

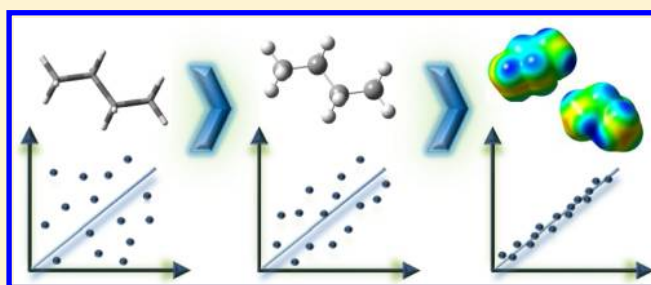
What is Wrong with Quantitative Structure–Property Relations Models Based on Three-Dimensional Descriptors?

M. Hechinger,[†] K. Leonhard,[‡] and W. Marquardt^{*,†}

[†]AVT-Process Systems Engineering and [‡]Chair of Technical Thermodynamics, RWTH Aachen University, 52064 Aachen, Germany

S Supporting Information

ABSTRACT: Quantitative structure–property relations (QSPR) employing descriptors derived from the three-dimensional (3D) molecular structure are frequently applied for property prediction in various fields of research. However, there is no common understanding of the necessary degree of detail to which molecular structure has to be known for reliable descriptor evaluation, but computational methods used vary from simplified molecular mechanics up to rigorous ab initio programs. In order to quantify the yet unknown error due to this heterogeneity, widely used 3D molecular descriptors from diverse fields of application are evaluated for molecular structures computed by different computational methods. The results clearly indicate that the widespread, exclusive use of the most stable molecular conformation as well as too simplistic computational methods yield systematically erroneous descriptor values with misleading information for the inferred structure–property relations. Thus, generating an awareness and understanding of this fundamental problem is considered an important first step to make 3D QSPR a generally accepted property prediction method.



■ INTRODUCTION

Quantitative structure–property relations (QSPR) have become a fundamental tool for property prediction in various scientific fields including chemistry, biology, pharmacology, and chemical engineering. Accordingly, relations between molecular structure and macroscopic quantities have been established in diverse areas ranging from thermophysics,^{1–7} carcinogenicity and toxicity,^{8–10} and catalytic activity¹¹ up to combustion kinetic properties^{12–15} and lubricity¹⁶ of biofuels. Quantitative structure–activity relations (QSAR) are routinely employed in drug design to identify molecules with high binding affinity to receptors in order to maximize biochemical activity.^{17–20} A very recent and comprehensive review on the theory and applications of QSPR, with particular emphasis on thermophysical properties, has been given by Katritzky et al.²¹

Regardless of the diversity of investigated properties, any QSPR and QSAR approach assumes that a macroscopic property of a chemical compound depends on one or more distinct molecular features, called descriptors, which are deduced from the molecular topology or geometry. The modeling methods underlying available tailored algorithms^{22,23} as well as ready-to-use software packages²⁴ work on large descriptor databases to identify the most relevant descriptors to predict a certain property. In the last decades, thousands of descriptors have been developed and tailored to numerous prediction problems.²⁵ They can be distinguished according to the level of detail in the representation of the molecular structure required to evaluate the descriptor. While zero- (0D), one- (1D), and two-dimensional (2D) descriptors rely on constitutional and topological information of molecular

structures only, 3D descriptors can only be evaluated if the full spatial molecular arrangement is known. Although 0D, 1D, and 2D descriptors are routinely applied in property prediction, many properties have been shown to require more detailed 3D information to properly capture the relevant molecular characteristics responsible for the macroscopic observations of interest.^{5,11,12,16,20,26–40} Hence, 3D descriptors seem to be indispensable for a reliable model-based design of chemical products involving complex product specifications. As a consequence, sufficiently accurate 3D representations of the molecular structure are needed to determine the 3D descriptors used in 3D QSPR modeling.

Unfortunately, there is yet no common understanding how accurate molecular structure has to be captured to obtain 3D descriptors at a precision which guarantees reliable property prediction by 3D QSPR. Until today, molecular representations used in the computation of 3D molecular descriptors are generated by means of a multitude of methods including rule-based tables of atom distances and angles^{40,41} as provided by the CORINA software,⁴² geometries optimized by means of molecular mechanics,^{12,32} a suite of semiempirical methods^{5–11,15,17,23,30,35,39} or molecules computed by ab initio quantum chemistry methods.^{1,10,11,16,20,26,31,36–38} While an early discussion⁴³ concludes that semiempirical methods sufficiently well represent the molecular structure to properly predict 3D descriptors, more recent comparisons between different computational methodologies^{44,45} but also between

Received: May 25, 2012

Published: July 9, 2012

variants of the same type of method⁴⁶ indicate significantly diverging descriptor values.

A related problem refers to the number and selection of conformers to be used in 3D QSPR modeling. While researchers in drug design typically consider conformational ensembles and even search for conformers with highest bioactivity,^{34,47,48} there are only few studies on 3D QSPR modeling for thermophysical property prediction which consider more than the most stable conformer for descriptor evaluation. They either neglect conformers with a weight of less than 10% at a single temperature,^{27,28} or they select a fixed number of identified conformers.²⁰ While this seems a reasonable pragmatic approach to reduce the computational effort, the prediction error resulting from these simplifications is unknown.

To the authors' knowledge, none of the published studies consider the effect of degenerate conformers or the influence of the thermodynamic contributions to the free energy of a conformer when computing ensemble averaged descriptor values. Karelson et al.⁴³ correctly state that even ab initio methods have yet only been applied to conformers at an energetic minimum corresponding to 0 K, such that temperature and entropic effects are not covered. Since almost all physical or chemical properties depend on temperature, a reliable prediction may therefore require a more rigorous description of molecular structure accounting for temperature effects.

The nonuniform and often not well documented representation of the spatial molecular structure not only detracts 3D QSPR models from becoming a generally accepted methodology for property prediction, but it also prohibits the comparison of different instances of 3D QSPR models reported in literature for the same property. In order to overcome this unsatisfactory situation, the present contribution systematically investigates the effect of molecular representation on the quality of 3D descriptors. We do not aim to establish a novel QSPR for a particular property but rather want to identify the level of detail in molecular representation which is necessary for reliable 3D QSPR model predictions. To this end, our investigations are deliberately restricted to *n*-alkanes, e.g., the homologous series from methane through *n*-heptane, as they are frequently used for the assessment of QSPR models.^{3,5} The discussion thus explicitly excludes macromolecules usually encountered in polymer, materials, or drug science, since they are too large to be investigated with the methods employed. Moreover, the 3D descriptors used in this work have been computed by the DragonX software,⁴⁹ since it is widely applied in QSPR for physical property predictions. Frequently used representatives from various classes of 3D descriptors are exemplarily evaluated for molecular geometries which have been determined by means of typically used computational methods of varying accuracy. The results clearly illustrate severe shortcomings in the methodologies currently employed throughout the QSPR literature, which do neither cover geometry optimization with sufficient rigor nor treat the temperature dependence of the 3D descriptors properly. The present contribution not only strives for a better understanding of the fundamental problems related to current 3D QSPR modeling, but rather aims at guidelines for an appropriate representation of molecular structure in 3D QSPR modeling in order to turn it into a reliable and accepted methodology for property prediction.

Principles of Boltzmann-Averaged Descriptor Calculation. Like any property derived from the molecular geometry, 3D descriptors have to be averaged over the different conformations in which a molecule exists.⁵⁰ The molecular descriptor D of a compound is given by the Boltzmann average of its conformer descriptors D_i^* ,

$$D(T) = \sum_i^N p_i(T) D_i^* \quad (1)$$

where the temperature-dependent Boltzmann weight $p_i(T)$ of conformer i is given by

$$p_i(T) = \frac{g_i \exp^{-\Delta A_i(T)/RT}}{\sum_i^N g_i \exp^{-\Delta A_i(T)/RT}} \quad (2)$$

Here, R refers to the universal gas constant, g_i to the degeneracy, and $\Delta A_i(T)$ to the relative free energy

$$\Delta A_i(T) = A_i(T) - \min_i A_i(T) \quad (3)$$

of conformer i . According to eq 3, the relative free energy $\Delta A_i(T)$ of conformer i equals the difference in total free energy between conformer i and the energetically lowest and hence most stable conformer of the same compound.

The total free energy $A_i(T)$ of conformer i is composed of the electronic energy contribution E_i^{elec} and a temperature-dependent, thermodynamic contribution to free energy A_i^{thermo} , which stems from statistical thermodynamics rather than from quantum mechanics:

$$A_i(T) = E_i^{\text{elec}} + A_i^{\text{thermo}}(T) \quad (4)$$

The electronic energy E_i^{elec} results from the conformers' electron configuration, which is obtained from a solution of the Schrödinger equation for the molecule in its equilibrium geometry, that in turn results from a geometry optimization of the conformer. It should be emphasized that only the electronic ground state energies have been considered for the *n*-alkanes investigated here, since the excited states are known to be much too high to contribute significantly. The thermodynamic contribution A_i^{thermo} consists of the translational free energy A_i^{trans} , the rotational free energy A_i^{rot} and the vibrational free energy A_i^{vib} :

$$A_i^{\text{thermo}}(T) = A_i^{\text{trans}} + A_i^{\text{rot}} + A_i^{\text{vib}} \quad (5)$$

The translational free energy A_i^{trans} depends only on the molecular weight and is therefore directly known, while the rotational free energy A_i^{rot} is based on molecular geometry and is thus available from geometry optimization. The vibrational free energy A_i^{vib} , however, requires an additional frequency calculation of the optimized geometry.

In addition to the relative free energies $\Delta A_i(T)$, the degeneracies g_i of the conformers need to be determined as well in order to evaluate eq 2. The concept of degeneracy is rarely discussed and usually neglected in 3D QSPR modeling. Degenerate conformers exhibit identical energies but different geometries. Figure 1 shows the conformers of *n*-butane as a simple illustrative example. A more detailed discussion on the correct identification of degeneracies is given by Fernández-Ramos et al.⁵¹ and Gilson and Irikura.⁵²

The determination of the Boltzmann-averaged molecular descriptors requires various quantities which partly depend on each other. The resulting computational workflow is illustrated

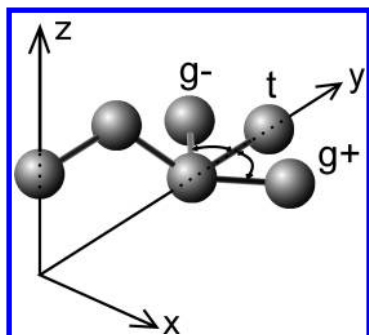


Figure 1. Degeneracy of *n*-butane conformers: The trans-butane conformer (t) is 1-fold degenerate, while the gauche conformers (g_{\pm}) are mirror images at the z - y -plane with an identical, but higher energy than the trans conformer. Hence, only one (g_{+} or g_{-}) is optimized, but it has to enter the Boltzmann distribution with a double weight, i.e., a 2-fold degeneracy ($g_1 = 1$, $g_2 = 2$ in eq 2).

in Figure 2. The first step is the identification of molecular conformations, which results in the type and number of conformers as well as their degeneracies g_i . The initially coarse geometries are then refined in a subsequent geometry optimization, which yields the conformers' geometries as well as the electronic, rotational, and translational energies. From these geometries, both, the molecular descriptors D_i^* and the vibrational frequencies of the identified conformers can be computed. With the vibrational free energies directly following from the vibrational frequencies, all quantities required to evaluate the Boltzmann-averaged descriptors from eqs 1–5 are available.

The various tasks of this workflow can be performed on different levels of rigor and hence accuracy. The following section presents the different computational approaches investigated within this work in order to identify suitable computational methods for the computation of reliable 3D descriptors.

Computational Methods for Geometry and Free Energy Calculation. The geometries and free energies of the homologous series of *n*-alkanes from methane to *n*-heptane have been determined with three different computational methods as summarized in Table 1. Up to five different types of calculations have been performed with three different methods, namely a molecular mechanics (MM), a semiempirical (SE), and an ab initio (AI) method. While the first two methods have been chosen since they are typically employed in QSPR and QSAR modeling across all disciplines, the AI methods are computationally expensive, but fairly reliable, such that very

accurate results can thus be used as a reference for comparison of the approximate methods.

For each computational method, a systematic search for conformers has initially been conducted with Spartan10.⁵⁴ While this approach is generally restricted to molecular structures with a limited number of degrees of freedom, such as those of the *n*-alkanes investigated here, it ensures a complete representation of the molecules' spatial configurations. Spartan10 is also able to compute the degeneracy g_i for each of the conformers, which—due to the systematic search—are identical for each of the computational methods MM, SE, and AI. For the conformational search by MM, the Merck molecular force field (MMFF94)⁵⁶ has been employed, whereas AM1⁵⁷ has been used to identify the conformers for the SE and AI methods. While a more rigorous theory could also have been employed for conformer search in case of the AI calculations, the AM1 geometries were considered to be of sufficient accuracy for the subsequent steps.

The identified conformational geometries have been further optimized with Gaussian09,⁵³ which, however, does not implement the MMFF94 force field. Since the consistent use of Gaussian09 is desirable for all methods (e.g., MM, SE, and AI) and since optimization and frequency calculation have to be performed with identical methods, the initial MMFF94 geometries computed with Spartan10 have been reoptimized with the Dreiding⁵⁸ force field in Gaussian09. Likewise, although the conformer geometries have already been determined by Spartan10 using AM1, a further geometry optimization based on AM1 has been performed in Gaussian09 to ensure consistency with the subsequent frequency calculation. While the refinement has been done by AM1 for the SE calculations, B3LYP^{59,60} has been used with the TZVP basis set for the AI method. This method has proven to yield reliable geometries which are an excellent starting point for the computation of single point energies by computationally more expensive methods.⁶¹ The B3LYP/TZVP calculations have been performed on an ultrafine grid with tight convergence criteria for all conformers.

For the AI method (cf. Table 1), a single-point energy calculation with the increased aug-cc-pVTZ basis set was performed on the optimized geometry, since the additional diffuse basis functions allow for a better representation of molecular electrostatics. Moreover, point charges have been fitted to the electrostatic potential $P^{\text{electrostatic}}$ according to the Chelp scheme.⁶² These charges are required to evaluate atomic charge-dependent descriptors. Although the electrostatic potential can also be derived by SE⁶³ and MM⁵⁶ methods, it has only been computed in the AI calculations, since common MM and SE approaches in 3D QSPR modeling generally omit

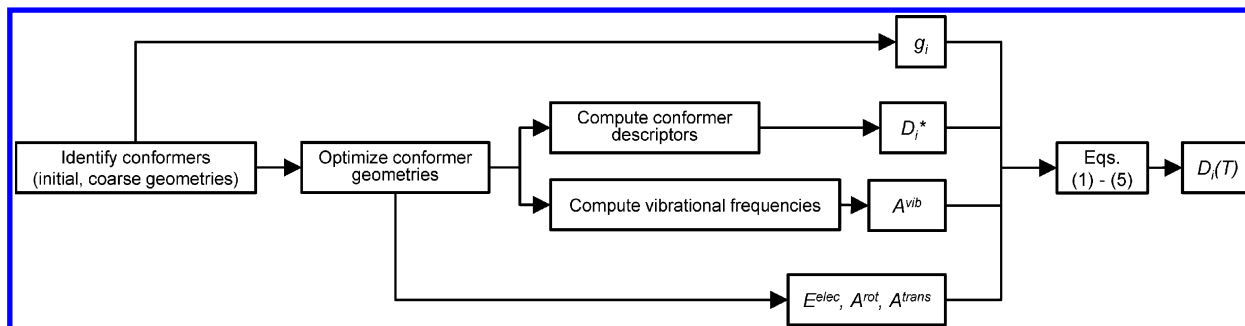


Figure 2. Workflow of Boltzmann-averaged descriptor evaluation.

Table 1. Performed Calculations and Applied Methods for Geometry and Energy Determination

	molecular mechanics (MM)	semiempirical (SE)	ab initio (AI)	calculated property
conformer search (Spartan10 ⁵⁴)	systematic/MMFF94	systematic/AM1	systematic/AM1	g_i
geometry optimization (Gaussian09 ⁵³)	DREIDING	AM1	B3LYP/TZVP	$E^{\text{elec}}, A^{\text{rot}}, A^{\text{trans}}$
single-point/electrostatics (Gaussian09 ⁵³)	na ^a	na ^a	B3LYP/aug-cc-pVTZ	$p^{\text{electrostatic}}$
frequency calculation (Gaussian09 ⁵³)	DREIDING	AM1	B3LYP/TZVP	A^{vib}
refined electronic energy (Turbomole ⁵⁵)	na ^a	na ^a	RI-MP2/aug-cc-pVTZ	$A^{\text{elec,refined}}$

^aIdentical to geometry optimization step.

Table 2. Computational Time (s) for the Evaluation of *n*-Alkanes by Means of Molecular Mechanics (MM), Semiempirical (SE), and Ab initio (AI) Methods as Specified in Table 1^a

	<i>n</i> -butane			<i>n</i> -pentane			<i>n</i> -hexane			<i>n</i> -heptane		
	MM	SE	AI	MM	SE	AI	MM	SE	AI	MM	SE	AI
conformer search	2	27	27	4	55	55	6	109	109	12	294	294
geometry optimization	2	4	822	3	17	5047	7	50	26640	17	185	137074
single point			1043			3330			7067			11734
frequency calculation	1	2	520	2	4	2291	6	10	10838	14	23	50964
refined energy			106			217			384			651
total	5	33	2518	9	76	10940	19	169	45038	33	502	200717

^aComputations have been performed on a single Intel Xeon CPU, 3.06 GHz, 24 GB RAM.

the use of an electrostatic potential and employ the computationally less demanding partial charges derived from Mullikens' population analysis.⁶⁴ Hence, in order to demonstrate the effect of using charges derived from the electrostatic potential, Mulliken charges are applied in MM and SE calculations. It should, however, be emphasized that several charge fitting schemes besides Chelp exist, which were shown to produce slightly different results.⁶⁵

The determination of the vibrational contribution to the overall free energy requires a frequency calculation, which has been performed for all methods (i.e., for MM, SE and AI). In all three cases, the same method used for geometry optimization has been employed to obtain meaningful results. The total thermodynamic contribution $A^{\text{thermo}}(T)$ (cf. eq 3) has then been evaluated between 0 and 1000 K in steps of 5 K for each conformer. For this purpose, the freqchk utility of Gaussian09 has been used, and the data have been tabulated for later interpolation.

Finally, the electronic energy E^{elec} computed in the AI route has been refined by an RI-MP2⁶⁶ calculation. Here the RI approximation⁶⁷ as implemented in Turbomole⁵⁵ has been used. While reducing the computational effort by a factor of up to 170 compared to standard MP2, the RI approximation has shown to provide accurate approximations of the normal MP2 energies.⁶⁷ Again, performing an energy computation by means of MP2 is inappropriate for the MM and SE methods and thus has been omitted. Accordingly, the energies resulting from geometry optimization have been used in the MM and SE calculations to reflect the typically employed approach in current QSPR and QSAR modeling.

Evaluation of Frequently Used 3D Descriptors. On the basis of the molecular information obtained from the computations described in the previous section, 3D molecular descriptors have been determined using DragonX⁴⁹ for the homologous series from methane through *n*-heptane. DragonX computes a total of 735 molecular descriptors from 6 different classes which depend on the 3D molecular structure or atomic charge. Due to the large variety of 3D descriptors provided by DragonX, this software has been employed in various fields of research including chromatography,⁶⁸ spectroscopy,^{69,70} and

catalysis⁷¹ but, in particular, also for predicting important thermophysical properties^{4,5,29,30,41,72–75} as they are frequently needed in chemical engineering applications and beyond. Still, the covered range of descriptors does obviously not properly reflect all physical properties of interest. For instance CodessaPro²⁴ provides a more diverse selection of electrostatic, molecular orbital-related, or quantum chemical descriptors compared to the 3D descriptors provided by DragonX. Furthermore, there even exist a number of other software packages to compute targeted molecular descriptors for specific applications. Despite this large number of different data sources, we have limited our attention to the broad range of descriptors provided by DragonX in this work for simplicity, since a representative study on the effects of molecular representation on the reliability of 3D descriptors for physical property prediction can still be carried out this way. The entire range of descriptors provided by DragonX has thus been evaluated, while commonly used 3D descriptors from different descriptor classes reported in QSPR modeling across the disciplines have been selected for a more detailed study. However, since the set of descriptors can be considered as representative, the conclusions of this work not only cover the types of descriptors investigated here. Rather, the methodological findings apply to any kind of 3D descriptor, regardless of its origin or molecular property it is supposed to describe.

RESULTS AND DISCUSSION

Geometry optimizations and energy calculations have been performed for all the *n*-alkanes from methane to *n*-heptane. However, since the smaller molecules exist in a single conformer only, the results are primarily discussed for *n*-butane to *n*-heptane. A detailed list of the optimized conformer energies and degeneracies is compiled in the Supporting Information. After a brief discussion of the computational aspects of the performed calculations, the general consequences of considering degeneracy and the thermodynamic contribution to free energy in the computation of Boltzmann-averaged properties are outlined before these effects are demonstrated for several widely employed 3D descriptors.

Computational Aspects. Spartan10 found too large values for the degeneracies g_i for some higher-energy conformations since the effect of topologically equivalent atoms on the symmetry of the nuclear wave function⁵² is not accounted for correctly. While this deficiency might reduce the accuracy of the Boltzmann-averaged descriptors at higher temperatures, it does not affect our general conclusions that the conformational spectrum has to be considered to properly capture the most significant effects of conformational range and degeneracy. Table 2 shows the time required for the computation of each task as well as the total time needed to compute all conformers by means of each method. The systematic conformer search dominates the computational effort for the MM and SE methods, while geometry optimization has the largest share (except for *n*-butane) in case of the AI method. Moreover, the SE method is about 10 times more expensive than the MM method, but still every molecule can be evaluated within a few minutes. In contrast, the total computational time required by the AI method for the computation of *n*-heptane amounts to about 56 h, e.g. 400 times the effort needed for the SE calculations. If the parallel mode is exploited on a multicore machine, the total time can be reduced to several minutes for *n*-butane and to several hours for *n*-heptane, respectively.

Effect of Degeneracy and Thermodynamic Contribution to Free Energy on Boltzmann Distributions. Before investigating the computational results for particular descriptors, the consequences of a simplified evaluation of Boltzmann-distributed quantities is briefly discussed with reference to eqs 2–4. As already mentioned above, the very few papers on QSPR modeling which consider conformers at all, neither account for degeneracies g_i nor for the thermodynamic contribution A_i^{thermo} when determining the Boltzmann weights of the conformers. If degeneracy is not considered (i.e., $g_i = 1 \forall i$), only a single conformer is accounted for for each energy. In case the thermodynamic contribution is neglected (i.e., $A_i = E_i^{\text{elec}}$), the investigation is restricted to 0 K. Figure 3 shows the

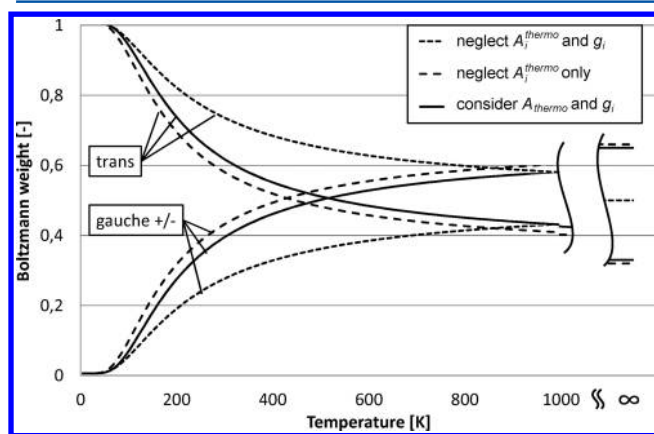


Figure 3. Boltzmann weights of *n*-butane conformers (cf. Figure 1) as a function of temperature, evaluated by the AI method; effect of thermodynamic contribution A_i^{thermo} and degeneracy g_i , cf. eq 2.

Boltzmann weights of *n*-butane over a wide temperature range to illustrate the effect of different degrees of simplification. While the energetically lower *trans* conformation of *n*-butane ($g_1 = 1$, cf. Figure 1) dominates at lower temperatures, the energetically higher *gauche* \pm conformations ($g_2 = 2$) become more significant with increasing temperature. However, incorrect degeneracies (i.e., $g_1 = g_2 = 1$) cause the *gauche*

conformation to become important only at much higher temperatures to even result in a conformer distribution at 1000 K which is opposite to the one considering correct degeneracies. This qualitative error is maintained up to infinite temperature, where the Boltzmann weights reach the asymptotic value of

$$p_i(T \rightarrow +\infty) = \frac{g_i}{\sum_i^N g_i} \quad (6)$$

which is $p_1 = p_2 = 0.5$ for incorrect and $p_1 = 1/3$ and $p_2 = 2/3$ for correct degeneracies, respectively.

If the thermodynamic contribution A_i^{thermo} to the free energy is accounted for in addition, a smaller but still considerable change in the conformer distribution can be observed. Although *n*-butane might show an unusually strong sensitivity to model simplification, the results clearly indicate the need to consider at least degeneracy when computing any Boltzmann-averaged quantity. Moreover, the low additional computational cost suggests to also include a systematic conformer search (cf. Table 2).

Effect of the Computational Method on the Boltzmann Distribution. Even if degeneracy and the thermodynamic contribution to the free energy are properly accounted for, the Boltzmann weights p_i may differ significantly if the computational method applied to determine the free energies in eqs 2–5 is varied. To illustrate this influence quantitatively, the Boltzmann distributions of *n*-butane and *n*-heptane conformers have been evaluated by the three different methods introduced above (cf. Table 1). The thermodynamic contribution A_i^{thermo} as well as the correct degeneracy g_i have been considered in all cases to compute the results given in Figure 4. The computational method obviously influences the Boltzmann distribution as shown in Figure 4a for *n*-butane. The results obtained by the MM method strongly differ from those determined using the SE and AI methods. Though the distribution determined by the SE method is close to the one determined by the AI method for *n*-butane, this does not hold for the two most important low-temperature conformers of *n*-heptane as shown in Figure 4b. The contribution of the most stable conformer at 0 K declines much slower for the MM and SE methods than for the more accurate AI method. Again, the SE method yields results which are closer to those of the AI method, but the relative deviation between SE and AI data still amounts up to 125% at the melting point T_m of *n*-heptane as shown in Figure 4b.

As stated above, many authors only account for the most stable conformer to determine the electronic energy E^{elec} , which in case of *n*-alkanes is always the straight chain configuration. However, the consideration of a temperature-dependent thermodynamic contribution $A^{\text{thermo}}(T)$ results in a temperature-dependent free energy (cf. eq 4), such that the lowest free energy conformer might change within the relevant temperature range. This is in fact the case for *n*-pentane, *n*-hexane, and *n*-heptane. Table 3 shows the identified transitions between most stable conformers exemplarily for *n*-heptane. The MM method does not detect any transition from the most stable straight chain structure (conformer number 1) for *n*-heptane in the entire temperature range from 0–1000 K. On the basis of the SE method, the free energy of conformer 2 (cf. the Supporting Information) falls below the value of the straight chain conformer 1 at around 605 K, whereas the AI method determines conformer 3 to become the most stable

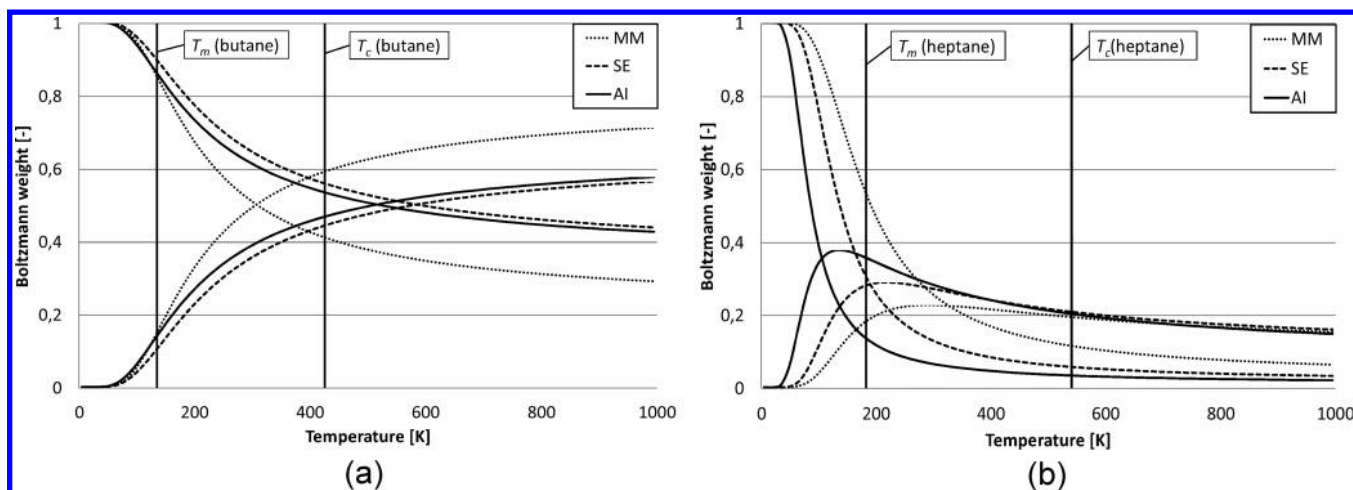


Figure 4. Comparison of Boltzmann distributions of *n*-butane (a) and *n*-heptane (b) for MM, SE, and AI methods (cf. Table 1). For *n*-heptane only the two most relevant low-temperature conformers are displayed for the sake of clarity. The melting point T_m and critical point T_c indicate the most relevant temperature range for QSPR and QSAR modeling.

Table 3. Detected Transitions and Associated Temperatures of the Most Stable *n*-Heptane Conformers^a

method	temperature [K]	lowest free energy conformer
MM	—	—
SE	605	1 → 2
AI	275	1 → 3

^aConformer numbers are ranked according to the electronic energies E^{elec} , with conformer 1 having the lowest value, cf. the Supporting Information.

configuration at around 275 K. Again, since the AI energies are considered most reliable, significant errors of the computationally less expensive methods have to be conjectured. As an erroneous detection of the most stable conformation affects the relative energies and Boltzmann weights of all other conformers (cf. eq 3), the deviations from true energies might lead to strongly deteriorated Boltzmann distributions.

Effect of the Computational Method on Selected 3D Descriptors. While the variations in the Boltzmann distribution discussed in the previous section obviously affect any 3D descriptor, a selection of frequently employed 3D descriptors

from various classes are exemplarily evaluated next. Since the chosen descriptors are representative, the conclusions are likewise applicable to other 3D descriptors and QSPR models beyond those discussed here. Again, the presentation is exemplarily restricted to *n*-pentane and *n*-hexane. Additional results for the remaining *n*-alkanes are given in the Supporting Information. All results have been derived by taking degeneracy and the thermodynamic contribution to free energy properly into account.

Figure 5a shows the RDF045u descriptor from the class of radial distribution functions (RDF),²⁵ which has previously been employed to model the aqueous solubility of organic compounds⁷⁶ and drug-like molecules.³³ Only the most stable conformer, either determined from rule-based tables⁷⁶ or by an SE-based geometry optimization,³³ has been considered in previous work to result in the descriptor value at 0 K in Figure 5a. However, the strong temperature dependence of the RDF045u descriptor clearly illustrates that exclusively considering the most stable conformer is only valid at very low temperatures. Since most QSPR and QSAR models are developed and employed at higher temperatures, considerable

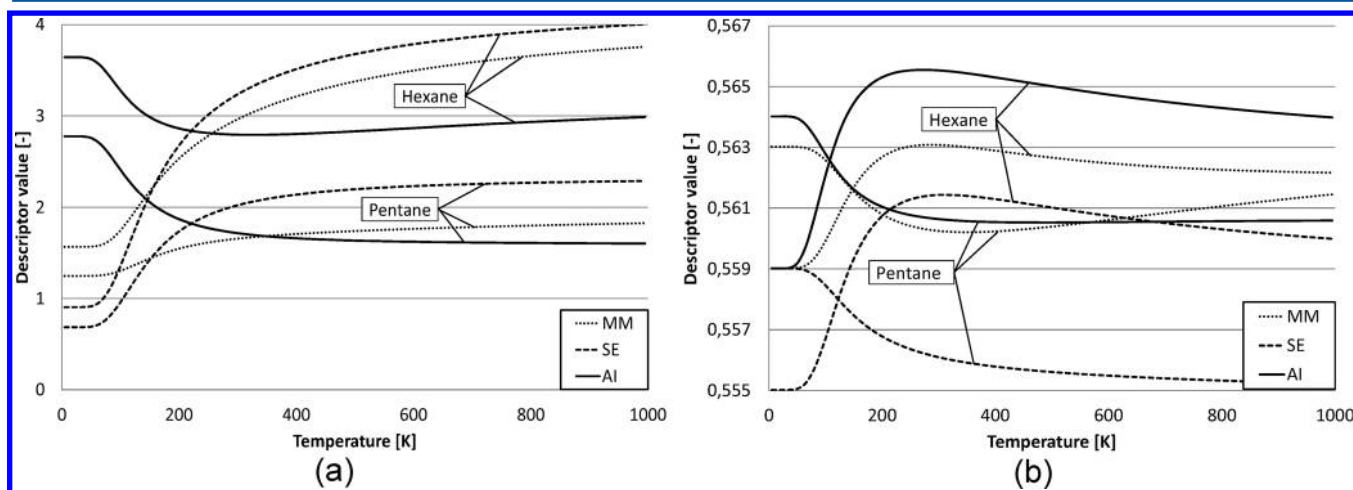


Figure 5. RDF descriptor RDF045u (a) and WHIM descriptor E1u (b) as a function of temperature, evaluated for *n*-pentane and *n*-hexane by the MM, SE, and AI methods.

prediction errors have to be expected. In addition, the descriptor values derived from the geometries computed by means of the MM and SE methods increase with increasing temperature, while the values obtained by the most reliable AI method decrease. Accordingly, strongly deteriorated descriptor values are likely to be obtained in case the computationally less demanding MM and SE methods are used.

The descriptor E1u of the WHIM class²⁵ investigated in Figure 5b has been employed in various models to predict phase equilibrium,⁷⁵ estrogen activity,⁷⁷ or the inhibition of aldose reductase for flavonoids.⁷⁸ Similar WHIM descriptors have been applied to model normal boiling points³⁰ or Hildebrandt solubility parameters⁷⁴ lately. In all cases only the most stable conformer obtained by either the MM or the SE method has been considered. This modeling approach is likely to be too simplistic given the strong temperature dependence of the WHIM descriptor as well as the difference in the absolute values computed by the MM and SE methods compared to the most reliable AI method as indicated in Figure 5b.

Further consequences of an insufficient representation of molecular structure are illustrated exemplarily by the charge-dependent descriptor qnmax,²⁵ which has for instance been used to model octanol–water partition coefficients³⁶ or retention indices of oils in gas chromatography.⁷⁹ As shown in Figure 6, the descriptor qnmax not only deviates by a factor

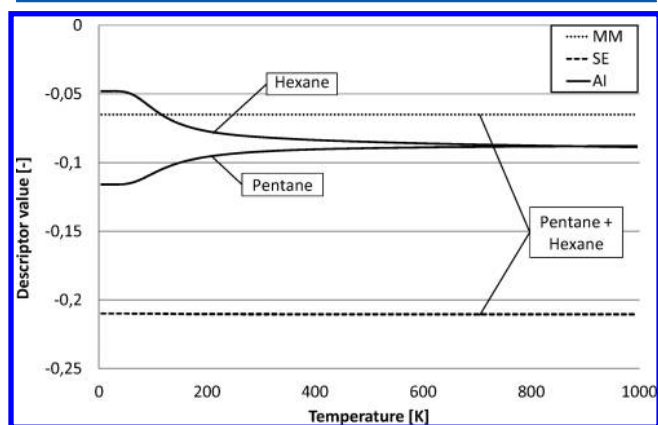


Figure 6. Maximum negative charge qnmax as a function of temperature, evaluated for *n*-pentane and *n*-hexane by the MM, SE, and AI methods.

of about 4 when evaluated by either the MM or the SE method. It is also unable to distinguish between *n*-pentane and *n*-hexane, which is due to the use of Mulliken charges, which are independent of the conformers' spatial arrangement. In contrast, the more reliable AI method resolves the conformational differences and thus provides a more rigorous description of the molecules' electrostatic properties.

Finally the temperature dependency of the DispV descriptor, a typical representative of the class of geometrical descriptors, is shown in Figure 7. It has recently been utilized to model melting points of various compounds by a 3D QSPR with descriptors computed by SE methods.⁷³ Since unsatisfactory predictions have been obtained by the DispV and several other descriptors, the authors conclude that the poor results might be due to deficiencies in the descriptor computation. The strong temperature dependence as well as the significant deviation resulting from the application of different computational

methods (cf. Figure 7) renders the use of SE methods as well as the temperature-independent descriptors likely reasons for the inadequate results. In addition, the location of the intersection of the *n*-pentane and *n*-hexane descriptors varies considerably for the MM, SE, and AI methods. Since such intersections result in a temperature-dependent ordering of the descriptors, the quality of the descriptors strongly affects their selection and hence the quality of any correlation-based QSPR modeling approach. This becomes even more obvious from Figure 7b, where the DispV values are evaluated for all *n*-alkanes ranging from methane (C₁) to *n*-heptane (C₇). While the descriptor values are almost independent of the computational method for those molecules existing in only a single conformer, the deviations increase with the chain length and reach a relative error of 40% for the SE method and of 70% for the MM method in case of *n*-heptane. These results are compared to those commonly obtained if the most stable conformer is based on a geometry computed by an SE method in Figure 7b. Whereas Boltzmann averaging results in a continuous descriptor behavior along the homologous series, the DispV values alternate between zero and nonzero if only the most stable straight chain conformer is considered.

These results demonstrate that the computational method used for the calculation of the geometrical representation of the molecule severely influences the quality of a QSPR model. Likewise, the consideration of only single conformers might lead to erroneous predictions based on structure–property relations. Since all QSPR and QSAR model identification methods correlate macroscopic properties with descriptor values,²¹ the choice of modeling methodology strongly affects the descriptor selection process as illustrated by the results displayed in Figures 5–7. Simplified models might result in nonphysical relationships which are likely not predictive and thus cannot serve for a reliable product design. Although the sensitivity of other 3D descriptors might be weaker than the sensitivity of the widely employed descriptors selected for illustration in this work, the possible existence of just a single sensitive descriptor in the data set strongly suggests a more rigorous representation of molecular structure in order to exploit the descriptors' true information content.

SUMMARY AND CONCLUSIONS

The effect of the type of molecular representation on the quality of 3D molecular descriptors frequently used in quantitative structure–property relations (QSPR) and quantitative structure–activity relations (QSAR) has been investigated. In particular, a systematic conformer search has been performed for a homologous series of *n*-alkanes ranging from methane to *n*-heptane to properly consider degeneracy. Optimized geometries and free energies of the identified conformers have been determined by means of three different and widely used computational methods including molecular mechanics as well as semiempirical and ab initio methods. The resulting Boltzmann distributions were used to evaluate commonly employed 3D descriptors over a wide temperature range.

The results clearly indicate that both the computational methods used as well as the way conformers and their degeneracies are treated have a major influence on the quality of the computed descriptor values. In contrast to the general assumption of a single most stable conformer, the Boltzmann-distributed descriptor evaluation introduces a temperature dependence of the 3D descriptors, which was shown to be

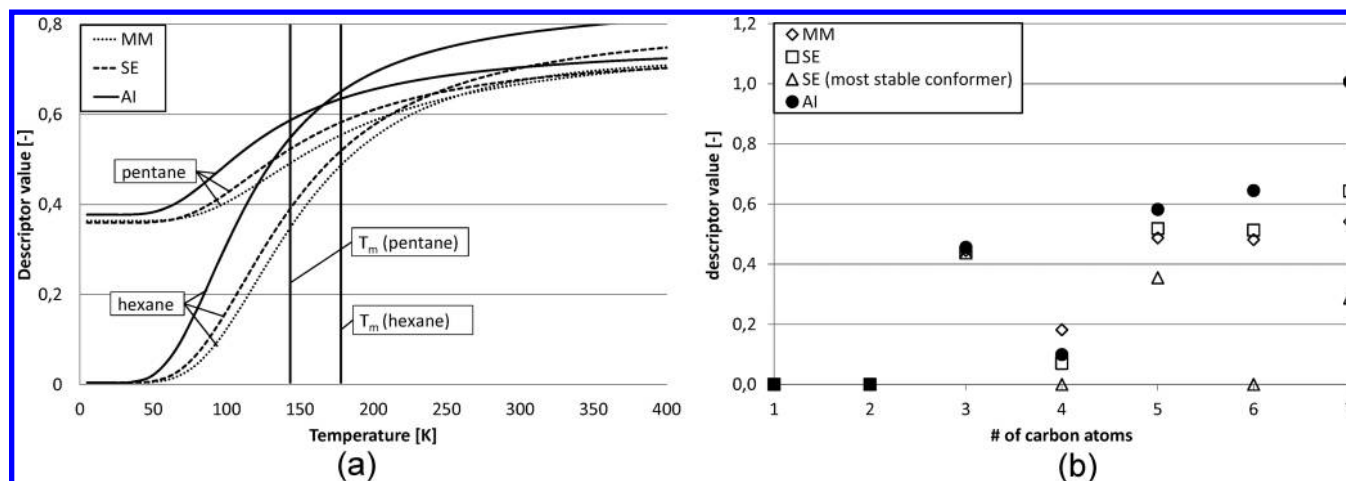


Figure 7. Geometrical descriptor DispV as a function of temperature (a), evaluated for *n*-pentane and *n*-hexane by the MM, SE, and AI methods, and as a function of the number of C atoms in the homologous series of alkanes ranging from methane (C_1) to *n*-heptane (C_7) at the melting point T_m of each compound.

significant for several representative and widely employed molecular descriptors. Accordingly, the widespread use of single conformers, which are optimized by either a molecular mechanics or a semiempirical method, runs at risk of not exploiting the real descriptors' information content. Thus, nonphysical structure–property relations with no or little predictive power may often result.

The investigations have only been performed for ideal gas phase energies in this work. A future extension to liquid phase energies requires the consideration of molecular interactions. Since such interactions can be quantified with available methods,^{80,81} the procedure described here can be regarded as a first step toward a more thorough description of molecular structure.

Moreover, the temperature dependence of 3D descriptors implies methodological changes in QSPR and QSAR modeling if a characteristic temperature such as the melting, boiling or critical temperature were predicted: Since such characteristic temperatures of a novel compound are unknown, but its 3D descriptors have to be evaluated at this very same temperature, iterative 3D QSPR and 3D QSAR modeling techniques have to be developed.

The high computational effort for structure optimization and energy calculation still restricts the applicability of the suggested approach to molecules of moderate size. Still, the suggested procedure can easily be applied to a wide range of industrially relevant molecules using established software tools. The computational procedure can be streamlined since the properties of any chemical substance have to be computed only once to be made available for later use. This way a reliable database of molecular properties can be established as it is common practice for thermophysical, reaction kinetic, and other data.⁸² Such an approach would lead to a transparent procedure for 3D QSPR modeling which constitutes a major stepping stone toward reliable model-based product design across the various domains of application.

■ ASSOCIATED CONTENT

● Supporting Information

Excel spreadsheet containing tables with electronic energies and degeneracies of the identified conformers, as well as computed descriptor data; structure files for all compounds and their

conformers discussed in this paper. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*E-mail address: wolfgang.marquardt@avt.rwth-aachen.de.
Tel.: +49-241-80-94668. Fax: +49-241-80-92326.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

This work was performed as part of the Cluster of Excellence "Tailor-Made Fuels from Biomass", which is funded by the Excellence Initiative by the German federal and state governments to promote science and research at German universities.

■ REFERENCES

- (1) Herndon, W.; Biedermann, P.; Agranat, I. Molecular structure parameters and predictions of enthalpies of formation for catacondensed and pericondensed polycyclic aromatic hydrocarbons. *J. Org. Chem.* **1998**, *63*, 7445–7448.
- (2) Zefirov, N.; Palyulin, V.; Oliferenko, A.; Ivanova, A.; Ivanov, A. Method for the Construction of Universal Structure-Property Relationship Models using the Example of Normal Boiling Temperature for a Wide Set of Organic Compounds. *Dokl. Chem.* **2001**, *381*, 356–358.
- (3) Brauner, N.; Shacham, M.; St.; Cholakov, G.; Stateva, R. Property prediction by similarity of molecular structures - practical application and consistency analysis. *Chem. Eng. Sci.* **2005**, *60*, 5458–5471.
- (4) Vatani, A.; Mehrpooya, M.; Gharagheizi, F. Prediction of standard enthalpy of formation by a QSPR model. *Int. J. Mol. Sci.* **2007**, *8*, 407–432.
- (5) Cholakov, G. S.; Stateva, R.; Brauner, N.; Shacham, M. Estimation of Properties of Homologous Series with Targeted Quantitative Structure-Property Relationships. *J. Chem. Eng. Data* **2008**, *53*, 2510–2520.
- (6) Shacham, M.; Cholakov, G.; Stateva, R.; Brauner, N. Quantitative Structure-Property Relationships for Prediction of Phase Equilibrium Related Properties. *Ind. Eng. Chem. Res.* **2010**, *49*, 900–912.
- (7) Katritzky, A.; Stoyanova-Slavova, I.; Tamm, K.; Tamm, T.; Karelson, M. Application of the QSPR Approach to the Boiling Points of Azeotropes. *J. Phys. Chem. A* **2011**, *115*, 3475–3479.

- (8) Sixt, S.; Altschuh, J.; Brüggemann, R. Quantitative structure-toxicity relationships for 80 chlorinated compounds using quantum chemical descriptors. *Chemosphere* **1995**, *30*, 2397–2414.
- (9) Helguera, A.; Cordeiro, M.; Natalia, D.; Perez, M.; Combes, R.; Gonzalez, M. QSAR modeling of the rodent carcinogenicity of nitrocompounds. *Bioorg. Med. Chem.* **2008**, *16*, 3395–3407.
- (10) Shamovsky, I.; Ripa, L.; Börjesson, L.; Mee, C. D.; Norden, B.; Hansen, P.; Hasselgren, C.; O'Donovan, M.; Sjö, P. Explanation for main features of structure-genotoxicity relationships of aromatic amines by theoretical studies of their activation pathways in CYP1A2. *J. Am. Chem. Soc.* **2011**, *133*, 16168–16185.
- (11) Occhipinti, G.; Bjørsvik, H.; Jensen, V. Quantitative Structure-Activity Relationships of Ruthenium Catalysts for Olefin Metathesis. *J. Am. Chem. Soc.* **2006**, *128*, 6952–6964.
- (12) Creton, B.; Dartiguelongue, C.; de Bruin, T.; Toulhoat, H. Prediction of the Cetane Number of Diesel Compounds Using the Quantitative Structure Property Relationship. *Energy Fuels* **2010**, *24*, 5396–5403.
- (13) Pan, Y.; Jiang, J.; Ding, X.; Wang, R.; Jiang, J. Prediction of Flammability Characteristics of Pure Hydrocarbons from Molecular Structure. *AIChE J.* **2010**, *56*, 690–701.
- (14) Hechinger, M.; Marquardt, W. Targeted QSPR for the prediction of the laminar burning velocity of biofuels. *Comput. Chem. Eng.* **2010**, *34*, 1507–1514.
- (15) Katritzky, A.; Fara, D. How chemical structure determines physical, chemical, and technological properties: An overview illustrating the potential of quantitative structure-property relationships for fuels science. *Energy Fuels* **2005**, *19*, 922–935.
- (16) Masuch, K.; Fatemi, A.; Murrenhoff, H.; Leonhard, K. A COSMO-RS based QSPR model for the lubricity of biodiesel and petrodiesel components. *Lubrication* **2011**, *23*, 249–262.
- (17) Kubinyi, H. QSAR and 3D QSAR in drug design Part I: methodology. *Drug Discovery Today* **1997**, *2*, 457–467.
- (18) Jonsdottir, S.; Jorgensen, F.; Brunak, S. Prediction methods and databases within chemoinformatics: emphasis on drugs and drug candidates. *Bioinformatics* **2005**, *21*, 2145–2160.
- (19) Perkins, R.; Fang, H.; Tong, W.; Welsh, W. Quantitative structure-activity relationship methods: Perspectives on drug discovery and toxicology. *Environ. Toxicol. Chem.* **2003**, *22*, 1666–1679.
- (20) Devereux, M.; Popelier, P.; McLay, I. Quantum Isostere Database: A Web-Based Tool Using Quantum Chemical Topology To Predict Bioisosteric Replacements for Drug Design. *J. Chem. Inf. Model.* **2009**, *49*, 1497–1513.
- (21) Katritzky, A.; Kuanar, M.; Slavov, S.; Hall, C. Quantitative Correlation of Physical and Chemical Properties with Chemical Structure: Utility for Prediction. *Chem. Rev.* **2010**, *110*, 5714–5789.
- (22) Shacham, M.; Brauner, N. The SROV program for data analysis and regression model identification. *Comput. Chem. Eng.* **2003**, *27*, 701–714.
- (23) Wold, S.; Sjöström, M.; Eriksson, L. PLS-regression: a basic tool of chemometrics. *Chemom. Intell. Lab. Syst.* **2001**, *58*, 109–130.
- (24) Katritzky, A.; Karelson, M.; Petrukhin, R. CODESSA Pro; CompuDrug International Inc.: Sedona, AZ, 2012.
- (25) Todeschini, R.; Consonni, V. *Molecular Descriptors for Chemoinformatics*, second ed.; Wiley-VCH: Weinheim, 2009.
- (26) Klamt, A.; Eckert, F.; Hornig, M.; Beck, M. E.; Bürger, T. Prediction of aqueous solubility of drugs and pesticides with COSMO-RS. *J. Comput. Chem.* **2002**, *23*, 275–281.
- (27) Dyekjaer, J.; Rasmussen, K.; Jonsdottir, S. QSPR models based on molecular mechanics and quantum chemical calculations. 1. Construction of Boltzmann-averaged descriptors for alkanes, alcohols, diols, ethers and cyclic compounds. *J. Mol. Model.* **2002**, *8*, 277–289.
- (28) Dyekjaer, J.; Jonsdottir, S. QSPR Models Based on Molecular Mechanics and Quantum Chemical Calculations. 2. Thermodynamic Properties of Alkanes, Alcohols, Polyols, and Ethers. *Ind. Eng. Chem. Res.* **2003**, *42*, 4241–4259.
- (29) Duchowicz, P. R.; Castro, E. A.; Fernandez, F. M.; Gonzalez, M. P. A new search algorithm for QSPR/QSAR theories: Normal boiling points of some organic molecules. *Chem. Phys. Lett.* **2005**, *412*, 376–380.
- (30) Duchowicz, P. R.; Fernandez, M.; Caballero, J.; Castro, E. A.; Fernandez, F. M. QSAR for non-nucleoside inhibitors of HIV-1 reverse transcriptase. *Bioorg. Med. Chem.* **2006**, *14*, 5876–5889.
- (31) Panek, J. J.; Jezierska, A.; Vracko, M. Kohonen Network Study of Aromatic Compounds Based on Electronic and Nonelectronic Structure Descriptors. *J. Chem. Inf. Model.* **2005**, *45*, 264–272.
- (32) Furusjö, E.; Svenson, A.; Rahmberg, M.; Andersson, M. The importance of outlier detection and training set selection for reliable environmental QSAR predictions. *Chemosphere* **2006**, *63*, 99–108.
- (33) Duchowicz, P. R.; Talevi, A.; Bruno-Blanch, L. E.; Castro, E. A. New QSPR study for the prediction of aqueous solubility of drug-like compounds. *Bioorg. Med. Chem.* **2008**, *16*, 7944–7955.
- (34) Bhonsle, J.; Huddler, D. Novel Method for Mining QSPR-Relevant Conformations. *Chem. Eng. Commun.* **2008**, *195*, 1396–1423.
- (35) Gharagheizi, F.; Tirandazi, B.; Barzin, R. Estimation of Aniline Point Temperature of Pure Hydrocarbons: A Quantitative Structure-Property Relationship Approach. *Ind. Eng. Chem. Res.* **2009**, *48*, 1678–1682.
- (36) Yang, S.-S.; Lu, W.-C.; Gu, T.-H.; Yan, L.-M.; Li, G.-Z. QSPR Study of n-Octanol/Water Partition Coefficient of Some Aromatic Compounds Using Support Vector Regression. *QSAR Comb. Sci.* **2009**, *28*, 175–182.
- (37) Kusic, H.; Rasulev, B.; Leszczynska, D.; Leszczynski, J.; Koprivanac, N. Prediction of rate constants for radical degradation of aromatic pollutants in water matrix: A QSAR study. *Chemosphere* **2009**, *75*, 1128–1134.
- (38) Goudarzi, N.; Goodarzi, M. QSPR Study of Partition Coefficient (Ko/w) of some Organic Compounds using Radial Basic Function-Partial Least Square (RBF-PLS). *J. Braz. Chem. Soc.* **2010**, *21*, 1776–1783.
- (39) Atabati, M.; Zarei, K.; Borhani, A. Predicting infinite dilution activity coefficients of hydrocarbons in water using ant colony optimization. *Fluid Phase Equilib.* **2010**, *293*, 219–224.
- (40) Rivera-Borroto, O. M.; Marrero-Ponce, Y.; Garcíade la Vega, J. M.; Grau-Abalo, R. d. C. Comparison of Combinatorial Clustering Methods on Pharmacological Data Sets Represented by Machine Learning-Selected Real Molecular Descriptors. *J. Chem. Inf. Model.* **2011**, *51*, 3036–3049.
- (41) 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.
- (42) Sadowski, J.; Gasteiger, J.; Klebe, G. Comparison of Automatic Three-Dimensional Model Builders Using 639 X-ray Structures. *J. Chem. Inf. Model.* **1994**, *34*, 1000–1008.
- (43) Karelson, M.; Lobanov, V.; Katritzky, A. Quantum-Chemical Descriptors in QSAR/QSPR Studies. *Chem. Rev.* **1996**, *96*, 1027–1043.
- (44) Jung, D.; Floyd, J.; Gund, T. A Comparative Molecular Field Analysis (CoMFA) Study Using Semiempirical, Density Functional, Ab Initio Methods and Pharmacophore Derivation Using DISCOtech on Sigma 1 Ligands. *J. Comput. Chem.* **2004**, *25*, 1385–1399.
- (45) Schüürmann, G. Quantum Chemical Descriptors in Structure-Activity Relationships – Calculation, Interpretation, and Comparison of Methods. In *Predicting Chemical Toxicity and Fate*, first ed.; Cronin, M. T. D., Livingstone, D. J., Eds.; CRC Press: Boca Raton, FL, 2004, pp 85–149.
- (46) Paster, I.; Shacham, M.; Brauner, N. Investigation of the Relationships between Molecular Structure, Molecular Descriptors, and Physical Properties. *Ind. Eng. Chem. Res.* **2009**, *48*, 9723–9734.
- (47) Cramer, D.; Patterson, D.; Bunce, J. Comparative Molecular Field Analysis (CoMFA). 1. Effect of Shape on Binding of Steroids to Carrier Proteins. *J. Am. Chem. Soc.* **1988**, *110*, 5959–5967.
- (48) Pissurlenkar, R.; Khedkar, V.; Iyer, R.; Coutinho, P. Ensemble QSAR: A QSAR Method Based on Conformational Ensembles and Metric Descriptors. *J. Comput. Chem.* **2011**, *32*, 2204–2218.
- (49) DragonX for Windows, version 1.4; Talete srl.: Milan Italy, 2009.

- (50) Lucas, K. *Molecular Models for Fluids*; Cambridge: New York, 2007.
- (51) Fernández-Ramos, A.; Ellingson, B.; Meana-Pañeda, R.; Marques, J.; Truhlar, D. Symmetry numbers and chemical reaction rates. *Theor. Chem. Acc.* **2007**, *118*, 813–826.
- (52) Gilson, M. K.; Irikura, K. K. Symmetry Numbers for Rigid, Flexible, and Fluxional Molecules: Theory and Applications. *J. Phys. Chem. B* **2010**, *114*, 16304–16317.
- (53) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. P.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, J. A. J.; Peralta, J. E.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, J. M.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Daniels, A. D.; Farkas, O.; Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. *Gaussian 09*, Revision A.02; Gaussian, Inc.: Wallingford, CT, 2009.
- (54) *Spartan10*, version 1.1.0; Wavefunction, Inc.: Irvine, CA, 2011.
- (55) *Turbomole*, version 6.0.2; Turbomole GmbH: Karlsruhe, Germany, 2009.
- (56) Halgren, T. Merck Molecular Force Field. II. MMFF94 van der Waals and Electrostatic Parameters for Intermolecular Interactions. *J. Comput. Chem.* **1996**, *17*, 520–550.
- (57) Dewar, M. J. S.; Zoebisch, E. G.; Healy, E. F.; Stewart, J. J. P. AM1: A New General Purpose Quantum Mechanical Molecular Model. *J. Am. Chem. Soc.* **1985**, *107*, 3902–3909.
- (58) Mayo, S. L.; Olafson, B. D.; Goddard, W. A. DREIDING: A generic force field for molecular simulations. *J. Phys. Chem.* **1990**, *94*, 8897–8909.
- (59) Becke, A. D. Density-functional thermochemistry. III. The role of exact exchange. *J. Chem. Phys.* **1993**, *98*, 5648–5652.
- (60) Lee, C.; Yang, W.; Parr, R. G. Development of the Colle-Salvetti correlation-energy formula into a functional of the electron density. *Phys. Rev. B* **1988**, *37*, 785–789.
- (61) Montgomery, J. A.; Frisch, M. J.; Ochterski, J. W.; Petersson, G. A. A complete basis set model chemistry. VI. Use of density functional geometries and frequencies. *J. Chem. Phys.* **1999**, *110*, 2822–2827.
- (62) Chirlian, L. E.; Francl, M. M. Atomic charges derived from electrostatic potentials: A detailed study. *J. Comput. Chem.* **1987**, *8*, 894–905.
- (63) Ferenczy, G.; Reynolds, A.; Richards, W. Semiempirical AM1 Electrostatic Potentials and AM1 Electrostatic Potential Derived Charges: A Comparison with Ab Initio Values. *J. Comput. Chem.* **1990**, *11*, 159–169.
- (64) Mulliken, R. S. Electronic Population Analysis on LCAOMO Molecular Wave Functions. I. *J. Chem. Phys.* **1955**, *23*, 1833–1840.
- (65) Sigfridsson, E.; Ryde, U. Comparison of methods for deriving atomic charges from the electrostatic potential and moments. *J. Comput. Chem.* **1998**, *19*, 377–395.
- (66) Møller, C.; Plesset, M. S. Note on an Approximation Treatment for Many-Electron Systems. *Phys. Rev.* **1934**, *46*, 618–622.
- (67) Weigend, F.; Köhn, A.; Hättig, C. Efficient use of the correlation consistent basis sets in resolution of the identity MP2 calculations. *J. Chem. Phys.* **2002**, *116*, 3175–3183.
- (68) Randić, M. Molecular Shape Profiles. *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 373–382.
- (69) Hemmer, M.; Steinhauer, V.; Gasteiger, J. Deriving the 3D structure of organic molecules from their infrared spectra. *Vib. Spectrosc.* **1999**, *19*, 151–164.
- (70) Schuur, J.; Selzer, P.; Gasteiger, J. The Coding of the Three-Dimensional Structure of Molecules by Molecular Transforms and Its Application to Structure-Spectra Correlations and Studies of Biological Activity. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 334–344.
- (71) Todeschini, R.; Lasagni, M. New molecular descriptors for 2D and 3D structures. Theory. *J. Chemom.* **1994**, *8*, 263–272.
- (72) Shacham, M.; Brauner, N.; Shore, H.; Benson-Karhi, D. Predicting temperature-dependent properties by correlations based on similarities of molecular structures: Application to liquid density. *Ind. Eng. Chem. Res.* **2008**, *47*, 4496–4504.
- (73) Hughes, L.; Palmer, D.; Nigsch, F.; Mitchell, J. Why Are Some Properties More Difficult To Predict than Others? A Study of QSPR Models of Solubility, Melting Point, and Log P. *J. Chem. Inf. Model.* **2008**, *48*, 220–232.
- (74) Goodarzi, M.; Duchowicz, P. R.; Freitas, M. P.; Fernandez, F. M. Prediction of the Hildebrand parameter of various solvents using linear and nonlinear approaches. *Fluid Phase Equilib.* **2010**, *293*, 130–136.
- (75) Eslamimanesh, A.; Gharagheizi, F.; Mohammadi, A. H.; Richon, D. Phase Equilibrium Modeling of Structure H Clathrate Hydrates of Methane + Water “Insoluble” Hydrocarbon Promoter Using QSPR Molecular Approach. *J. Chem. Eng. Data* **2011**, *56*, 3775–3793.
- (76) Yan, A.; Gasteiger, J. Prediction of Aqueous Solubility of Organic Compounds Based on a 3D Structure Representation. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 429–434.
- (77) Liu, H.; Papa, E.; Gramatica, P. QSAR Prediction of Estrogen Activity for a Large Set of Diverse Chemicals under the Guidance of OECD Principles. *Chem. Res. Toxicol.* **2006**, *19*, 1540–1548.
- (78) Mercader, A. G.; Duchowicz, P. R.; Fernandez, F. M.; Castro, E. A.; Bennardi, D. O.; Autino, J. C.; Romanelli, G. P. QSAR prediction of inhibition of aldose reductase for flavonoids. *Bioorg. Med. Chem.* **2008**, *16*, 7470–7476.
- (79) Noorizadeh, H.; Farmany, A.; Noorizadeh, M. Application of gas-PLS and GA-KPLS calculations for the prediction of the retention indices of essential oils. *Quim. Nova* **2011**, *34*, 1398–1404.
- (80) Klamt, A.; Eckert, F.; Hornig, M.; Beck, M. E.; Bürger, T. Prediction of aqueous solubility of drugs and pesticides with COSMO-RS. *J. Comput. Chem.* **2002**, *23*, 275–281.
- (81) Duffy, E.; Jorgensen, W. Prediction of Properties from Simulations: Free Energies of Solvation in Hexadecane, Octanol, and Water. *J. Am. Chem. Soc.* **2000**, *122*, 2878–2888.
- (82) NIST Chemistry WebBook, <http://webbook.nist.gov> (accessed June 26, 2012).