve odds of 'in silico' antibiotic discovery. *J. Chem. Inf. Model.* **2006**, *46*, 1214–1222.
- (40) MOE, version 2005.06; Chemical Computing Group: Montreal, Quebec, Canada, 2005.
- (41) Zheng, W.; Tropsha, A. Novel variable selection quantitative structure--property relationship approach based on the k-nearest-neighbor principle. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 185–194.
- (42) Vapnik, V. N. *The Nature of Statistical Learning Theory*, 2nd ed.; Springer-Verlag: New York, 2000.
- (43) Kovatcheva, A.; Golbraikh, A.; Oloff, S.; Xiao, Y. D.; Zheng, W.; Wolschann, P.; Buchbauer, G.; Tropsha, A. Combinatorial QSAR of ambergris fragrance compounds. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 582–595.
- (44) Oloff, S.; Mailman, R. B.; Tropsha, A. Application of validated QSAR models of D1 dopaminergic antagonists for database mining. *J. Med. Chem.* **2005**, *48*, 7322–7332.
- (45) Golbraikh, A.; Shen, M.; Xiao, Z.; Xiao, Y. D.; Lee, K. H.; Tropsha, A. Rational selection of training and test sets for the development of validated QSAR models. *J. Comput.-Aided Mol. Des.* **2003**, *17*, 241–253.
- (46) Karelson, M.; Maran, U.; Wang, Y. L.; Katritzky, A. R. QSPR and QSAR models derived using large molecular descriptor spaces. A review of CODESSA applications. *Collect. Czech. Chem. Commun.* **1999**, *64*, 1551–1571.
- (47) Varnek, A.; Kireeva, N.; Tetko, I. V.; Baskin, I. I.; Solov'ev, V. P. Exhaustive QSPR studies of a large diverse set of ionic liquids: How accurately can we predict melting points? *J. Chem. Inf. Model.* **2007**, *47*, 1111–1122.
- (48) Grubbs, F. E. Procedures for Detecting Outlying Observations in Samples. *Technometrics* **1969**, *11*, 1–21.
- (49) *Software for multilinear regression analysis and variable subset selection by Genetic Algorithm*, version 1.0 beta; Talet s.r.l.: Milan, Italy, 2004.
- (50) Leardi, R.; Boggia, R.; Terrile, M. Genetic Algorithms As A Strategy for Feature-Selection. *J. Chemom.* **1992**, *6*, 267–281.
- (51) Todeschini, R.; Consonni, V.; Maiocchi, A. The K correlation index: theory development and its application in chemometrics. *Chemom. Intell. Lab. Syst.* **1999**, *46*, 13–29.
- (52) Wehrens, R.; Putter, H.; Buydens, L. M. C. The bootstrap: a tutorial. *Chemom. Intell. Lab. Syst.* **2002**, *54*, 35–52.
- (53) *Unscrambler*, version 9.1; Camo Process AS: Oslo, Norway, 2005.
- (54) Martens, H.; Næs, T. *Multivariate calibration*; Wiley: Chichester, U.K., 1998.
- (55) Wold, S.; Sjostrom, M.; Eriksson, L. PLS-regression: a basic tool of chemometrics. *Chemom. Intell. Lab. Syst.* **2001**, *58*, 109–130.
- (56) Tetko, I. V. Associative neural network. *Neural Proc. Lett.* **2002**, *16*, 187–199.
- (57) Tetko, I. V. Neural network studies. 4. Introduction to associative neural networks. *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 717–728.
- (58) Tetko, I. V. Efficient partition of learning data sets for neural network training. *Neural Networks* **1997**, *10*, 1361–1374.
- (59) Tetko, I. V. Neural Network Studies. 1. Comparison of Overfitting and Overtraining. *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 826–833.
- (60) Witten, I. H.; Frank, E. *Data Mining: Practical machine learning tools and techniques*, 2nd ed.; Morgan Kaufmann: San Francisco, CA, 2005.
- (61) Netzeva, T. I.; Worth, A.; Aldenberg, T.; Benigni, R.; Cronin, M. T.; Gramatica, P.; Jaworska, J. S.; Kahn, S.; Klopman, G.; Marchant, C. A.; Myatt, G.; Nikolova-Jeliazkova, N.; Patlewicz, G. Y.; Perkins, R.; Roberts, D.; Schultz, T.; Stanton, D. W.; van de Sandt, J. J.; Tong, W.; Veith, G.; Yang, C. Current status of methods for defining the applicability domain of (quantitative) structure-activity relationships. The report and recommendations of ECVAM Workshop 52. *Altern. Lab. Anim.* **2005**, *33*, 155–173.
- (62) Tetko, I. V.; Bruneau, P.; Mewes, H. W.; Rohrer, D. C.; Poda, G. I. Can we estimate the accuracy of ADME-Tox predictions. *Drug Discovery Today* **2006**, *11*, 700–707.
- (63) Atkinson, A. C. *Plots, transformations and regression*; Clarendon Press: Oxford, U.K., 1985.
- (64) Eriksson, L.; Jaworska, J.; Worth, A. P.; Cronin, M. T.; McDowell, R. M.; Gramatica, P. Methods for reliability and uncertainty assessment and for applicability evaluations of classification- and regression-based QSARs. *Environ. Health Perspect.* **2003**, *111*, 1361–1375.
- (65) Renner, S.; Schneider, G. Scaffold-hopping potential of ligand-based similarity concepts. *ChemMedChem* **2006**, *1*, 181–185.
- (66) Willett, P.; Winterman, V. A Comparison of Some Measures for the Determination of Intermolecular Structural Similarity Measures of

- Intermolecular Structural Similarity. *Quant. Struct.-Act. Relat.* **1986**, 5, 18–25.
- (67) Klopman, G.; Saiakhov, R.; Rosenkranz, H. S.; Hermens, J. L. M. Multiple Computer-Automated Structure Evaluation program study of aquatic toxicity I: Guppy. *Environ. Toxicol. Chem.* **1999**, 18, 2497–2505.
- (68) Klopman, G.; Saiakhov, R.; Rosenkranz, H. S. Multiple computer-automated structure evaluation study of aquatic toxicity II. Fathead minnow. *Environ. Toxicol. Chem.* **2000**, 19, 441–447.
- (69) Tetko, I. V.; Solov'ev, V. P.; Antonov, A. V.; Yao, X.; Doucet, J. P.; Fan, B.; Hoonakker, F.; Fourches, D.; Jost, P.; Lachiche, N.; Varnek, A. Benchmarking of linear and nonlinear approaches for quantitative structure-property relationship studies of metal complexation with ionophores. *J. Chem. Inf. Model.* **2006**, 46, 808–819.
- (70) Tetko, I. V.; Villa, A. E.; Livingstone, D. J. Neural network studies. 2. Variable selection. *J. Chem. Inf. Comput. Sci.* **1996**, 36, 794–803.
- (71) Gramatica, P. Principles of QSAR models validation: internal and external. *QSAR Comb. Sci.* **2007**, 26, 694–701.
- (72) Yang, C.; Richard, A. M.; Cross, K. P. The Art of Data Mining the Minefields of Toxicity Databases to Link Chemistry to Biology. *Curr. Comput.-Aided Drug Des.* **2006**, 2, 135–150.

CI700443V