

Prediction of Molecular Solvation Free Energy Based on the Optimization of Atomic Solvation Parameters with Genetic Algorithm

Hongsuk Kang, Hwanho Choi, and Hwangseo Park*

Department of Bioscience and Biotechnology, Sejong University, 98 Kunja-Dong, Kwangjin-Ku, Seoul 143-747, Korea

Received October 23, 2006

We propose an improved solvent contact model to estimate the solvation free energy of an organic molecule from individual atomic contributions. The modification of the solvation model involves the optimization of three kinds of parameters in the solvation free energy function: atomic fragmental volume, maximum atomic occupancy, and atomic solvation parameters. All of these atomic parameters for 24 atom types are developed by the operation of a standard genetic algorithm in such a way as to minimize the difference between experimental and calculated solvation free energies. The data set for experimental solvation free energies is divided into a training set of 131 compounds and a test set of 24 compounds. Linear regressions with the optimized atomic parameters yield fits with the squared correlation coefficients (r^2) of 0.89 and 0.86 for the training set and for the test set, respectively. Overall, the results indicate that the improved solvent contact model with the newly developed atomic parameters would be a useful tool for rapid calculation of molecular solvation free energies in aqueous solution.

INTRODUCTION

Molecular solvation free energy or, in similar words, molecular solubility in aqueous solution is a very important quantity in a variety of chemical and biological processes including equilibria and kinetics for metabolic reactions, intermolecular association, and structural changes.^{1,2} The solubility is also an important property to be considered in drug discovery because it influences the bioactivity of a drug at the site of action, which renders the estimation of the aqueous solubility one key task at an early state of drug discovery.³ Furthermore, computing the solvation free energy has been a challenge for structure-based drug design because it is well appreciated that the desolvation effect plays a significant role in determining the binding mode and the binding affinity of the protein–ligand complex.^{4–7} However, the experimental measurement of solubility is a very time-consuming procedure, which prevents its use for the purpose of screening a large number of compounds. Now the need to estimate the differences between solubilities of structurally related compounds becomes more urgent with the development of combinatorial chemistry in recent years, further motivating the development of a reliable computational method to predict solvation free energies of organic molecules.

However, solvation free energy has been considered as one of the most calculation-difficult energy terms due mainly to the complexity of solvent–solute interactions.⁸ Many computational methods for solubility prediction have nonetheless been explored since the pioneering work of Yalkowsky and Valvani.⁹ More specifically, numerous statistical modeling methods have been examined using the molecular descriptors such as fragmental substructures^{10,11} and the

topological parameters^{12–14} and based on various theoretical frameworks including accessible surface area model,¹⁵ artificial neural network,^{16–18} QSPR,^{19–21} multiple linear regression,^{22–24} general solubility equation,²⁵ and free energy perturbation approaches.^{26,27} The continuum electrostatic models of solvation have also been proposed to deal with molecular solvation free energy.²⁸ Among them, the simplest one was to adjust the dielectric constant in a distant-dependent fashion to model electrostatic screening by solvent.^{29,30} A more precise model for predicting molecular solvation free energies involved the solution of the Poisson–Boltzmann equation to calculate electrostatic potentials around a molecule under investigation.³¹ In the continuum models, however, the structural change of a solute upon solvation could not be taken into account, which has limited their usefulness to the solutes of simple ions and small molecules.

In the early 1990s, Stouten et al. suggested a solvation model for a protein molecule by extending the solvent contact model proposed by Colonna-Cesari and Sander.^{32,33} The three key parameters in this model were the maximum atomic occupancy, the atomic fragmental volume, and the atomic solvation parameters representing the solvation free energy per unit of volume.³³ Under the assumption that the solvation free energy of an amino acid residue would be given by the sum over atomic contributions, they obtained the atomic parameters for six atom types (C, N, O, N⁺, O[−], and S) using the standard linear least-squares procedures with the experimental solvation free energies of amino acids. This simple solvation model proved to be very successful in estimating the structural properties of a protein as well as in saving computation time in molecular dynamics simulations when compared to the explicit solvent model.³³ Due to such a small number of atom types and the exclusive use of amino acids in the training set, however, some proper modifications need to be made in order for the extended solvent contact model

* Corresponding author phone: +82-2-3408-3766; fax: +82-2-3408-3334; e-mail: hspark@sejong.ac.kr.

to be also useful in predicting solvation free energies of organic compounds.

In the present study, we further improve the Stouten et al.'s solvent contact model by extending the parameter space to cope with as many atom types as commonly encountered in normal organic compounds. All of the atomic parameters in the solvation free energy function are optimized by the operation of a standard genetic algorithm (GA) using the experimental solvation free energy data. It will be shown that the improved solvent contact model with the newly developed atomic parameters can be an appropriate tool for predicting solvation free energies of organic molecules in aqueous solution.

COMPUTATIONAL METHODS

Data Set. In the optimization of atomic parameters with genetic algorithm to calculate molecular solvation free energies, we worked with a chemical library containing 145 molecules for which experimental solvation energies have been reported.³⁴ These total 145 molecules were divided into a training set of 131 compounds and a test set of 24 compounds. All of the compounds included in the two sets were then subjected to the CORINA program to generate their 3-D coordinates in the Sybyl MOL2 format.³⁵ As implemented in CORINA, only a single conformation of each molecule was generated based on the conformational parameters derived from the X-ray structures of small molecules. The 3-D structures obtained in this way have been shown to be similar to the molecular geometries optimized with the semiempirical AM1 calculations including solvation effects,³⁶ which indicates the reasonableness of the molecular structures derived with CORINA.

Definition of Atom Types. Different atom types should have different contributions to solvation free energy in the present solvation model under investigation. We used 24 basic atom types for the elements commonly found in organic molecules. The atom type of a given atom in a molecule was differentiated according to element, hybridization state, and chemical environment around the atom under consideration. Considering the portability and the simplicity of implementation of the classifications, all atom types were designated in the same fashion as in the Sybyl MOL2 format. The phosphorus atom was excluded in the present parametrization because the proper experimental data for solvation free energy have not been at hand. Upon the availability of the additional solvation data in the future, however, we will also optimize the atomic solvation parameters for various atom types of the phosphorus atom.

Optimization of Atomic Volume Parameters with Genetic Algorithm. Three kinds of atomic parameters need to be optimized in order to calculate the solvation free energy of a molecule based on the solvent contact model. Among them, the atomic volume parameter V_j represents the fragmental volume of atom j in a molecule. Because the V_j values exhibited a bad convergent behavior in the simultaneous optimization of the three kinds of parameters, they were optimized with the operation of an independent genetic algorithm as detailed below.

The total volume of a molecule should be determined prior to the parametrization of V_j values. For this purpose, each molecule in the training and the test sets was placed in a

3-D box whose length, width, and height correspond to the maximum distances along the three axes defining the coordinate system of the van der Waals volume of the molecule. Monte Carlo simulations involving random selections of a point in the predefined 3-D box were then carried out to calculate the total volume of the molecule (V_{mol}) embedded in the box. In this simulation, V_{mol} could be obtained by the volume of the box (V_{box}) multiplied by the ratio of the number of trials to select a point in the molecular van der Waals volume (N_{hits}) to the total number of trials (N_{trials}). Thus, we have

$$V_{\text{mol}} = V_{\text{box}} \times \frac{N_{\text{hits}}}{N_{\text{trials}}} \quad (1)$$

With the calculated V_{mol} values in hand, the atomic volume parameters were optimized with the standard genetic algorithm. A generation was defined with 100 vectors comprising the V_j parameters, followed by the removal of 50 with a bias toward preserving the most fit with the lowest error. The empty 50 vectors were then filled with point mutations to alter the value of one of the parameters with probability 0.01 and with cross breeds with probability 0.6 to select some parameters from one vector to replace the elements of another vector of the top 50. The 50 newly created vectors were then evaluated together with the top 50. This cycle was repeated as many times as desired. In the evaluation of the 100 vectors, we used a gradient-based minimization method on the error hypersurface (F_V). This hypersurface is defined by the sum of the absolute values of the differences between the calculated V_{mol} value of a molecule and the sum of V_j values in the molecule.

$$F_V = \sum_k^{\text{molecules}} |V_{\text{mol}}^k - \sum_j^{\text{atoms}} V_j| \quad (2)$$

Calculation of Solvation Parameters with Genetic Algorithm. The solvent contact model to calculate molecular solvation free energy is based on several fundamental assumptions. First, the solvation free energy of a molecule k can be approximated by the sum of individual atomic contributions.

$$\Delta G_{\text{calc}}^k = \sum_i^{\text{atoms}} \Delta G_{\text{sol}}^i \quad (3)$$

Second, the individual solvation energy of an atom i in the molecule can be given by the product of the atomic solvation parameter (S_i) and the degree of its exposure to bulk solvent (F_i).

$$\Delta G_{\text{sol}}^i = S_i F_i \quad (4)$$

Third, the atomic degree of exposure is approximated as the percentage of the unoccupied volume around the atom in the molecule. The occupied volume around the atom i (O_i) can then be determined by summing the atomic volume parameters representing the fragmental volumes of all other atoms in the molecule multiplied by a suitable envelope function.

$$O_i = \sum_{j \neq i}^{\text{atoms}} V_j e^{-\frac{r_{ij}^2}{2\sigma^2}} \quad (5)$$

Here, the Gaussian envelope function is employed with the variable r_{ij} representing the distance between the centers of atoms i and j in the molecule. Because F_i is the difference between the maximum occupancy of atom i (O_i^{\max}) in a molecule³³ and O_i , the solvation free energy of a molecule k can be expressed as

$$\Delta G_{\text{calc}}^k = \sum_i^{\text{atoms}} S_i (O_i^{\max} - \sum_j^{i \neq j} V_j e^{-\frac{r_{ij}^2}{2\sigma^2}}) \quad (6)$$

Therefore, the two atomic parameters (S_i and O_i^{\max}) need also to be optimized in addition to V_j to estimate the molecular solvation free energy based on eq 6. These parametrizations were carried out by operating the genetic algorithm with the same procedure as in the optimization of atomic volume parameters. We used a gradient-based minimization method on the error hypersurface defined by the sum of the absolute values of the differences between the calculated and experimental solvation free energies. Formally this fitness function is defined as

$$F_s = \sum_{i=1}^{\text{molecules}} |\Delta G_{\text{exp}}^i - \Delta G_{\text{calc}}^i| \quad (7)$$

RESULTS AND DISCUSSION

Listed in Table 1 are the optimized atomic volume (V_j), maximum atomic occupancy (O_i^{\max}), and atomic solvation parameters (S_i) for the 24 atom types that are necessary to depict 131 molecules in the training set. It should be noted that the atomic solvation parameters optimized in this study differ from those obtained by Stouten et al. in that their use should be limited to small organic molecules because no biomolecule was included in the training set. The V_j values calculated in this study are very different from the atomic volumes of isolated atoms. The reason lies in that each V_j value represents the average of the contributions of the atom with type j to the van der Waals volumes of various sizes and shapes the molecules in the training set can have. This indicates that the V_j values may exhibit a strong dependence on the molecules comprising the training set, whereas the atomic volume of an isolated atom has to be a constant value.

On the other hand, the calculated S_i parameters reveal a trend consistent with general atomic properties. We note, for example, that the S_i values get more negative in going from sp^3 to sp^2 and sp in the case of carbon atoms. This indicates that atomic solvation would be more favorable with the increase of the s-character in the hybridization of atomic orbitals of a carbon atom. This is not surprising because the increase in s-character raises the electronegativity of carbon atom, which would have an effect of increasing the stability in aqueous solution by facilitating the hydrophilic interactions with solvent molecules. In the case of nitrogen atom, on the contrary, the S_i values are shown to be more negative in the order of $\text{sp}^3 \geq \text{sp}^2 \geq \text{sp}$, which exhibits the same trend as in the order of basicity. The greater basicity of a nitrogen atom with the lower s-character is attributed, in general, to the

Table 1. Atomic Fragmental Volume (V_j), Maximum Atomic Occupancy (O_i^{\max}), and Atomic Solvation Parameters (S_i) Optimized with Genetic Algorithm

| atom type | description | V_j (Å ³) | O_i^{\max} (Å ³) | S_i (kcal/mol Å ³) |
|-----------|--|-------------------------|--------------------------------|----------------------------------|
| C.3 | sp ³ carbon | 14.919 | 320.000 | 0.667 |
| C.2 | sp ² carbon | 10.612 | 334.588 | 0.510 |
| C.1 | sp carbon | 11.612 | 323.765 | 0.039 |
| C.ar | aromatic carbon | 13.790 | 361.882 | −0.667 |
| N.3 | sp ³ nitrogen | 10.484 | 323.765 | −22.000 |
| N.2 | sp ² nitrogen | 12.534 | 324.471 | −16.353 |
| N.1 | sp nitrogen | 14.548 | 366.353 | −12.588 |
| N.ar | aromatic nitrogen | 9.354 | 333.882 | −10.470 |
| O.3 | sp ³ oxygen in hydroxyl group | 8.870 | 365.882 | −6.941 |
| O.4 | sp ³ oxygen in ether group | 8.870 | 320.000 | −5.059 |
| O.5 | sp ³ oxygen in ester group | 8.870 | 320.000 | −12.588 |
| O.2 | sp ² oxygen | 14.435 | 336.941 | −12.353 |
| O.co2 | carboxylate oxygen | 11.774 | 330.353 | −12.588 |
| S.3 | sp ³ sulfur | 14.314 | 346.353 | −12.118 |
| S.2 | sp ² sulfur | 12.836 | 341.882 | −15.882 |
| F | fluorine | 7.419 | 320.000 | 0.784 |
| Cl | chlorine | 17.580 | 320.000 | −2.941 |
| Br | bromine | 21.032 | 331.529 | −2.745 |
| I | iodine | 30.191 | 320.000 | −5.098 |
| H | hydrogen bonded to carbon | 7.903 | 367.059 | 0.000 |
| H.2 | hydrogen bonded to nitrogen | 7.257 | 365.412 | 0.784 |
| H.3 | hydrogen bonded to oxygen | 6.935 | 361.176 | −11.177 |
| H.4 | hydrogen bonded to sulfur | 7.486 | 320.235 | 7.843 |

reduced electronegativity that is responsible for the increase in the tendency of the lone electron pair to react with proton in aqueous solution. It is thus apparent that the nitrogen atom with higher basicity has a greater tendency to be stabilized by the establishment of hydrogen bond interactions with solvent molecules, which can be invoked to explain its more negative atomic solvation parameter than the less basic nitrogen atom.

The sp^3 oxygens of hydroxy, ether, and ester groups are subdivided to the three different atom types and parametrized independently, due to a significant difference in their solvation free energies. It is noted in Table 1 that the sp^3 oxygen in the ester group has a much more negative S_i value than those in hydroxy and ether groups. With respect to the hydrogen atom, different atom types are defined depending on the atom to which it is attached. As a result, we obtained significantly different S_i values among the four different atom types. The definition of differentiating atom types on the basis of chemical environment around an atom is necessary to obtain the accurate S_i parameters. In this regard, it has been observed that the squared correlation coefficient (r^2) values increase by 0.05 and 0.06 due to the subdivisions of the atom types of sp^3 oxygen and hydrogen, respectively.³⁷

Table 2 lists the calculated solvation free energies for the training and the test sets with the optimized atomic parameters shown in Table 1, in comparison with the experimental ones. It is seen that the calculated solvation free energies compare reasonably well with the experimental results except for several cases. Twenty-one out of 131 compounds in the training set have deviations from the corresponding experimental values greater than 1 kcal/mol, while 9 out of 24 compounds have such deviations in the test set. The greatest deviation from the experimental solvation free energy occurs in the case of compound 54 (1,3-dioxolane): its aqueous solubility is underestimated by 2.19 kcal/mol. Judging from the poor atomic solvation of the sp^3 carbon as shown in Table 1, the underestimation of the solubility of 1,3-dioxolane is

Table 2. Experimental and Calculated Solvation Free Energies (in kcal/mol) of (a) the 131 Compounds in the Training Set and (b) the 24 Compounds in the Test Set

| no. | molecule name | expt | calc | no. | molecule name | expt | calc |
|------------------|-------------------------------------|-------|-------|-----|------------------------------------|--------|--------|
| (a) Training Set | | | | | | | |
| 1 | cyclopentane | 1.22 | 0.75 | 67 | methanol | -5.14 | -5.79 |
| 2 | methylcyclopentane | 1.62 | 0.85 | 68 | methane thiol | -1.26 | -1.26 |
| 3 | methylcyclohexane | 1.73 | 0.92 | 69 | ethanol | -4.96 | -5.31 |
| 4 | ethylene | 1.30 | 0.31 | 70 | 2,2,2-trifluoroethanol | -4.35 | -4.57 |
| 5 | 1-pentene | 1.69 | 0.72 | 71 | 1-propanol | -4.92 | -4.96 |
| 6 | cyclopentene | 0.57 | 0.72 | 72 | allyl alcohol | -5.10 | -5.10 |
| 7 | 2-methyl-2-butene | 1.33 | 0.71 | 73 | 1,1,1-trifluoro-2-propanol | -4.21 | -4.16 |
| 8 | cyclohexene | 0.37 | 0.81 | 74 | 1-butanol | -4.78 | -4.71 |
| 9 | 4-methyl-1-pentene | 1.93 | 0.81 | 75 | 2-methyl-1-propanol | -4.57 | -4.59 |
| 10 | 1-methylcyclohexene | 0.68 | 0.91 | 76 | 1-pentanol | -4.55 | -4.53 |
| 11 | 1-octene | 2.20 | 1.06 | 77 | cyclohexanol | -5.02 | -4.01 |
| 12 | 1,3-butadiene | 0.57 | 0.55 | 78 | 4-methyl-2-pentanol | -3.79 | -3.89 |
| 13 | 2-methyl-1,3-butadiene | 0.69 | 0.68 | 79 | phenol | -6.62 | -6.20 |
| 14 | 1,5-hexadiene | 1.02 | 0.80 | 80 | 4-bromophenol | -7.20 | -6.79 |
| 15 | propyne | -0.48 | 0.20 | 81 | thiophenol | -2.58 | -2.27 |
| 16 | 1-pentyne | 0.01 | 0.50 | 82 | 2-cresol | -5.94 | -5.73 |
| 17 | 1-heptyne | 0.61 | 0.75 | 83 | 4-hydroxybenzaldehyde | -10.61 | -9.19 |
| 18 | 1-nonyne | 1.06 | 0.97 | 84 | 4- <i>tert</i> -butylphenol | -6.00 | -5.22 |
| 19 | benzene | -0.90 | -1.09 | 85 | acetaldehyde | -3.55 | -3.37 |
| 20 | <i>p</i> -xylene | -0.82 | -0.68 | 86 | butanal | -3.22 | -2.81 |
| 21 | propylbenzene | -0.54 | -0.54 | 87 | heptanal | -2.71 | -2.35 |
| 22 | 2-butylbenzene | -0.46 | -0.41 | 88 | <i>trans</i> -2-butenal | -4.28 | -2.95 |
| 23 | naphthalene | -2.45 | -1.60 | 89 | <i>trans</i> -2-octenal | -3.48 | -2.35 |
| 24 | fluoromethane | -0.22 | 0.43 | 90 | <i>trans-trans</i> -2,4-hexadienal | -4.70 | -2.65 |
| 25 | trifluoromethane | 0.82 | 0.86 | 91 | benzaldehyde | -4.08 | -4.15 |
| 26 | chloromethane | -0.54 | -0.67 | 92 | acetone | -3.85 | -3.00 |
| 27 | trichloromethane | -1.04 | -2.25 | 93 | 2-pentanone | -3.56 | -2.40 |
| 28 | bromomethane | -0.80 | -0.64 | 94 | 2-octanone | -2.92 | -1.89 |
| 29 | tribromomethane | -2.16 | -2.16 | 95 | acetophenone | -4.64 | -3.77 |
| 30 | iodomethane | -0.90 | -1.30 | 96 | acetic acid | -6.78 | -6.89 |
| 31 | chlorofluoromethane | -0.79 | -0.44 | 97 | butyric acid | -6.44 | -6.01 |
| 32 | chloroethane | -0.64 | -0.45 | 98 | methylacetate | -3.66 | -4.18 |
| 33 | bromoethane | -0.70 | -0.44 | 99 | ethylacetate | -3.12 | -3.80 |
| 34 | iodoethane | -0.73 | -1.06 | 100 | isopropylacetate | -2.68 | -3.42 |
| 35 | 1,2-dichloroethane | -1.75 | -1.23 | 101 | isobutylacetate | -2.39 | -3.23 |
| 36 | 1,2-dibromoethane | -2.13 | -1.19 | 102 | propylbutyrate | -2.31 | -2.82 |
| 37 | 1-chloro-2-bromoethane | -1.98 | -1.21 | 103 | methyloctanoate | -2.07 | -2.80 |
| 38 | 2-chloro-1,1,1-trifluoroethane | 0.06 | 0.17 | 104 | methylbenzoate | -4.34 | -4.76 |
| 39 | 1-chloropropane | -0.36 | -0.28 | 105 | ethylamine | -4.67 | -4.94 |
| 40 | 1-bromopropane | -0.57 | -0.27 | 106 | butylamine | -4.43 | -4.24 |
| 41 | 1-iodopropane | -0.62 | -0.87 | 107 | hexylamine | -4.09 | -3.87 |
| 42 | <i>cis</i> -1,2-dichloroethylene | -1.19 | -1.36 | 108 | dimethylamine | -4.34 | -5.18 |
| 43 | <i>trans</i> -1,2-dichloroethylene | -0.77 | -1.14 | 109 | diethylamine | -4.12 | -4.16 |
| 44 | 3-chloropropene | -0.58 | -0.34 | 110 | pyrrolidine | -5.54 | -4.33 |
| 45 | chlorobenzene | -1.02 | -1.79 | 111 | piperidine | -5.17 | -3.87 |
| 46 | bromobenzene | -1.48 | -1.76 | 112 | dipropylamine | -3.70 | -3.39 |
| 47 | 1,2-dichlorobenzene | -1.38 | -2.46 | 113 | hexamethyleneimine | -4.97 | -3.47 |
| 48 | 1,4-dibromobenzene | -2.33 | -2.41 | 114 | trimethylamine | -3.27 | -4.88 |
| 49 | <i>p</i> -bromotoluene | -1.41 | -1.54 | 115 | triethylamine | -3.07 | -3.40 |
| 50 | 1-bromo-2-ethylbenzene | -1.20 | -1.30 | 116 | <i>n</i> -methylpyrrolidine | -4.02 | -4.04 |
| 51 | <i>o</i> -bromocumene | -0.86 | -1.12 | 117 | <i>n</i> -methylpiperidine | -3.94 | -3.58 |
| 52 | dimethyl ether | -1.92 | -2.97 | 118 | ethylenediamine | -9.88 | -9.90 |
| 53 | dimethyl sulfide | -1.56 | -3.24 | 119 | acetonitrile | -3.94 | -4.01 |
| 54 | 1,3-dioxolane | -4.14 | -1.95 | 120 | butyronitrile | -3.69 | -3.42 |
| 55 | methylpropyl ether | -1.69 | -2.29 | 121 | nitroethane | -3.76 | -3.51 |
| 56 | tetrahydrofuran | -3.51 | -2.31 | 122 | 2-nitropropane | -3.18 | -3.18 |
| 57 | dioxane | -5.11 | -5.05 | 123 | nitrobenzene | -4.17 | -4.57 |
| 58 | methyl <i>tert</i> -butyl ether | -2.24 | -1.89 | 124 | 3-nitrotoluene | -3.50 | -4.30 |
| 59 | 2-methyltetrahydrofuran | -3.34 | -1.97 | 125 | pyridine | -4.75 | -3.55 |
| 60 | tetrahydropyran | -3.16 | -1.99 | 126 | 4-methylpyridine | -4.99 | -3.26 |
| 61 | dipropyl ether | -1.17 | -1.63 | 127 | 4-ethylpyridine | -4.78 | -2.98 |
| 62 | 1,2-diethoxyethane | -3.30 | -3.35 | 128 | 2,4-dimethylpyridine | -4.92 | -2.85 |
| 63 | di- <i>n</i> -butyl ether | -0.84 | -1.19 | 129 | 2-methylpyrazine | -5.58 | -5.58 |
| 64 | anisole | -1.05 | -3.08 | 130 | 2-ethyl-3-methoxypyrazine | -4.45 | -5.47 |
| 65 | thioanisole | -2.76 | -3.91 | 131 | 2-isobutyl-3-methoxypyrazine | -3.73 | -4.69 |
| 66 | 2,2'-dichlorodiethyl sulfide | -3.97 | -3.91 | | | | |
| (b) Test Set | | | | | | | |
| 1 | 2-methylpentane | 2.56 | 1.28 | 13 | 3-hexanol | -3.73 | -5.30 |
| 2 | <i>cis</i> -1,2-dimethylcyclohexane | 1.60 | 1.71 | 14 | 4-nitrophenol | -10.74 | -12.55 |
| 3 | 1-hexene | 1.73 | 1.19 | 15 | hexanal | -2.85 | -2.92 |
| 4 | 2,3-dimethyl-1,3-butadiene | 0.40 | 1.11 | 16 | 2-butanone | -3.76 | -3.35 |
| 5 | toluene | -0.77 | -1.23 | 17 | methylformate | -2.82 | -5.40 |
| 6 | <i>tert</i> -butylbenzene | -0.44 | -0.59 | 18 | ethylpropionate | -2.83 | -4.76 |
| 7 | dichloromethane | -1.42 | -1.67 | 19 | isoamylacetate | -2.24 | -4.33 |
| 8 | 1,3-dibromopropane | -1.99 | -1.18 | 20 | propylamine | -4.56 | -5.41 |
| 9 | chloroethylene | 0.50 | -0.40 | 21 | dibutylamine | -3.38 | -5.13 |
| 10 | 1,4-dichlorobenzene | -1.02 | -3.33 | 22 | 1-nitropropane | -3.38 | -3.88 |
| 11 | diethyl sulfide | -1.45 | -3.34 | 23 | 4-ethylpyridine | -4.66 | -4.28 |
| 12 | diisopropyl ether | -0.50 | -2.75 | 24 | 2-isobutylpyrazine | -5.11 | -7.10 |

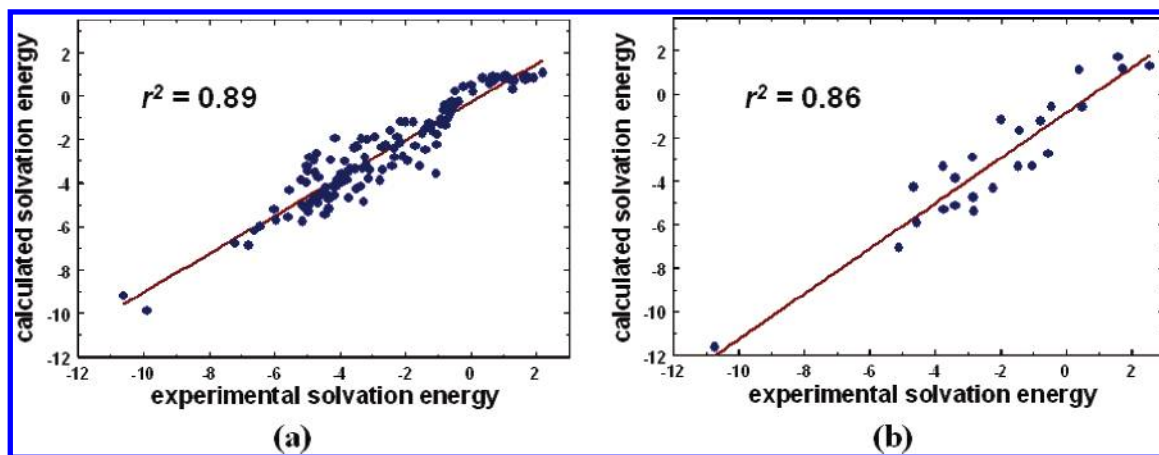


Figure 1. Correlation between experimental versus calculated solvation free energies for (a) 131 molecules in the training set and (b) 24 molecules in the test set. All energy values are given in kcal/mol.

due most likely to the underestimation of its molecular polarity in the present solvation model. Indeed, its three sp^3 carbons should be discriminated from those in the other molecules and parametrized separately because they are bonded to three or four electronegative atoms and therefore are expected to be extraordinarily electron-deficient.

One of the limitations of the solvation models based on atomic contributions has been their incapability to explain the difference in solvation free energies between stereoisomers. This is because the atomic composition has been considered as the only determinant for molecular solvation free energy. In the present solvation model, in contrast, the effects of 3-D atomic distribution are also reflected in the calculation of molecular solvation free energy. For example, the isomers with the same atomic composition can be discriminated in the present solvation model due to the differences between the isomers in r_{ij} 's in the envelope function as shown in eq 6. As a consequence, the *cis* and *trans* isomers of 1,2-dichloroethylene (compounds 42 and 43) are assigned different solvation free energies within the present solvation model, which is consistent with the experimental results.

The correlation between the experimental and the calculated solvation free energies are illustrated in Figure 1. With the test set of 24 molecules, we obtain the r^2 value of 0.86, which is a little worse than the fitting with the training set of 131 molecules. The accuracy of the present method in predicting molecular solvation free energy is similar to that of the QSPR model trained with 775 compounds in which some druglike properties of organic compounds computed from their 2-D structures were used as molecular descriptors.³⁸ The quality of the present solvation model is also comparable to that of the artificial neural network (ANN) model reported by Liu and So which was trained with 1033 compounds using 19 adjustable variables.³⁹ The comparisons thus indicate that our GA-based parametrization method would be more efficient in estimating molecular solvation free energies than the QSPR and the ANN models because the former could be trained with significantly fewer molecules than the latter to achieve a similar correlation. Most probably, such an enhanced efficiency is due to the direct use of 3-D structures in the parametrizations rather than 1-D or 2-D molecular descriptors as in the other methods.

The present GA-based solvation model involving the atomic parametrizations has additional advantages over the

traditional statistical models. First, 3-D molecular structures have only to be provided prior to the optimizations of the individual atomic parameters. The computational cost for molecular solvation free energy can therefore be saved to a significant extent when compared to the other methods in which molecular descriptors need to be calculated to construct a statistical model for solvation. Such a computational acceleration enables the present solvation model to be an appropriate tool to cope with large chemical libraries. Second, the solvation free energy function given in eq 6 and the newly developed parameters can be incorporated into the potential energy function of a molecule in aqueous solution as an implicit solvation model. This effective solvation term is likely to be efficient in terms of both saving computational cost for atomistic simulations and exploring structural properties of organic compounds just as Stouten et al.'s previous solvent contact model revealed such an efficiency and accuracy in molecular dynamics simulation of proteins in aqueous solution.³³ Finally, as mentioned above, the accuracy of the present GA-based solvation model can be enhanced in a straightforward way by subdividing the atom types according to the chemical environment around the atom of interest. The atomic parameters of some atom types could not be determined in this study due to the unavailability of corresponding experimental data. We will extend the atomic solvation parameter space to cover a more variety of atom types in molecules with the availability of more experimental data in the future.

CONCLUSIONS

We have shown the outperformance of the modified solvent contact model involving the GA-based atomic parametrizations in predicting molecular solvation free energies in aqueous solution. The present solvation model is based only on 3-D molecular coordinates with no additional molecular descriptors being required to calculate solvation free energy. Using the newly developed atomic parameters for 24 atom types with genetic algorithm, the solvation model was able to predict the experimental solvation free energies of a variety of organic compounds with the r^2 values of 0.89 and 0.86 for training and test sets, respectively. Considering the efficiency in energy calculation and the simplicity in model refinement by subdividing atom types, we expect that the present solvation model will be a

new useful tool for rapid calculation of molecular solvation free energies.

ACKNOWLEDGMENT

This work was partially supported by the grant from the Stroke Oriental Medicine Project (M1052701000005N2701-00000) of the Ministry of Science and Technology of Korea. The authors would also like to acknowledge the support from KISTI (Korea Institute of Science and Technology Information) under "The Eighth Strategic Supercomputing Support Program" with Dr. Sang Min Lee as the technical supporter. The use of the computing system of the Supercomputing Center is greatly appreciated.

REFERENCES AND NOTES

- Honig, B.; Nicholls, A. Classical Electrostatics in Biology and Chemistry. *Science* **1995**, *268*, 1144–1149.
- Cramer, C. J.; Truhlar, D. G. Implicit Solvation Models: Equilibria, Structure, Spectra, and Dynamics. *Chem. Rev.* **1999**, *99*, 2161–2200.
- Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings. *Adv. Drug Delivery Rev.* **1997**, *23*, 3–25.
- Pei, J.; Wang, Q.; Zhou, J.; Lai, L. Estimating Protein-Ligand Binding Free Energy: Atomic Solvation Parameters for Partition Coefficient and Solvation Free Energy Calculation. *Proteins* **2004**, *57*, 651–664.
- Zou, X.; Sun, Y.; Kuntz, I. D. Inclusion of Solvation in Ligand Binding Free Energy Calculations Using Generalized-Born Model. *J. Am. Chem. Soc.* **1999**, *121*, 8033–8043.
- Shoichet, B. K.; Leach, A. R.; Kuntz, I. D. Ligand Solvation in Molecular Docking. *Proteins* **1999**, *34*, 4–16.
- Wang, H.; Ben-Naim, A. A. Possible Involvement of Solvent-Induced Interactions in Drug Design. *J. Med. Chem.* **1996**, *39*, 1531–1539.
- Jorgensen, W. L.; Duffy, E. M. Prediction of Drug Solubility from Structure. *Adv. Drug Delivery Rev.* **2002**, *54*, 355–366.
- Yalkowsky, S.; Valvani, S. Solubility and Partitioning I: Solubility of Nonelectrolytes in Water. *J. Pharm. Sci.* **1980**, *69*, 912–922.
- Clark, M. Generalized Fragment-Substructure Based Property Prediction Method. *J. Chem. Inf. Model.* **2005**, *45*, 30–38.
- Pepe, G.; Guiliani, G.; Loustalet, S.; Halfon, P. Hydration Free Energy a Fragmental Model and Drug Design. *Eur. J. Med. Chem.* **2002**, *37*, 865–872.
- Laffort, P.; Hericourt, P. A Simplified Molecular Topology to Generate Easily Optimized Values. *J. Chem. Inf. Model.* **2006**, *46*, 1723–1734.
- Engkvist, O.; Wrede, P. High-Throughput, in silico Prediction of Aqueous Solubility Based on One- and Two-Dimensional Descriptors. *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 1247–1249.
- Huuskonen, J. Estimation of Aqueous Solubility for a Diverse Set of Organic Compounds Based on Molecular Topology. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 773–777.
- Hou, T.; Qiao, X.; Zhang, W.; Xu, X. Empirical Aqueous Solvation Models Based on Accessible Surface Areas with Implicit Electrostatics. *J. Phys. Chem. B* **2002**, *106*, 11295–11304.
- 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.
- Wegner, J.; Zell, A. Prediction of Aqueous Solubility and Partition Coefficient Optimized by a Genetic Algorithm Based Descriptor Selection Method. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1077–1084.
- Bruneau, P. Search for Predictive Generic Model of Aqueous Solubility using Bayesian Neural Nets. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1605–1616.
- Katritzky, A. R.; Oliferenko, A. A.; Oliferenko, P. V.; Petrukhin, R.; Tatham, D. B.; Maran, U.; Lomaka, A.; Acree, W. E., Jr. A General Treatment of Solubility. 2. QSPR Prediction of Free Energies of Solvation of Specified Solutes in Ranges of Solvents. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1806–1814.
- McElroy, N. R.; Jurs, P. C. Prediction of Aqueous Solubility of Heteroatom-Containing Organic Compounds from Molecular Structure. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1237–1247.
- Yaffe, D.; Cohen, Y.; Espinosa, G.; Arenas, A.; Giralto, F. A Fuzzy ARTMAP Based on Quantitative Structure-Property Relationships (QSPRs) for Predicting Aqueous Solubility of Organic Compounds. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1177–1207.
- Catana, C.; Gao, H.; Orrenius, C.; Stouten, P. F. W. Linear and Nonlinear Methods in Modeling the Aqueous Solubility of Organic Compounds. *J. Chem. Inf. Model.* **2005**, *45*, 170–176.
- Butina, D.; Gola, J. M. R. Modeling Aqueous Solubility. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 837–841.
- Hou, T. J.; Xia, K.; Zhang, W.; Xu, X. J. ADME Evaluation in Drug Discovery. 4. Prediction of Aqueous Solubility Based on Atom Contribution Approach. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 266–275.
- Ran, Y.; Yalkowsky, S. H. Prediction of Drug Solubility by the General Solubility Equation. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 354–357.
- Pitera, J. W.; van Gunsteren, W. F. One-Step Perturbation Method for Solvation Free Energies of Polar Solutes. *J. Phys. Chem. B* **2001**, *105*, 11264–11274.
- Reddy, M. R.; Singh, U. C.; Erion, M. D. Development of a Quantum Mechanics-Based Free Energy Perturbation Method: Use in the Calculation of Relative Solvation Free Energies. *J. Am. Chem. Soc.* **2004**, *126*, 6224–6225.
- Orozco, M.; Luque, F. J. Theoretical Methods for the Description of the Solvent Effect in Biomolecular Systems. *Chem. Rev.* **2000**, *100*, 4187–4225.
- Pickersgill, R. W. A Rapid Method of Calculating Charge-Charge Interaction Energies in Proteins. *Protein Eng.* **1988**, *2*, 247–248.
- Mehler, E. L.; Solmajer, T. Electrostatic Effects in Proteins: Comparison of Dielectric and Charge Models. *Protein Eng.* **1991**, *4*, 903–910.
- Gilson, M. K.; Sharp, K. A.; Honig, B. H. Calculating the Electrostatic Potential of Molecules in Solution: Method and Error Assessment. *J. Comput. Chem.* **1988**, *9*, 327–335.
- Colonna-Cesari, F.; Sander, C. Excluded Volume Approximation to Protein-Solvent Interaction. The Solvent Contact Model. *Biophys. J.* **1990**, *57*, 1103–1107.
- Stouten, P. F. W.; Frömmel, C.; Nakamura, H.; Sander, C. An Effective Solvation Term Based on Atomic Occupancies for Use in Protein Simulations. *Mol. Simul.* **1993**, *10*, 97–120.
- Viswanadhan, V. N.; Ghose, A. K.; Singh, U. C.; Wendoloski, J. J. Prediction of Solvation Free Energy of Small Organic Molecules: Additive-Constitutive Models Based on Molecular Fingerprints and Atomic Constants. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 405–412.
- Sadowski, J.; Gasteiger, J.; Klebe, G. Comparison of Automatic Three-Dimensional Model Builders Using 639 X-ray Structures. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 1000–1008.
- Goller, A. H.; Hennemann, M.; Keldenich, J.; Clark, T. In Silico Prediction of Buffer Solubility Based on Quantum-Mechanical and HQSAR- and Topology-Based Descriptors. *J. Chem. Inf. Model.* **2006**, *46*, 648–658.
- The subdivisions of the atom types of the sp³ oxygen and the hydrogen lead to the increases in the *r*² values from 0.84 to 0.89 and from 0.83 to 0.89, respectively.
- Cheng, A.; Merz, K. M., Jr. Prediction of Aqueous Solubility of a Diverse Set of Compounds Using Quantitative Structure-Property Relationships. *J. Med. Chem.* **2003**, *46*, 3572–3580.
- Liu, R.; So, S.-S. Development of Quantitative Structure-Property Relationship Models for Early ADME Evaluation in Drug Discovery. 1. Aqueous Solubility. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1633–1639.

CI600453B