

Toward Novel Universal Descriptors: Charge Fingerprints

Frank R. Burden,^{‡,§} Mitchell J. Polley,[†] and David A. Winkler^{*,†}

CSIRO Molecular & Health Technologies, Private Bag 10, Clayton South MDC, Clayton, Victoria 3169, Australia, School of Chemistry, Monash University, Clayton, Victoria 3168, Australia, and SciMetrics Limited, 548 Canning Street, Carlton North, Victoria 3054, Australia

Received August 19, 2008

Although there are a myriad of molecular descriptors for QSAR described in the literature, many descriptors contain similar information as others or are information poor. Recent work has suggested that it may be possible to discover a relatively small pool of ‘universal’ descriptors from which subsets can be drawn to build a diverse variety of models. We describe a new type of descriptor of this type, the charge fingerprint. This descriptor family can build good QSAR models of a diverse range of physicochemical and biological properties and can be calculated quickly and easily. It appears to be useful for modeling large data sets and has potential for screening large virtual libraries.

INTRODUCTION

Although Quantitative Structure–Activity Relationship (QSAR) methods have been used for over 40 years,^{1,2} it has only recently become clear that a set of “universal descriptors” may be achievable. The motivation for the search for these descriptors was that the current myriad of molecular descriptors are often highly correlated and describe similar information, and many are very information-poor. The premise is that by devising better descriptors, rather than more of them, it may be possible to create a relatively small pool of information-rich descriptors from which suitable subsets could be drawn for a wide variety of modeling applications. Even a small pool would have sufficient flexibility to model a vast range of properties. For example, with a pool size of 50, and 20 being used in a model, there are approximately 10^{32} different combinations of descriptors. Such sparse descriptor sets would reduce feature selection requirements and allow the construction of models with optimum prediction (generalization) ability. If this can be achieved the problem then largely becomes one of selecting in a supervised (context-dependent) manner the best set of descriptors for a given problem.

Wildman and Crippen proposed descriptors based on atomic contributions³ in the same year that Stanton proposed the use of BCUT diversity metrics⁴ in QSAR and QSPR studies.⁵ Shortly after, Ertl and co-workers published their Topological Polar Surface Area (TPSA) descriptor, which had great utility for modeling bioavailability,⁶ and Labute showed binned surface property density having general applicability in QSAR.⁷ In 2002 Visco et al. described their molecular signature descriptor, which is a systematic coding of a molecular graph over a set of defined atom types.⁸ More recently Gallegos et al. published a pair of universal descriptors; one a quantum similarity measure and the other a fragment self-similarity measure.⁹ Clark and colleagues

also reported work on surface property descriptors^{10,11} derived from molecular orbital properties, and Sun generated models of several important partition coefficients using descriptors akin to substructural fingerprints.¹²

Atomic partial charges have been used in previous QSAR models generally with limited success (e.g., see ref 13). There are also a relatively large number of composite charge-related descriptors such as the most negative atom charge, the most positive atom charge, various sums of charges on atoms, and such like that have been used in QSAR models.¹⁴ These studies have generally used charges on specific, relevant atoms as descriptors. Recently, Dominy and Shakhnovich¹⁵ used a clustering of atom property pairs (atomic charge and Born radius) to derive improved knowledge-based potentials for docking and scoring, showing that these properties are important for modeling the interactions of small molecules with protein targets.

We have studied the information-rich descriptor problem in our group and are developing descriptors that we hope will approach the ‘universal’ ideal. In this paper, we introduce a new universal descriptor set for QSAR studies: the charge fingerprint. We describe how it is computed and its relationship to other descriptors and illustrate applicability to QSAR using a diverse range of modeled properties.

MATERIALS AND METHODS

Data Sets. We generated QSAR models using five relatively large and diverse data sets to assess the efficacy of descriptors and choose bin numbers and boundaries. These data sets were as follows:

- set of 245 benzodiazepines acting at the GABA_A receptor.¹⁶
- set of 503 compounds exhibiting acute toxicity toward a ciliate *Tetrahymena pyriformis*.¹⁷
- set of 13,474 experimental log P values from the PHYSPROP database.¹⁸
- set of 1412 very diverse inhibitors of farnesyl transferase.¹⁹

* Corresponding author e-mail: dave.winkler@csiro.au.

[‡] Monash University.

[§] SciMetrics.

[†] CSIRO Molecular & Health Technologies.

Table 1. Bin Definitions (esu) for Elements Included in Charge Fingerprints

element	low bin	medium bin	high bin
H	<0.1	0.1–0.2	>0.2
C	<0.1	0.1–0.2	>0.2
N	<–0.2	–0.2 to 0.0	>0.0
O	<–0.3	–0.3 to –0.2	>–0.2
Si	<0.1	0.1–0.2	>0.2
P	<–0.2	–0.2 to 0.0	>0.0
S	<0.0	0.0–0.2	>0.2

- cyclooxygenase 2 (COX-2) inhibition data set of 454 compounds from the Rao and Stockfish²⁰

These sets were chosen to compare the efficacy of the descriptors to model drug properties (COX-2, benzodiazepine, and farnesyl transferase inhibitor sets), toxicity (ciliate acute toxicity set), and physical properties and octanol–water partition (dependent on lipophilic properties). Although the data set sizes differed considerably, any significant deficiencies in the ability of the descriptors to model certain kinds of properties would become evident.

Software. Descriptor and regression calculations were carried out using a purpose-developed software package called MolSAR, written in the Python programming language. MolSAR can rapidly calculate a range of molecular descriptors and perform several types of linear and nonlinear regression. MolSAR, including the charge fingerprints and Bayesian neural net, is available as part of the Bio-RAD KnowItAll package (<http://www.bio-rad.com>).

Descriptors. The charge fingerprint is a simple but powerful type of molecular description. The calculation of the fingerprints is relatively straightforward. Atom charges are computed using electronegativity equalization methods for each structure. We explored both Gasteiger's method of charge equalization²¹ and also the more recent electronegativity equalization method (EEM) based on Sanderson's equations, as reported by Mortier²² and Bultinck,²³ which has its origins in density functional theory. Other methods such as semiempirical molecular orbital methods, DFT, or *ab initio* methods can also be used to calculate atom charges if the bin boundaries are set appropriately.

Once the atom charges are calculated, the charges for each element type are used to populate bins. For each of the elemental types represented in the data set, the value of the charge is compared to bin-boundaries (Table 1) and a tally kept for each bin (Figure 1). The vector of bin occupancies for all element types represents the charge fingerprint. This results in a vector or fingerprint of length typically 20–25. The number of and location of charge bins for each atom type was determined by calculating histograms of charge using several large, diverse data sets. The sets were first screened to remove atom types poorly represented in the set.

We assessed the efficacy of the charge fingerprint descriptors alone and in combination with atomistic and Burden eigenvalue descriptors. The atomistic descriptors have been described previously²⁴ and are simple counts of the numbers of atoms in each hybridization state for each element in the molecule, plus numbers of rings of different sizes in the molecule. The eigenvalue descriptors have also been described previously.^{25,26} They are derived from the highest five and lowest five eigenvalues of the modified adjacency matrix for the molecules. The modified adjacency matrix has

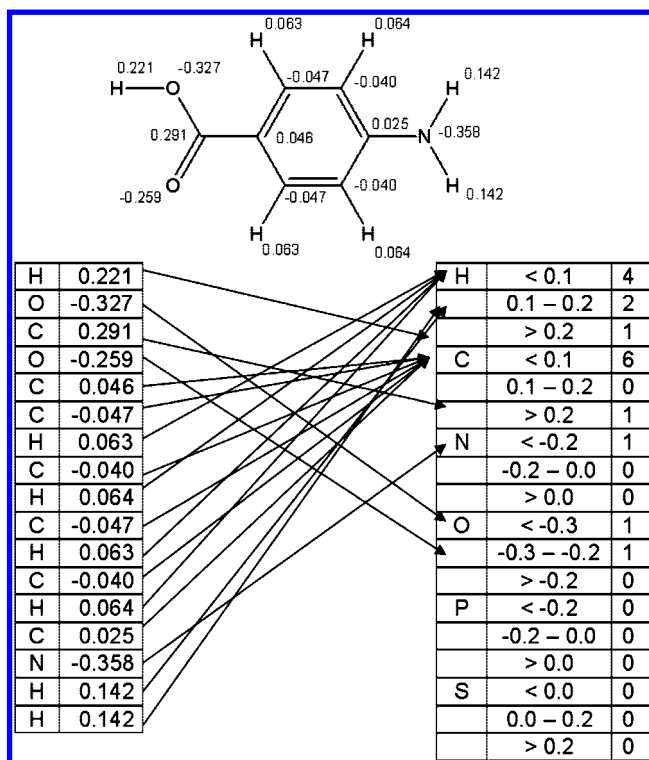


Figure 1. Example calculation of a charge fingerprint. *p*-Aminobenzoic acid with Gasteiger–Marsili charges computed (top). List of atoms with computed charges (left), with arrows indicating assignment to elemental charge bins (right), which yields the final charge fingerprint (far right column).

the same dimensions as the number of heavy atoms in the molecules, with off-diagonal elements being the inverse square root of the order of the bond joining two atoms, and the diagonal elements relating to atomic properties. In this work we did not use the empirically derived CIMI diagonal elements, rather, as a result of a theoretical analysis of the method, we used atomic hardnesses as diagonal elements in the matrix. All descriptors were mean centered and scaled to unit variance before use in QSAR modeling.

Regression Methods. We built QSAR models from the data sets using a Bayesian regularized artificial neural network (BRANN)¹⁶ and multiple linear regression (MLR) for comparison. The BRANN is a back-propagation artificial neural network that adjust the weights so that the log(evidence) is a maximum rather than minimizing the errors of a validation set prediction. This prevents overtraining and returns the most generalizable network given the architecture (number of hidden nodes). It is also robust against overfitting since the log(evidence) maximization uses the minimum necessary set of weights (the effective weights), and we have shown¹⁶ that increasing the number of hidden nodes, and hence the number of initial weights, does not increase the number of effective weights. We used three neurodes in the hidden layer of the neural network and stopped training at the maximum of the evidence. We used 80% of the data set for training the models and 20% as an independent test set never used in training. The test set was chosen using k-means clustering, with one example taken from each cluster to create the test set. For the two large physical property data sets, the test set was chosen randomly from the data set due to the time required for the k-means clustering. Although Bayesian regularized neural networks do not strictly re-

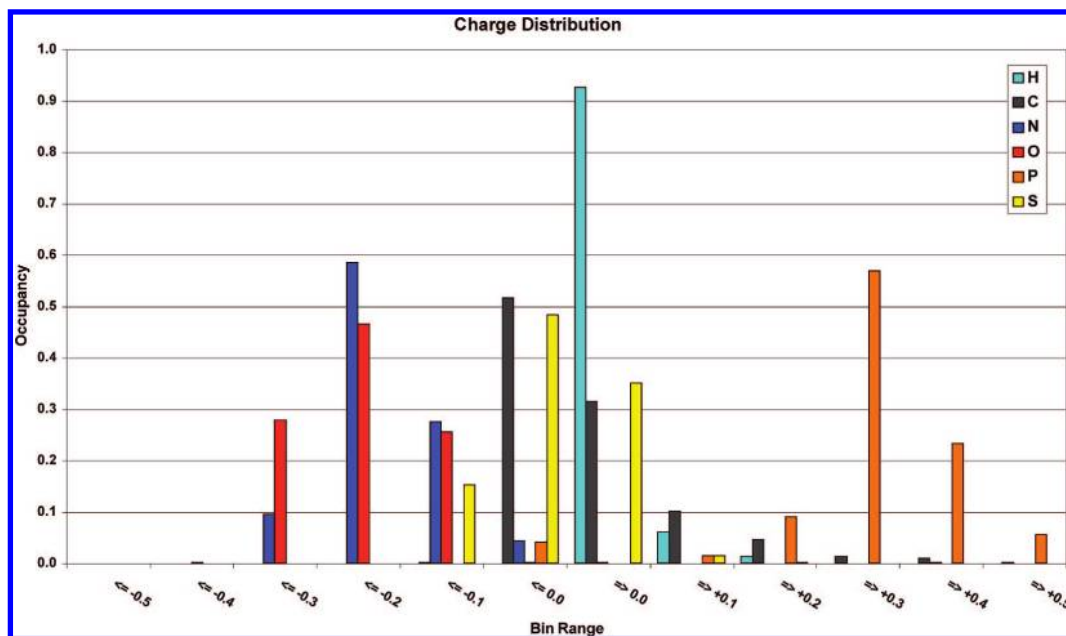


Figure 2. Binned charge distributions for elements included in charge fingerprints.

quire a test set,¹⁶ we used one to compare the relative efficacy of different descriptors in the neural network and MLR models. We use test sets rather than cross-validation methods, such as leave-one-out (LOO), as the preferred and more rigorous method of assessing predictivity of the models.

RESULTS AND DISCUSSION

Bin Numbers and Boundaries. The bin boundaries were determined through a combination of empiricism and heuristics. Initial definitions were originally derived from distributions of charges calculated on the PHYSPROP data set.¹⁸ They were subsequently optimized manually to provide peak performance across a more diverse set of data sets encompassing biological activity and physicochemical properties. We found very similar numbers of bins were required for the two charge calculation methods, and neither was superior to the other when data sets were modeled using them. Consequently we used the Gasteiger–Hückel charge calculation method, as it was applicable to a larger number of elements than the EEM method. Inspection of the charge histograms for each element type showed that three bins were sufficient to capture the variation in the charge values.

As can be seen from analysis of several data sets, the bin definitions remain appropriate for the distributions of computed charges observed in elements included in charge fingerprints (Figure 2). It was thought that halogens could be included in the fingerprints, but analysis of the charge distributions of halogen atoms in the available data sets displayed narrow ranges. All but fluorine displayed 99% occupancy of the same charge bin. Fluorine itself was split between two bins, but it was felt that there would not be enough information contained within to warrant three more descriptor columns. Final consideration went to a combined halogen atom type, unifying all the halogens under one set of three bins. However this would have merely served to allow discrimination between fluorine and non-fluorine halogens, which is better served by a simpler atomistic descriptor.²³

Correlation with Atom Types. Clustering of elemental atom charges into bins may intuitively be expected to correlate strongly with atom types based on hybridization states of atoms (i.e., with the atomistic descriptors also used in the QSAR modeling in this work). Dominy and Shakhnovich¹⁵ published a paper, while our manuscript was in preparation that is relevant. These authors raised this descriptor family correlation issue but pointed out that ‘because Gasteiger charges (and other EEM charge calculations) are generated through partial equilibration of initial formal charges, these descriptors retain information relevant to the local bonded connectivity’. They further note that while some similarity to atom types is expected, the similarity depends on the number of atom types (bins in our case) employed. It is also likely that there may be significant similarity between charge fingerprints and atomistic descriptors in a linear correlation, but higher order and interaction (cross) terms may generate important contributions in nonlinear models.

We assessed the relationship between charge fingerprints and atomistic descriptors in two ways: a simple, linear correlation matrix and a combination of these descriptors in QSAR models. In the latter case, if the charge fingerprints and atomistic descriptors contained very similar information, combinations of these descriptors would not improve QSAR models substantially. If they encode substantially different information, improved models should result when these descriptors are used together.

The linear correlations between the atomistic descriptors and charge fingerprints are shown in Table 2 for the vapor pressure data set, which was the largest and most diverse. The correlation matrices for the other data sets did not differ substantially from this example. It is clear that there are only a few places where the correlation is high, notably for hydrogen atoms in the lowest charge bin (mainly aromatic and aliphatic hydrogen atoms) that were also the most numerous. The carbon atom atomistic descriptors describing its hybridization state and/or number of connections did not

Table 2. Correlation Matrix for Charge Fingerprints and Atomistic Descriptors^a

	H	C _{sp}	C _{sp2}	C _{sp3}	C _{ar}	N _{sp}	N _{sp2}	N _{sp3}	N ⁺ _{sp3}	N _{ar}	O _{sp2}	O _{sp3}	O _{ar}
H low	0.97												
H mid	0.08												
H high	0.22												
C low		-0.01	0.17	0.63	0.48								
C mid		0.05	0.31	0.30	-0.03								
C high		-0.04	0.25	0.18	-0.01								
N low						-0.06	0.47	0.79	-0.04	0.02			
N mid						0.34	0.60	0.33	-0.01	-0.01			
N high						-0.02	0.06	0.08	0.39	-0.01			
O low											0.09	0.92	0.00
O mid											0.73	0.30	0.06
O high											0.50	0.02	0.07

^a Higher correlations are in bold.

correlate substantially with any of the carbon charge bins (<0.65 maximum). There was higher correlation between the nitrogen Nsp3 (three connections) and the lowest nitrogen bin (0.79) and also between the Nsp2 atom type and the middle nitrogen charge bin (0.60). There is a higher correlation between the oxygen atomistic descriptors and the oxygen atom charge bins, which probably reflects the more limited range of environments for oxygen atoms in most molecules, compared with carbon and nitrogen. However, it is clear from the correlation matrices, and from the QSAR modeling below, the two descriptors families do encode different information in nonlinear space. In the context of Dominy and Shakhnovich's work, the charge fingerprints are performing the function of their clustered charges, and the atomistic descriptors the function of the Born radii. The eigenvalue descriptors are encoding information about molecular connectivity, polarizability, and topology.

QSAR Modeling. Table 3 shows a summary of the QSAR models built using the three descriptors families alone and in all possible combinations. The atomistic and charge fingerprint descriptors are more successful than the eigenvalue descriptors in building QSAR models on their own. Improvement of the eigenvalue descriptors is an ongoing research interest in our group.²⁵

It is clear that combinations of the descriptor families yield more significant QSAR models than the families alone. The combination of the atomistic and charge fingerprint descriptors substantially improves model quality indicating that these two descriptor families do not contain the same information. The best combination of descriptors gives very good QSAR models, as shown by the ability of the models to predict an independent test set not used in training. Clearly, in most cases the charge fingerprint descriptors are not sufficiently information rich on their own to be used as a universal descriptor. However, the combination of the A, B, and C descriptors has properties approaching that ideal. Most of the models show that the relationship between the descriptors and the response variable is nonlinear to some degree as the MLR models have substantially lower statistical significance than the nonlinear Bayesian neural net models. It is also clear that combinations of descriptor families are capable of building very good QSAR models.

Descriptor Interpretability. The ideal QSAR model would be robust, sparse, predictive, and interpretable. In many cases such an ideal is not achievable with current descriptors and response variable mapping methods, although much effort is being expended in approaching this ideal.

Table 3. Summary of QSAR Models Built Using Several Data Sets and Various Combinations of Descriptors^a

	data set	descriptors	SEE	r ² train	SEP	r ² test	
benzodiazepine		A	0.200	0.605	0.193	0.557	C
		B	0.220	0.443	0.206	0.514	C
		C	0.165	0.745	0.219	0.505	C
	MLR	AB	0.187	0.668	0.184	0.602	C
		AC	0.184	0.659	0.207	0.501	C
		BC	0.195	0.621	0.201	0.578	C
toxicity	MLR	ABC	0.171	0.720	0.191	0.581	C
		ABC	0.184	0.698	0.199	0.588	C
		A	0.110	0.722	0.075	0.858	C
		B	0.127	0.591	0.122	0.593	C
		C	0.123	0.629	0.123	0.584	C
	MLR	AB	0.107	0.737	0.088	0.801	C
log P		AC	0.089	0.827	0.079	0.851	C
		BC	0.115	0.684	0.106	0.735	C
		ABC	0.092	0.814	0.086	0.815	C
	MLR	ABC	0.106	0.768	0.092	0.817	C
		A	0.060	0.846	0.062	0.841	R
		B	0.101	0.416	0.102	0.418	R
farnesyl transferase		C	0.068	0.788	0.068	0.787	R
	MLR	AB	0.062	0.830	0.060	0.834	R
		AC	0.054	0.873	0.054	0.867	R
		BC	0.070	0.785	0.071	0.783	R
		ABC	0.057	0.856	0.056	0.864	R
		ABC	0.063	0.846	0.062	0.853	R
COX-2 inhibition		A	0.130	0.751	0.150	0.668	C
		B	0.175	0.485	0.170	0.416	C
		C	0.139	0.712	0.145	0.678	C
	MLR	AB	0.117	0.805	0.143	0.704	C
		AC	0.124	0.776	0.144	0.696	C
		BC	0.151	0.642	0.153	0.643	C
MLR		ABC	0.113	0.818	0.130	0.763	C
		ABC	0.135	0.751	0.152	0.688	C
		A	0.134	0.590	0.135	0.631	C
	COX-2 inhibition	B	0.166	0.392	0.136	0.621	C
		C	0.118	0.673	0.124	0.661	C
		AB	0.135	0.647	0.135	0.578	C
MLR		AC	0.116	0.715	0.114	0.623	C
		BC	0.127	0.688	0.119	0.680	C
		ABC	0.109	0.797	0.119	0.668	C
		ABC	0.126	0.746	0.137	0.620	C

^a All models were built using a Bayesian regularized neural network, with MLR statistics shown for comparison. A - atomistic descriptors, B - Burden (eigenvalue) descriptors, C - charge fingerprints. C and R in the last column refer to k-means cluster or random selection of test set (20% of data set). SEE and SEP represents the standard error of estimation and the standard error of prediction, respectively (scaled 0–1). r² represents the squared correlation coefficient for training and test sets.

Consequently, QSAR modeling tends to be divided into two classes, depending on the intended outcome of the study. Interpretative modeling often uses linear modeling tools,

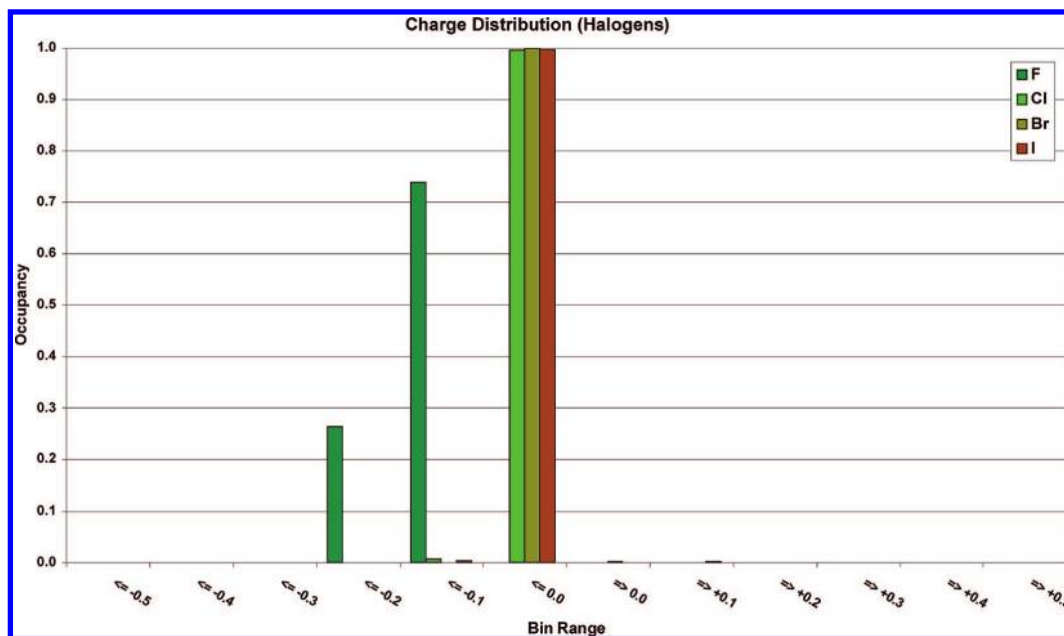


Figure 3. Binned charge distributions for halogens.

chemically relevant and interpretable descriptors, and smaller, more congeneric data sets that have usually been measured to a higher degree of accuracy. Predictive QSAR aims to model large, chemically diverse data sets that are often noisy or contain missing data. They often use computationally derived descriptors (for speed) and flexible, nonlinear structure–activity mapping methods such as neural networks and aim to be as predictive as possible so that new synthesis candidates can be assessed prior to synthesis or so that large databases or virtual libraries can be screened for hits. These predictive approaches are increasingly being used for AD-MET modeling for these reasons.²⁷

The charge fingerprint descriptors are clearly useful in building predictive models. While it may be possible to interpret the fingerprints in certain cases, this is generally not easy because the topological information is lost when the charges from all similar atoms in the molecule are hashed into the fingerprint. Charges on heteroatoms, which may be expected to correlate to some degree with hydrogen bond donor–acceptor properties, are likely to be the most interpretable descriptors.

Comparison with Other Models. The QSAR models produced using our fingerprint descriptors can be compared with other models reported in the literature for the same or similar data sets. A number of authors have presented QSAR studies for relatively small data sets of structurally similar benzodiazepines, and we have not compared these with our study that used a much larger, more chemically diverse data set. Burden²⁸ and Burden et al.^{29,30} compared Gaussian processes with MLR and neural networks for their ability to model a similar benzodiazepine data set and a smaller toxicity data set. They used different descriptors with almost twice as many being used in the final models as in the current study. They obtained similar training statistics for the benzodiazepine data set and slightly better test set r^2 values (0.71). For the toxicity test set they obtained models of almost identical statistical quality as those reported here, albeit with a data set almost half the size.

Several octanol–water (log P) models based on data from the physprop database have been reported in the literature.

Faulon³¹ used novel signature molecular and MolConnZ descriptors to model log P. He used a very small test set of 123 compounds (1% of the training set) and obtained training set r^2 of 0.88 and test set r^2 of 0.72 using MolConnZ descriptors and training set r^2 of 0.92 and test set r^2 of 0.77 using signature descriptors. These models have substantially lower predictivity than our best neural net models (training set r^2 = 0.86 and test set r^2 = 0.86). Sun¹² reported models of similar quality to those reported here using a ‘universal’ molecular descriptor based on atom type classifications. He reported a training set r^2 value of 0.91 and cross-validated q^2 value of 0.89 comparable with our more stringent test set r^2 value of 0.86. Tetko et al.³² used E-state descriptors and a backpropagation neural net ensemble to model the physprop log P data. They obtained training set r^2 values of 0.95 and test set values of 0.94 but omitted between 90 and 188 outliers. In contrast our models did not omit any outliers but would be expected to perform substantially better if this many outliers were omitted. Molnar et al.³³ also used a neural network together with atomic fragmental descriptors to model log P. They used up to 130 descriptors and 12 hidden layer neurodes and obtained r^2 values of 0.94 and 0.91–0.92 for the training and test sets, respectively. These models have higher statistical significance than those reported here but require a large number of descriptors and a more complex neural network architecture that takes longer to train. As these descriptors are new, it is not clear whether they are more generally applicable across different types of response data as our descriptors appear to be.

Rao and Stockfisch²⁰ used their PUMP-RP recursive partitioning method to analyze the COX-2 data set. They employed three types of fingerprint descriptors (ISIS public keys, DAYLIGHT Fingerprints, and Cerius) to develop QSAR models. They developed a classification model that denoted compounds as active and inactive rather than a quantitative model. They employed 89% of the data set as a training set and 11% as a test set. The results of their classification model were difficult to compare with our

continuous model is a meaningful way. No other QSAR models using this data set have been reported to our knowledge.

CONCLUSIONS

It is clear that the charge fingerprints are useful in QSAR modeling and that they encode different types of information to other descriptors studied and used in combination with them. They are easily and rapidly calculated for very large data sets and are therefore suitable for screening virtual libraries. They can be combined with other easily calculated descriptors to build very good QSAR models. They are also intuitively reasonable as recent work on improved knowledge-based docking and scoring potentials shows. They represent another step in the path toward robust, sparse molecular descriptors of wide, if not universal, applicability.

REFERENCES AND NOTES

- Hansch, C.; Fujita, T. ρ - σ - π Analysis: A method for the correlation of biological activity and chemical structure. *J. Am. Chem. Soc.* **1964**, *86*, 1616–1626.
- Leo, A.; Hansch, C.; Church, C. Comparison of parameters currently used in the study of structure-activity relationships. *J. Med. Chem.* **1969**, *12*, 766–771.
- Wildman, S. A.; Crippen, G. M. Prediction of Physicochemical Parameters by Atomic Contributions. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 868–873.
- Pearlman, R. S.; Smith, K. M. Novel Software Tools for Chemical Diversity. *Perspect. Drug. Discovery* **1998**, *9/10/11*, 339–353.
- Stanton, D. S. Evaluation and Use of BCUT Descriptors in QSAR and QSPR Studies. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 11–20.
- Ertl, P.; Rohde, B.; Selzer, P. Fast Calculation of Molecular Polar Surface Area as a Sum of Fragment-Based Contributions and Its Application to the Prediction of Drug Transport Properties. *J. Med. Chem.* **2000**, *43*, 3714–3717.
- LaBute, P. A widely applicable set of descriptors. *J. Mol. Graphics Modell.* **2000**, *18*, 464–477.
- Visco, D. P., Jr.; Pophale, R. S.; Rintoul, M. D.; Faulon, J.-L. Developing a methodology for an inverse quantitative structure-activity relationship using the signature molecular descriptor. *J. Mol. Graphics Modell.* **2002**, *20*, 429–438.
- Gallegos, A.; Carbo-Dorca, R.; Poncet, R.; Waisser, K. Similarity approach to QSAR. Application to antimycobacterial benzoxazines. *Int. J. Pharm.* **2004**, *269*, 51–60.
- Ehresmann, B.; de Groot, M. J.; Alex, A.; Clark, T. New Molecular Descriptors based on local properties at the molecular surface and a boiling-point model derived from them. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 658–668.
- Clark, T. QSAR and QSPR based solely on surface properties. *J. Mol. Graphics Modell.* **2004**, *22*, 519–525.
- Sun, H. A Universal Molecular Descriptor System for Prediction of LogP, LogS, LogBB, and Absorption. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 748–757.
- Wang, D.-F.; Wiest, O.; Helquist, P.; Lan-Hargest, H.-Y.; Wiech, N. L. QSAR Studies of PC-3 cell line inhibition activity of TSA and SAHA-like hydroxamic acids. *Bioorg. Med. Chem. Lett.* **2004**, *14*, 707–711.
- Liu, X.; Wang, B.; Huang, Z.; Han, S.; Wang, L. Acute toxicity and quantitative structure-activity relationships of α -branched phenylsulfonyl acetates to *Daphnia magna*. *Chemosphere* **2003**, *50*, 403–408.
- Dominy, B. N.; Shakhnovich, E. I. Native Atom Types for Knowledge-Based Potentials: Application to Binding Energy Prediction. *J. Med. Chem.* **2004**, *47*, 4538–4558.
- Burden, F. R.; Winkler, D. A. Robust QSAR Models Using Bayesian Regularized Artificial Neural Networks. *J. Med. Chem.* **1999**, *42*, 3183–3187.
- Serra, J. R.; Jurs, P. C.; Kaiser, K. L. E. Linear Regression and Computational Neural Network Prediction of Tetrahymena Acute Toxicity for Aromatic Compounds from Molecular Structure. *Chem. Res. Toxicol.* **2001**, *14*, 1535–1545. Data courtesy of TerraBase Inc.
- Physical Properties Database PHYSPROP, Syracuse Research Corporation. <http://www.syrres.com/esc/physprop.htm> (accessed March 30, year).
- Polley, M. J.; Winkler, D. A.; Burden, F. R. Broad-based QSAR of farnesyltransferase inhibitors using a Bayesian regularised neural network. *J. Med. Chem.* **2004**, *47*, 6230–6238.
- Rao, S. N.; Stockfisch, T. P. Partially unified multiple property recursive partitioning (PUMP-RP) analyses of cyclooxygenase (COX) inhibitors. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1614–1622.
- Gasteiger, J.; Marsili, M. Iterative Partial Equalization Of Orbital Electronegativity - A Rapid Access To Atomic Charges. *Tetrahedron* **1980**, *36*, 3219–3228.
- Mortier, W. J.; Ghosh, S. K.; Shankar, S. Electronegativity-equalization method for the calculation of atomic charges in molecules. *J. Am. Chem. Soc.* **1986**, *108*, 4315–4320.
- Bultinck, P.; Langenaeker, W.; Lahorte, P.; De Proft, F.; Geerlings, P.; Waroquier, M.; Tollenaere, J. P. The Electronegativity Equalization Method I: Parameterization and validation for atomic charge calculations. *J. Phys. Chem. A* **2002**, *106*, 7887–7894.
- Winkler, D. A.; Burden, F. R.; Watkins, A. J. R. Atomistic topological indices applied to benzodiazepines using various regression methods. *Quant. Struct.-Act. Relat.* **1998**, *17*, 14–19.
- Burden, F. R. A chemically intuitive molecular index based on the eigenvalues of a modified adjacency matrix. *Quant. Struct.-Act. Relat.* **1997**, *16*, 309–314.
- Burden, F. R.; Winkler, D. A. New QSAR methods applied to structure-activity mapping and combinatorial chemistry. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 236–242.
- Winkler, D. A. Neural networks in ADME and toxicity prediction. *Drugs Future* **2004**, *29*, 1043–1057.
- Burden, F. R. Quantitative Structure-Activity Relationship Studies Using Gaussian Processes. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 830–835.
- Burden, F. R.; Ford, M. G.; Whitley, D. C.; Winkler, D. A. Use of Automatic Relevance Determination in QSAR Studies Using Bayesian Neural Networks. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 1423–1430.
- Burden, F. R.; Winkler, D. A. A QSAR Model for the Acute Toxicity of Substituted Benzenes towards *Tetrahymena pyriformis* using Bayesian Regularized Neural Networks. *Chem. Res. Toxicol.* **2000**, *13*, 436–440.
- Faulon, J.-L. The Signature Molecular Descriptor. 1. Using Extended Valence Sequences in QSAR and QSPR Studies. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 707–720.
- Tetko, I. V.; Tanchuk, V. Y.; Villa, A. E. P. 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.
- Molnár, L.; Kesaru, G. M.; Papp, A.; Gulyás, Z.; Darvas, F. A Neural Network Based Prediction of Octanol-Water Partition Coefficients using Atomic5 Fragment Descriptors. *Bioorg. Med. Chem. Lett.* **2004**, *14*, 851–853.

CI800290H