



Extended solvent-contact model for protein solvation: Test cases for dipeptides

Hwanho Choi¹, Hongsuk Kang¹, Hwangseo Park^{*}

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

ARTICLE INFO

Article history:

Accepted 13 February 2013

Available online 16 March 2013

Keywords:

Solvation energy function ; Self-solvation

Genetic algorithm

Protein

Dipeptide

ABSTRACT

Solvation effects are critically important in the structural stabilization and functional optimization of proteins. Here, we propose a new solvation free energy function for proteins, and test its applicability in predicting the solvation free energies of dipeptides. The present solvation model involves the improvement of the previous solvent-contact model assuming that the molecular solvation free energy could be given by the sum over the individual atomic contributions. In addition to the existing solvent-contact term, the modified solvation free energy function includes the self-solvation term that reflects the effects of intramolecular interactions in the solute molecule on solute–solvent interactions. Four kinds of atomic parameters should be determined in this solvation model: atomic fragmental volume, maximum atomic occupancy, atomic solvation, and atomic self-solvation parameters. All of these parameters for 16 atom types are optimized with a standard genetic algorithm in such a way to minimize the difference between the solvation free energies of dipeptides obtained from high-level quantum chemical calculations and those predicted by the solvation free energy function. The solvation free energies of dipeptides estimated from the new solvation model are in good agreement with the quantum chemical results. Therefore, the optimized solvation free energy function is expected to be useful for examining the structural and energetic features of proteins in aqueous solution.

© 2013 Elsevier Inc. All rights reserved.

1. Introduction

Because most biological processes occur in aqueous solution, protein–solvent interactions have a significant effect on the protein–protein and protein–ligand interactions as well as on the structural stability of proteins. The role of water molecules in optimizing the structure and function of proteins is not surprising because more than 30% of amino acid residues are ionizable [1]. Actually, the strength of interactions with solvent molecules affects the ionization equilibria for the titratable residues, which in turn determine the long-range electrostatic forces associated with intramolecular and intermolecular interactions of proteins. Therefore, the solvation free energy of a protein needs to be estimated with accuracy to understand its structural and functional properties in solution. Despite the necessity for a precise description for protein solvation, current experimental techniques such as osmotic stress [2] and far-infrared laser vibration–rotation tunneling spectroscopy [3] have provided only limited information about protein–solvent interactions.

Complementary to the experimental techniques, computational methods have drawn a particular interest as a tool for coping

with protein solvation problems because they can describe the solute–solvent interactions directly from a molecular perspective [4,5]. However, the solvation free energy has been considered one of the most calculation-difficult energy terms due to the complexity in solvent–solute interactions [6]. Although the explicit solvent models may provide a straightforward way to calculate the solvation free energies of proteins, a high computational cost has made it difficult for them to be employed in practical applications. Therefore, various implicit solvation models with high efficiency have been developed as alternatives, which includes solvent-accessible surface area (SASA) model [7–9], appropriately defined first solvation shell model [10], and modified generalized Born model [11]. Poisson–Boltzmann (PB) equation approaches have also been applied in investigating the solvation effects in a realistic way [12]. However, a high computational cost for the finite-difference algorithm has limited the applicability of PB model. To reduce the computational burden, therefore, the potential of mean force methods have also been developed to estimate the solvation free energy of proteins at the expense of some accuracy [13].

In the early 1990s, Stouten et al. proposed a solvation free energy function for proteins based on the solvent-contact model [14,15]. Under the assumption that the solvation free energy of proteins could be given by the summation over the atomic contributions, they optimized the solvation free energy function by classifying all of the atoms in amino acids into six atom types. Although this simple solvation model proved to be very successful in estimating

^{*} Corresponding author. Tel.: +82 2 3408 3766; fax: +82 2 3408 4334.

E-mail address: hspark@sejong.ac.kr (H. Park).

¹ Both authors equally contribute to this work.

the structural properties of proteins in solution, some modifications were required to enhance the accuracy of the solvation free energy function. In the previous work, we improved Stouten et al.'s solvation model to obtain a solvation free energy function suitable for amino acids and organic molecules. By defining 16 and 23 atom types, we were able to obtain the solvation free energy functions that could predict the experimental solvation free energies of amino acids and organic compounds with the associated squared correlation coefficients (r^2) of 0.94 and 0.86, respectively [16,17].

Actually, most of previous theoretical investigations for protein solvation were based on the group additivity that assumed a linear relationship between the solvation free energy and the volume of hydration shell. However, this assumption has been inappropriate to deal with the solvation problems involving strong intramolecular interactions between solute atoms [18–20], which is called the self-solvation effects. In the present study, we examine these self-solvation effects on protein solvation with the aim to propose an accurate solvation free energy function that can reflect the nonlinearity inherent in protein–solvent interactions. More specifically, a new solvation free energy function for protein involving the self-solvation effects are established and tested for predicting the solvation free energies of dipeptides. This new solvation free energy function is expected to be useful for investigating structural and energetic features of proteins in aqueous solution.

2. Theory and computational methods

Construction of the new solvation free energy function: The solvent-contact model to calculate the solvation free energy of proteins is based on several assumptions. First, the solvation free energy (ΔG_{sol}) can be approximated by the sum of individual atomic contributions as follows.

$$\Delta G_{\text{sol}} = \sum_i^{\text{atoms}} \Delta G_{\text{sol}}^i \quad (1)$$

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

$$\Delta G_{\text{sol}}^i = S_i F_i \quad (2)$$

Third, the atomic degree of exposure (F_i) can be obtained by measuring the unoccupied volume around the atom of interest. The occupied volume around the atom i (O_i) indicates the region to which the approach of solvent molecule is forbidden due to the occupation by the rest of protein atoms. O_i can be determined by summing the atomic volume parameters (V_j 's) representing the fragmental volumes of all the other atoms multiplied by a suitable envelope function, $E(r_{ij})$, with respect to the distance between the centers of atoms i and j .

$$O_i = \sum_{j \neq i}^{\text{atoms}} V_j E(r_{ij}) \quad (3)$$

Here, the Gaussian envelope function is employed with the variable r_{ij} representing the interatomic distance and the σ value of 3.5 Å. Although the screened Coulomb potential was also proposed as an envelope function, it was excluded in this study because its performance was shown to be worse than the Gaussian function [17]. Because F_i is the difference between the maximum occupancy of atom i (O_i^{max}) and O_i [15], the solvation free energy of a protein can be expressed in the following form:

$$\Delta G_{\text{sol}} = \sum_i^{\text{atoms}} S_i \left(O_i^{\text{max}} - \sum_{j \neq i}^{\text{atoms}} V_j e^{-\frac{r_{ij}^2}{2\sigma^2}} \right). \quad (4)$$

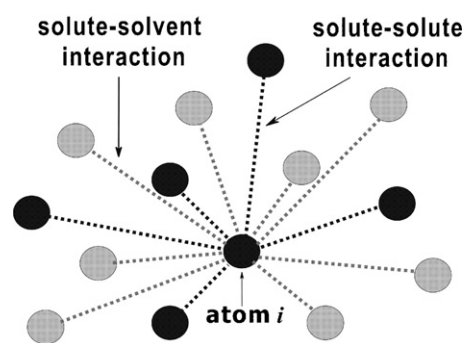


Fig. 1. Schematic diagram for the interactions of a protein atom i in solution. The black and gray circles indicate protein and solvent atoms, respectively. In this example, the atom i interacts with eight solvent atoms and five the other protein atoms.

Although this solvation model proved to be successful in predicting the solvation free energies of small organic molecules, its applicability to the problem of protein solvation remains unclear due to the neglect of self-solvation effects. Indeed, it has been demonstrated that the effects of the intramolecular non-bonded interactions should be reflected in the solvation free energy function to properly describe the structural and functional properties of large molecules such as proteins in solution [18–20]. For example, the intramolecular hydrogen bond and van der Waals interactions established in the occupied volume can affect the strength of the protein–solvent interactions through the change in electron distribution in the outer region exposed to bulk solvent. The presence of these self-solvation effects in protein solvation can thus be attributed to the intramolecular stabilization/destabilization of the atoms in the solvent-exposed region by the atoms in the occupied volume.

Fig. 1 exemplifies a typical pattern for the interactions of protein atoms in solution. As illustrated, a solute atom can be stabilized not only by the interactions with solvent molecules (solvation) but also by those with the rest of solute atoms (self-solvation). Therefore, the solvation energy function in Eq. (4) should be insufficient to fully describe the stabilization of a protein in solution because it contains a protein–solvent interaction term only and lacks the self-solvation term. To define the self-solvation energy of a protein in solution, we assume that it can be obtained by the summation of all atomic contributions. The self-solvation energy of an atom i (ΔG_{self}^i) can then be approximated by the product of the atomic self-solvation parameter (P_i) and the occupied volume around the atom i as follows.

$$\Delta G_{\text{self}}^i = P_i O_i \quad (5)$$

Here, P_i describes the stabilization energy of a protein atom i per unit volume due to the intramolecular interactions with the other protein atoms in solution. This self-solvation term can be appended to the solvation term in Eq. (4) to obtain the improved solvation free energy function for proteins, which can be expressed in the following form.

$$\Delta G_{\text{sol}} = \sum_i^{\text{atoms}} \left\{ S_i \left(O_i^{\text{max}} - \sum_{j \neq i}^{\text{atoms}} V_j e^{-\frac{r_{ij}^2}{2\sigma^2}} \right) + P_i \sum_{j \neq i}^{\text{atoms}} V_j e^{-\frac{r_{ij}^2}{2\sigma^2}} \right\} \quad (6)$$

The first and second terms in Eq. (6) correspond to the contribution from protein–solvent interactions and that from the intramolecular interactions between protein atoms to the stabilization of the protein in solution, respectively. Thus, this new solvation free energy function places an emphasis on the fact that a protein solute can be stabilized in solution as a consequence of the coordination between the protein–solvent interactions and the intramolecular interactions among the protein atoms. The four key

atomic parameters in the present solvation model include the maximum atomic occupancy (O_i^{\max}), atomic fragmental volume (V_j), and the atomic solvation (S_i) and self-solvation (P_i) energies per unit volume. Actually, the extent of contribution from the intra-molecular interactions to protein solvation can be determined by P_i parameters. The four kinds of atomic parameters should thus be optimized for all possible atom types to estimate the solvation free energy of proteins. This solvation model is similar to the morphometric approach [21] in the sense that protein solvation free energy can be estimated directly from the potential energy function and four parameters. It should be noted that the number of parameters can be reduced to three including V_j , $S_i O_i^{\max}$, and $(S_i - P_i)$ if Eq. (6) is rearranged effectively. Nonetheless, we treated S_i , O_i^{\max} , and P_i as the independent parameters to evaluate the extents of the contributions of both solvation and desolvation terms to the molecular solvation free energy.

Data set: Unlike organic molecules, it is very difficult to collect the experimental solvation free energy data for proteins that are necessary for the optimization of the parameters in Eq. (6). Even the experimental data for polypeptides (di-, tri-, and tetrapeptides, for example) are unavailable at present due in a large part to the difficulty in the precise measurement of the vapor pressures of polypeptides that are required to convert the solubility data to solvation free energy energies [22]. The reference dataset should therefore be constructed prior to the optimization of the atomic parameters in the solvation free energy function. In the absence of experimental data, the protein solvation free energies obtained from high-level quantum chemical calculations can be alternative for the reference dataset to complete the solvation free energy function. Indeed, the solvation free energies calculated with quantum chemical solvation model such as polarizable continuum model (PCM) and its variants [23] proved to be in a good agreement with the experimental results [24,25]. In the present study, we limited the scope of our interest to the solvation free energy of dipeptides due to the high computational cost for the PCM calculations of polypeptides.

To select the best computational method to obtain the reference dataset for dipeptides, we compared the solvation free energies of amino acids computed with various PCM methods based on density functional theory (DFT) calculations and those measured from dynamic vapor pressure measurements [26]. The optimization of the atomic parameters in Eq. (6) were preceded by the calculation of solvation free energies of dipeptides with the selected DFT-based PCM method to construct the reference dataset. The neutral form of each amino acid was considered only in the formation of dipeptides because the solvation free energy of the ionized forms had inevitably been overestimated in quantum chemical calculations [23]. In case of histidine, the two isomeric forms involving a hydrogen atom at N_δ or N_ϵ positions are considered in the construction of dipeptides. Atomic coordinates of all possible 441 dipeptides were then extracted from protein data bank (PDB), which were in turn divided into 401 and 40 dipeptides at random to construct the training and test sets, respectively. The test set was constructed by the random selection of 40 dipeptides in such a way that all 21 amino acids and 16 atom types were included. Solvation free energies of all 441 dipeptides were then calculated at B3LYP/6-31G* level with PCM solvation model to obtain the reference dataset.

Definition of atom types: Different atom types should have different contributions to solvation free energy in the present solvation model. We used 16 basic atom types for the elements commonly found in amino acids. The atom type of a given atom in a dipeptide was thus differentiated according to element, hybridization state, and chemical environment around the atom. Considering the portability and simplicity for implementing the classifications, we designated all atom types in the same fashion as in Sybyl MOL2 format. Some MOL2 atom types needed to be split

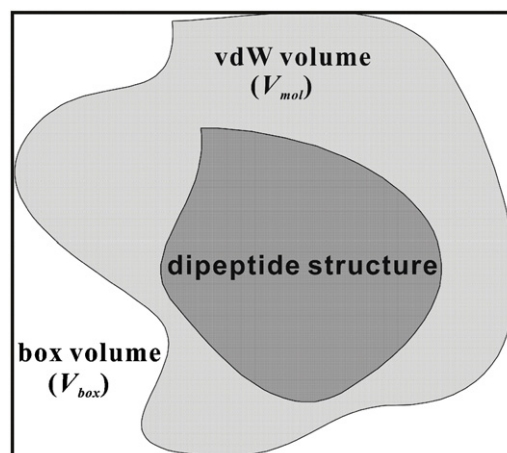


Fig. 2. Embedment of a dipeptide in a box to calculate its total volume (V_{mol}).

into a few subtypes on the basis of chemical environment around the atoms. For example, we subdivided the atom type of hydrogen depending on the atom to which it was attached.

Optimization of atomic volume parameters with genetic algorithm: As mentioned above, four atomic parameters should be determined for a given atom type in order to complete the solvation free energy function. Among them, the atomic volume parameter V_j represents the fragmental volume of atom j in a dipeptide. Because these V_j values exhibited a bad convergent behavior in the simultaneous optimization of the four kinds of parameters, they were optimized separately with the genetic algorithm as detailed below. To optimize the V_j parameters, the van der Waals volume (V_{mol}) and the sum of V_j values were assumed to be the real volume and the estimated volume of a dipeptide, respectively. Therefore, the determination of the V_{mol} value of each dipeptide is prerequisite for the optimization of V_j parameters for varying atom types. For this purpose, each dipeptide was placed in a 3-D box as illustrated in Fig. 2. The length, width, and height of this box correspond to the maximum distances along the three axes defining the coordinate system of the van der Waals volume of the dipeptide. To construct the molecular van der Waals volumes of dipeptides, the atomic radii of carbon, nitrogen, oxygen, sulfur, and hydrogen atoms are set equal to 1.53, 1.45, 1.36, 1.70, and 1.08, respectively. Monte Carlo simulations involving the random selections of a point in the predefined 3-D box were then carried out to calculate the V_{mol} value of the dipeptide embedded in the box. In this simulation, V_{mol} could be obtained by the product of the box volume (V_{box}) and the ratio of the number of trials to select a point in the van der Waals volume (N_{hits}) to the total number of trials (N_{trials}). Thus, we have

$$V_{\text{mol}} = V_{\text{box}} \times \frac{N_{\text{hits}}}{N_{\text{trials}}} \quad (7)$$

The van der Waals volumes of dipeptides in the test set calculated with Eq. (7) were similar to those calculated with the SYBYL program of version 8.1.1 with the associated r^2 value of 0.99 and mean relative deviation of 11.8 Å³.

With the calculated V_{mol} values in hand, the atomic volume parameters were optimized with the standard genetic algorithm. This started with the definition of a generation defined with 100 vectors comprising the V_j parameters for all atom types, which was followed by the removal of 50 with a bias toward preserving the most fit with the lowest error. The total number of vectors was selected to be 100 as the compromise between the convergence and the extensiveness of parameter space. 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

Table 1Solvation free energies (in kcal/mol) of amino acids calculated with PCM, CPCM, and SCIPCM methods at the B3LYP/6-31G* level in comparison with the experimental data.^a

Amino acid	$\Delta G_{\text{sol}}^{\text{exp}}$	$\Delta G_{\text{sol}}^{\text{PCM}}$	$\Delta G_{\text{sol}}^{\text{CPCM}}$	$\Delta G_{\text{sol}}^{\text{SCIPCM}}$	d^{PCM}	d^{CPCM}	d^{SCIPCM}
Ala	1.94	−0.06	−0.07	−0.12	1.99	2.00	2.06
Asn	−9.54	−12.4	−15.27	−9.76	2.87	5.74	0.23
Cys	−1.22	−4.75	−5.01	−2.81	3.53	3.79	1.59
Gln	−9.25	−11.86	−14.58	−9.02	2.62	5.34	0.22
Ile	2.10	−0.25	−0.33	−0.24	2.35	2.43	2.35
Leu	2.25	−0.31	−0.41	−0.29	2.56	2.66	2.53
Met	−1.46	−3.33	−3.42	−2.37	1.87	1.96	0.91
Phe	−0.74	−3.33	−3.95	−1.80	2.59	3.21	1.06
Ser	−4.97	−6.31	−7.89	−4.47	1.34	2.92	0.49
Thr	−4.80	−6.71	−8.47	−4.43	1.91	3.67	0.37
Trp	−5.81	−9.88	−11.07	−5.98	4.08	5.27	0.18
Tyr	−6.02	−10.35	−11.98	−6.09	4.33	5.96	0.07
Val	1.96	−0.19	−0.25	−0.13	2.15	2.21	2.09
Arg	−10.71	−15.11	−20.04	−10.91	4.41	9.34	0.21
Lys	−4.30	−4.83	−5.82	−3.04	0.53	1.52	1.26
His	−10.06	−13.6	−15.15	−8.10	3.54	5.09	1.95
Asp	−6.57	−10.41	−12.4	−6.79	3.84	5.83	0.22
Glu	−6.35	−10.24	−12.19	−6.47	3.89	5.84	0.12

^a d^{PCM} , d^{CPCM} , and d^{SCIPCM} are given by $|\Delta G_{\text{sol}}^{\text{exp}} - \Delta G_{\text{sol}}^{\text{PCM}}|$, $|\Delta G_{\text{sol}}^{\text{exp}} - \Delta G_{\text{sol}}^{\text{CPCM}}|$, $|\Delta G_{\text{sol}}^{\text{exp}} - \Delta G_{\text{sol}}^{\text{SCIPCM}}|$, respectively.

elements of another vector of the top 50. The 50 new vector created in these ways were then evaluated together with the top 50. This cycle was repeated as many times as desired. To evaluate the 100 vectors, we used the error hypersurface (F_V) defined by the sum of the absolute values of the differences between the calculated V_{mol} value and the sum of V_j values of the dipeptide.

$$F_V = \sum_k^{\text{dipeptides}} \left| V_{\text{mol}}^k - \sum_j^{\text{atoms}} V_j \right|. \quad (8)$$

Optimization of atomic solvation and self-solvation parameters: In addition to V_j , three remaining atomic parameters (S_i , O_i^{max} , and P_i) in Eq. (6) should be determined for each atom type to obtain the complete form of solvation free energy function for dipeptides. These parameterizations were carried out by operating the genetic algorithm with the same procedure as in the optimization of V_j parameters. To optimize the parameters, the error hypersurface was defined by the sum of the absolute values of the differences between the solvation free energies calculated with the DFT-based PCM calculations (ΔG_{QM}^k) and those obtained with Eq. (6) (ΔG_{calc}^k). This fitness function can be written as

$$F_S = \sum_k^{\text{dipeptides}} |\Delta G_{\text{QM}}^k - \Delta G_{\text{calc}}^k|. \quad (9)$$

Actually the O_i^{max} parameters in Eq. (6) were assumed to be the same as those obtained in the parameterization of Eq. (4) because of a bad convergent behavior of the F_S value in the simultaneous optimization of O_i^{max} , S_i , and P_i parameters. During the operation of the genetic algorithm, the F_S value became convergent to 0.276 kcal/mol after around 10,000 iterations.

3. Results and discussion

To select the best computational method to obtain the reference dataset for solvation free energies of dipeptides, DFT calculations at B3LYP level of theory involving various PCM models and basis sets were tested for how accurately they could reproduce the experimental solvation free energies of amino acids. Three quantum mechanical solvation models (PCM, CPCM, and SCIPCM) and nine basis sets were taken into account in these test calculations including 6-31G, 6-31G*, 6-31G**, 6-31+G*, 6-31+G**, 6-311G*, 6-311+G*, 6-311+G**, and 6-311++G*. The 6-31G* basis set produced the solvation free energies of amino acids most similarly to the

experimental results in all three cases of quantum mechanical solvation models under consideration. Table 1 lists the solvation free energies of amino acids obtained with PCM, CPCM, and SCIPCM calculations at the B3LYP/6-31G* level of theory in comparison with the experimental data. The SCIPCM results appear to compare most reasonably with the experimental data with the mean absolute deviation of 1.00 kcal/mol as compared to 2.80, and 4.15 kcal/mol for PCM and CPCM results, respectively. Despite the best agreement with experimental data for the energy values, however, the differences in solvation free energies between amino acids seem to be substantially underestimated in some cases of SCIPCM results. For example, the difference between lysine and methionine amounts to only 0.67 kcal/mol as compared to 2.84 and 1.50 kcal/mol in the experimental data and in PCM results, respectively. Further evaluations are therefore required to choose the most appropriate quantum mechanical solvation model to construct the reference dataset for solvation free energies of dipeptides.

Fig. 3 shows the correlation diagrams between the experimental solvation free energies of amino acids and those obtained from the three different PCM calculations performed at B3LYP/6-31G* level of theory. We see that the solvation free energies estimated with the three computational methods compare reasonably well with the corresponding experimental data. Among them, the PCM results exhibit the best fit with the experimental solvation free energies of amino acids with the associated r^2 value of 0.962. Furthermore, the correlation diagram between experimental and PCM results has the slope closest to one (0.867), which indicates that the differences in solvation free energies between amino acids can be estimated most accurately by PCM calculations. On the basis of these validation results, we selected the PCM method in combination with B3LYP/6-31G* level of theory to obtain the solvation free energies of 441 dipeptides that should serve as the yardsticks to evaluate the accuracy of the predicted solvation free energies with Eq. (6).

Prior to the calculation of solvation free energies of dipeptides, their geometries were fully optimized at B3LYP/6-31G* level of theory from the initial structures extracted from the high-resolution (<1.8 Å) crystal structures of various proteins in PDB. Each of these 441 optimized dipeptide structures served as an input for the calculation of solvation free energies at B3LYP/6-31G* level of theory with the PCM method during which the structures of dipeptides were reoptimized. The calculated solvation free energies range from −55.6 to −19.6 kcal/mol, the maximum and minimum of which were found in Arg-Arg and Val-Pro dipeptides, respectively. It is noteworthy that the absolute values of minimum and

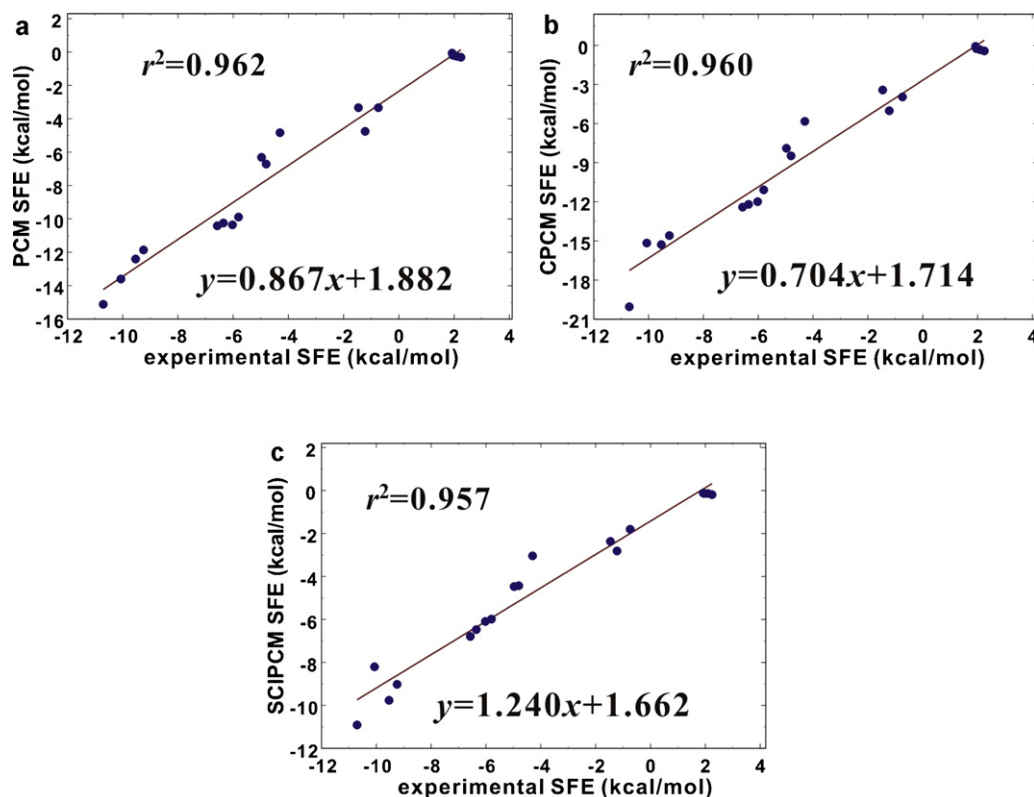


Fig. 3. Correlation diagrams for the experimental solvation free energies (SFEs) of amino acids and those calculated from B3LYP/6-31G* level of theory based on (a) PCM, (b) CPCM, and (c) SCIPCM methods.

maximum solvation free energies of dipeptides appear to be larger than the twice of those for amino acids, which is consistent with the nonlinearity in protein–solvent interactions.

With the calculated solvation free energy data in hand, we first evaluate the previous solvation model in Eq. (4) that neglects the self-solvation effects. This solvation model was shown to be useful in predicting the solvation free energies of small organic molecules in the previous studies [16,17]. We now address its applicability to the problem of protein solvation in the present study. Listed in Table 2 are the optimized atomic volume (V_j), maximum atomic occupancy (O_i^{\max}), and atomic solvation parameters (S_i) in Eq. (4) for the 16 atom types that were optimized with 400 dipeptides in the training set. It should be noted that the number of atom types decreases from 23 for organic compounds [16] to 16 for dipeptides in this study because some atom types (for example, halogens and

the atoms involving triple bonds) can be excluded in describing the amino acids. Therefore, the applicability of these atomic parameters should be limited to amino acids and polypeptides because no organic molecule was included in the training set. It is also noteworthy that various atom types have similar O_i^{\max} values because they were implicitly optimized in the form of $S_i O_i^{\max}$. On the other hand, the optimized V_j values exhibit a large difference with varying atom types even in the case of the same element. For example, the V_j value of trigonal planar nitrogen (N.pl3) appears to be larger than that of sp^3 nitrogen (N.3) by a factor of two. Such a substantial difference in V_j parameters between the atoms with similar atomic radii is actually not surprising because each V_j value represents the average of atomic contributions with type j to the van der Waals volumes of the dipeptides with various molecular sizes and conformations.

Table 2
Optimized atomic fragmental volume (V_j), maximum atomic occupancy (O_i^{\max}), and atomic solvation parameters (S_i) in the solvation model without self-solvation effects.

Atom type	Description	V_j (\AA^3)	O_i^{\max} (\AA^3)	S_i (kcal/mol \AA^3)
C.3	sp^3 carbon	10.159	367.1	3.365
C.2	sp^2 carbon	8.571	365.9	2.143
C.ar	Aromatic carbon	8.995	361.9	0.048
C.cat	Carbocation	22.857	365.4	−0.476
N.3	sp^3 nitrogen	5.397	356.5	−13.889
N.2	sp^2 nitrogen	6.984	346.4	−22.381
N.am	Amidic nitrogen	5.397	333.8	−28.571
N.pl3	Trigonal planar nitrogen	10.357	334.6	−22.381
O.3	sp^3 oxygen in hydroxyl group	8.571	330.4	−11.231
O.2	sp^2 oxygen	9.365	324.5	−12.857
O.co2	Carboxylate oxygen	10.952	361.2	−15.238
S.3	sp^3 sulfur	10.556	392.2	−4.762
H.1	Hydrogen bonded to carbon	2.715	321.8	1.429
H.2	Hydrogen bonded to nitrogen	3.205	320.0	−5.281
H.3	Hydrogen bonded to oxygen	3.142	320.0	−7.460
H.4	Hydrogen bonded to sulfur	1.571	320.2	−1.745

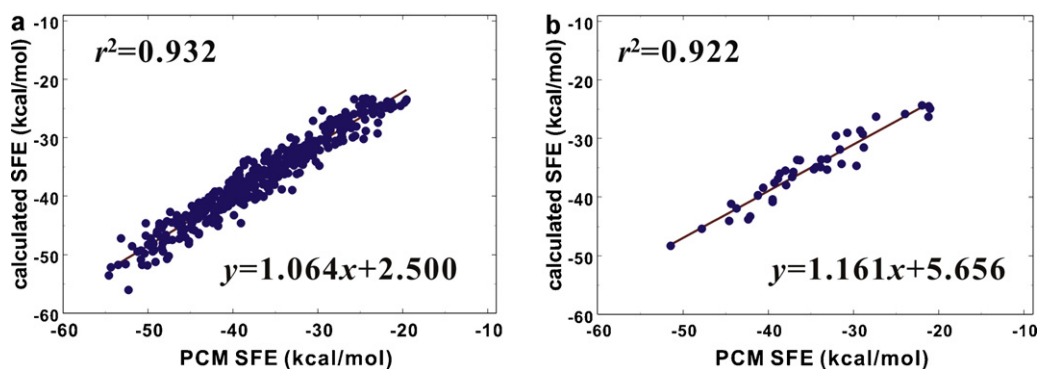


Fig. 4. Correlation diagrams for PCM solvation free energies (SFEs) versus those obtained with the solvation free energy function without self-solvation term for (a) 401 dipeptides in the training set and (b) 40 dipeptides in the test set. All energy values are given in kcal/mol.

Actually, the V_j parameters are expected to become more diverse by the splitting of the atom types according to the chemical environments. For example, a central sp^3 carbon atom can be subdivided to C.3.1, C.3.2, and C.3.3, and C.3.4 according to the number of substituents. The atomic parameter space extended in this way seems to be capable of discriminating the atoms with different steric hindrances for the accessibility of solvent molecules. In this regard, we find that the V_j parameters of C.3.1, C.3.2, and C.3.3, and C.3.4 are optimized to 11.341, 9.365, 7.977, and 6.342, respectively, which indicates the decrease of V_j values with the increase in the number of substituents at the central sp^3 carbon. This can be understood because a highly substituted solute atom should have a small solvent-exposed volume due to the increased steric hindrance by the neighboring groups.

The optimized S_i parameters exhibit a trend consistent with general atomic properties. We note, for example, that the neutral carbon and their attached hydrogen atoms have positive S_i values, which indicates the unfavorable interactions with solvent molecules. This is consistent with the immiscibility of hydrocarbons in water. It is also noteworthy that the S_i values for nitrogen and oxygen atoms get more negative as their atomic hybridization changes from sp^3 (N.3 and O.3) to sp^2 (N.2, O.2, and the similar planar atom types such as N.am, N.pl3, and O.co2). This indicates that the increase of s-character in the hybridization of atomic orbitals would have an effect of strengthening the attractive interactions with water molecules. This can be understood by noting the fact that the increase in the s-character of the hybrid orbitals of a central atom leads to the increase in the polarity of its chemical bonds, which culminates in the promotion of solute–solvent interactions in polar solvents. In case of hydrogen atoms, S_i values appear to get more negative in the order of electronegativity of the heavy atom to which the hydrogen atom of interest is attached. This is not surprising because the increase in electronegativity causes to raise the acidity of the central atom, which would have an effect of strengthening the hydrogen-bond interactions with solvent molecules.

The correlations between the solvation free energies calculated with the DFT-based PCM method and those obtained with Eq. (4) are illustrated in Fig. 4. With the test set comprising 40 dipeptides, we obtain the r^2 value of 0.922, which is a little worse than the fitting with the training set including 401 dipeptides (0.932). The solvation free energies of dipeptides obtained with Eq. (4) are thus found to be as accurate as those of amino acids calculated from free energy simulations with GROMOS96 force field [27]. The accuracy of the solvation model ignoring the self-solvation effects is also comparable to those of the quantitative structure–property relationship (QSPR) and the artificial neural network (ANN) models trained with a large number of molecules and descriptors [28,29]. The merit of the GA-based solvent-contact model in estimating the solvation free energies, when compared to QSPR and ANN methods,

lies in the direct use of 3-D structures of dipeptides in the optimization of atomic parameters instead of the molecular descriptors that should be calculated separately prior to the parameterization. However, the comparisons also indicate that the accuracy should be improved in order for the solvation free energy function to be used in estimating the solvation free energies of polypeptides in place of the existing methods.

We next turn to the new solvation model in which the self-solvation effects are taken into account. The atomic parameters optimized with genetic algorithm are summarized in Table 3. Because S_i values of the neutral carbon atoms remain positive, their interactions with solvent molecules are also expected to be repulsive in the modified solvation model. However, the negative values of their optimized atomic self-solvation parameters (P_i values) imply that even the neutral carbon atoms can make a contribution to the stabilization of dipeptides in aqueous solution. This stabilization effect should apparently be attributed to the attractive intramolecular hydrophobic interactions between nonpolar solute atoms. The involvement of such intramolecular hydrophobic interactions in stabilizing the solutes in solution was also implicated in the structural studies on the effects of inter-residue interactions on protein folding [30]. Despite the abundance of neutral carbon atoms, however, the low absolute values of their P_i parameters indicate that the intramolecular interactions between carbon atoms would be insufficient in themselves to be the major driving force to stabilize the proteins in solution.

It is also a common feature with the previous solvation model that most of nitrogen atoms have more negative S_i parameters than oxygen atoms, which indicates that the former would interact with solvent molecules in more attractive fashion than the latter. This is not surprising because nitrogen is generally a better hydrogen bond acceptor than oxygen because the relatively low electronegativity of the former makes its lone-pair electrons more unstable than those of the latter. On the other hand, most nitrogen atoms appear to have positive P_i values except for the trigonal planar nitrogen (N.pl3) in contrast to the negative values for all oxygen atoms. This implies that nitrogen atoms in dipeptides would be stabilized predominantly by the interactions with solvent molecules whereas oxygen atoms can be involved in both solute–solvent and intramolecular interactions to stabilize the dipeptides in solution. More noteworthy, the S_i parameters of the hydrogen atoms bonded to nitrogen and oxygen are found to be positive in the present solvation model while their corresponding P_i values appear to be negative. These results indicate the preference for the formation of intramolecular hydrogen bonds over the intermolecular solute–solvent ones in the case that a group in the dipeptide solute should play a role of hydrogen bond donor. With respect to the preference for the intramolecular hydrogen bonds in dipeptides, we find that only phenolic moiety is less basic than water among

Table 3Atomic fragmental volume (V_j), maximum atomic occupancy (O_i^{\max}), atomic solvation (S_i), and atomic self-solvation parameters optimized with genetic algorithm.

Atom type	Description	V_j (\AA^3)	O_i^{\max} (\AA^3)	S_i (kcal/mol \AA^3)	P_i (kcal/mol \AA^3)
C.3	sp ³ carbon	10.159	367.1	1.238	−2.127
C.2	sp ² carbon	8.571	365.9	2.778	−2.016
C.ar	Aromatic carbon	8.995	361.9	0.095	−1.920
C.cat	Carbocation	22.857	365.4	7.524	9.254
N.3	sp ³ nitrogen	5.397	356.5	−16.587	2.857
N.2	sp ² nitrogen	6.984	346.4	−26.825	1.763
N.am	Amidic nitrogen	5.397	333.8	−21.317	5.331
N.pl3	Trigonal planar nitrogen	10.357	334.6	−16.778	−5.175
O.3	sp ³ oxygen in hydroxyl group	8.571	330.4	−18.064	−3.651
O.2	sp ² oxygen	9.365	324.5	−15.048	−8.857
O.co2	Carboxylate oxygen	10.952	361.2	−12.133	−10.232
S.3	sp ³ sulfur	10.556	392.2	−10.242	−5.937
H.1	Hydrogen bonded to carbon	2.715	321.8	0.841	−1.529
H.2	Hydrogen bonded to nitrogen	3.205	320.0	0.095	−8.921
H.3	Hydrogen bonded to oxygen	3.142	320.0	0.540	−5.317
H.4	Hydrogen bonded to sulfur	1.571	320.2	2.417	1.032

the hydrogen bond accepting groups in the five hydrogen bonds shown in Fig. 5. This indicates that the remaining four can be maintained in aqueous solution, which is consistent with the difficulty of water molecules in playing the role of hydrogen bond donor with respect to the dipeptide groups. The roles of acceptor for the intramolecular hydrogen bonds seem to be played in a predominant part by oxygen and trigonal planar nitrogen atoms in dipeptides instead of sp³, sp², and amidic nitrogens, which can be inferred from their respective negative and positive P_i values. Thus, the optimized P_i values of polar atoms further exemplify the importance of intramolecular interactions in stabilizing the dipeptides in solution, and necessitate the inclusion of the self-solvation term in the solvation free energy function.

The limited role of solvent molecules in the hydrogen-bond stabilization of dipeptides is consistent with the fact that a strong hydrogen bond is more difficult to be established in solution than in the gas phase due to the role of rupturing or weakening the hydrogen bonds played by water molecules [31,32]. The extent of such a negative solvent effect should be greater in the intermolecular solute–solvent hydrogen bonds than in the intramolecular ones because the former is exposed to bulk solvent in a larger part than the latter. Indeed, the intramolecular hydrogen bonds have a better chance to be protected in solution than the solute–solvent hydrogen bonds due to the presence of the more neighboring solute

atoms that can limit the approach of solvent molecules. Fig. 5 shows the structures of Lys-Glu and Tyr-Trp dipeptides in aqueous solution obtained by the geometry optimization at B3LYP/6-31G* level of theory with the PCM solvation model. In the optimized Lys-Glu dipeptide, we see that the dipeptide can be stabilized in solution by establishing two strong intramolecular N···H···O hydrogen bonds. These intramolecular hydrogen bonds seem to make a larger contribution to the self-solvation effects than the intramolecular hydrophobic interactions, which can be inferred from the much more negative values of the P_i parameters of oxygen and N.pl3 atoms than those of carbon atoms (Table 3). Three intramolecular hydrogen bonds are observed in the optimized structure of Tyr-Trp dipeptide (Fig. 5b) in addition to the face-to-face intramolecular hydrophobic interactions between the side-chain phenyl and indole rings that would play the role of positioning the intramolecular hydrogen bonds. This exemplifies the involvement of both the intramolecular hydrogen bonds and the intramolecular van der Waals interactions in the stabilization of dipeptides in solution.

The correlations between the PCM solvation energies and those calculated using Eq. (6) are illustrated in Fig. 6. As in Fig. 4b, the correlation slope for the test set is close to 1 (1.130) while the intercept amounts to 2.114. This indicates that the solvation free energy values can be overestimated in the present solvation model to a similar

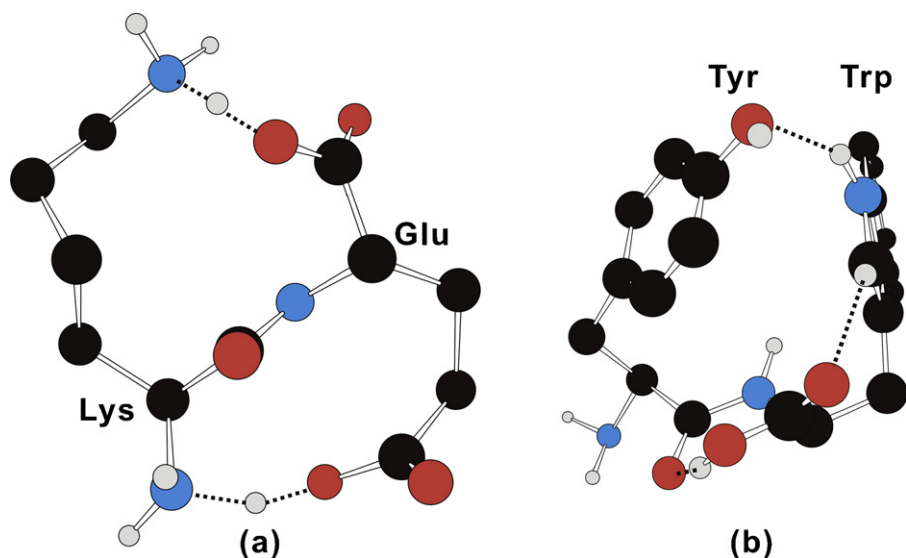


Fig. 5. The structures of (a) Lys-Glu and (b) Tyr-Trp dipeptides optimized at B3LYP/6-31G* level of theory based on the PCM solvation model. Each dotted line indicates a hydrogen bond. Hydrogen atoms attached to carbons are omitted for visual clarity.

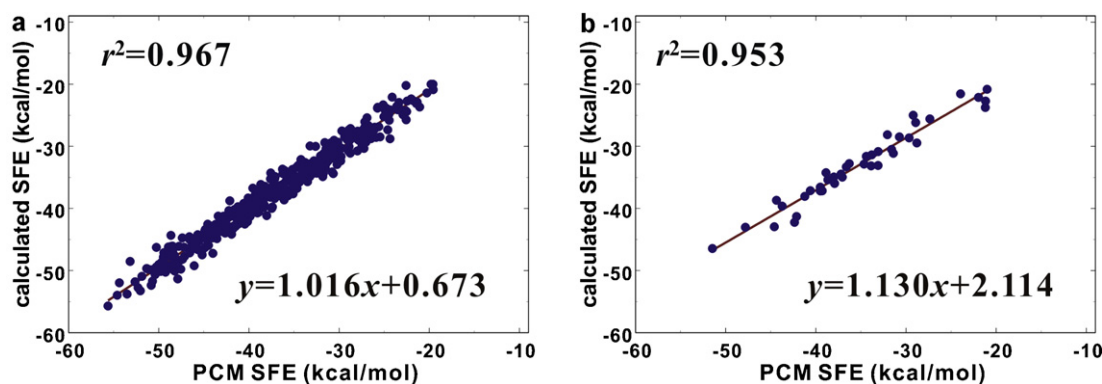


Fig. 6. Correlation diagrams for PCM solvation free energies (SFEs) versus those obtained with the solvation free energy function with self-solvation term for (a) 401 dipeptides in the training set and (b) 40 dipeptides in the test set. All energy values are given in kcal/mol.

extent for various dipeptides. We note that the r^2 values for training and test sets increase from 0.932 and 0.922 in the results obtained with the previous method (Fig. 4) to 0.967 and 0.953 in the present solvation model, respectively. Apparently, these improvements should be attributed to the inclusion of self-solvation term in the solvation free energy function. The accuracy of the present solvation model is actually comparable to those of the quantum chemical dielectric continuum solvation model (COSMO) [33] and molecular dynamics free energy perturbation method [34,35]. These comparisons indicate that our GA-based solvent-contact model would be more efficient in estimating the molecular solvation free energies than the sophisticated quantum chemical method and statistical simulations with all-atom models because the former can produce the solvation free energy values simply from a potential function. In comparison with the results for the previous solvation model that neglected the self-solvation effects, the largest improvements in solvation free energies are observed for the dipeptides involving the amino acids that can establish the strong intramolecular hydrogen bonds or intramolecular van der Waals contacts. For example, the differences between the estimated solvation free energies with the potential function and the PCM results for Arg-Asp and Tyr-Trp dipeptides decrease from 3.64 and 3.23 kcal/mol in the previous solvation model to 0.02 and 0.03 kcal/mol, respectively, due to the inclusion of the self-solvation term in the solvation free energy function. These substantial improvements further exemplify the importance of intramolecular interactions to stabilize the dipeptides in solution, which was also proposed for the structural stability of proteins in solution [30].

To further compare the accuracies of the solvation free energy functions with and without the self-solvation term, we calculated the mean absolute deviation (MAD), maximum absolute deviation (XAD), mean relative absolute deviation (MRAD), and maximum relative absolute deviation (XRAD) values between the solvation free energies of the dipeptides in the test set obtained with the PCM method and those calculated using Eqs. (4) and (6). As shown in Table 4, MAD and MRAD values decrease from 1.87 and 5.67 to 1.13 kcal/mol and 3.35%, respectively, due to the inclusion of the self-solvation term in the solvation free energy function. These results confirm the necessity for the self-solvation term to enhance the accuracy of solvation free energy function.

To address the possibility of the present solvation model being used for organic molecules, we examined the accuracy of the solvation free energy function in Eq. (6) using the 26 molecules in a druglike dataset and in a library of polychlorinated benzenes, which had also been used in the validation of the reference interaction site models (RISMs) for solvation [36,37]. Missing atomic parameters for these organic compounds were estimated from those of similar atom types in Table 3. As shown in Fig. 7, the r^2 value between

the experimental and calculated solvation free energies amounts to 0.838. The decrease in the r^2 value for the organic compounds is not surprising because the 16 atom types used in this study should be insufficient to describe all their atoms that have more complex chemical environment than amino acids. The accuracy of the present solvation model in predicting the solvation free energies of organic compounds seems to be enhanced in a straightforward way by subdividing the atom types according to chemical environment around the atoms in the molecules. Thus, the solvation free energy function in Eq. (6) is expected to be also useful for predicting the solvation free energies of organic compounds after the extension and optimization of the atomic parameters.

As an additional benchmark of the present solvation model, we also calculated the solvation free energy of the native structure of the immunoglobulin binding domain of streptococcal protein G comprising 56 amino acids [38] (PDB entry: 2GB1) using Eq. (6). The calculated solvation free energy amounts to -1075.9 kcal/mol, which is similar to that (-1088.4 kcal/mol) obtained by thermodynamic analysis of protein folding in aqueous solution [39]. Therefore, the present solvation model is likely to be also useful for predicting the solvation free energies of proteins.

The earlier solvation models such as solvent accessible surface area (SASA), hydrophobicity scales, and group additivity models proved to be incapable of explaining the solvation properties of proteins due to the assumption that the stability of proteins in solution should be determined by the extent of the solvent-exposed regions. This linearity criterion for solute–solvent interactions is actually inapplicable to proteins because the buried regions of a protein can also contribute to its structural stability in solution. Our modified

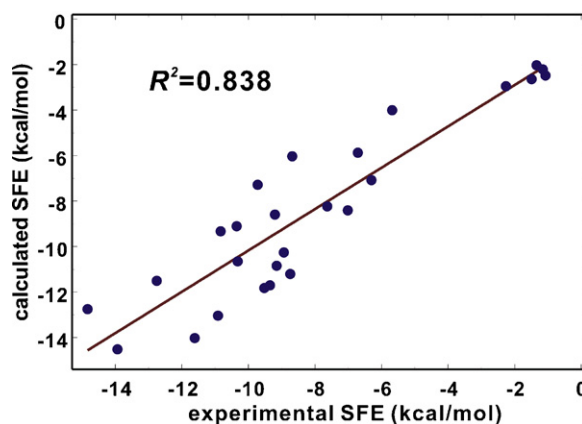


Fig. 7. Correlation diagram between the experimental and calculated solvation free energies for 26 organic compounds.

Table 4
Mean absolute deviation (MAD), maximum absolute deviation (XAD), mean relative absolute deviation (MRAD), and maximum relative absolute deviation (XRAD) between the solvation free energies of the dipeptides in the test set obtained with the PCM method and those calculated with the two solvation models.

Solvation model	MAD (kcal/mol)	XAD (kcal/mol)	MRAD (%)	XRAD (%)
Solvation term only	1.87	5.09	5.67	19.37
Solvation and self-solvation terms	1.13	5.67	3.35	18.86

solvent-contact model confirmed that the solvation free energies of dipeptides could be estimated by combining the contributions from solvent-exposed and self-solvation regions, the relative importance of which should be dependent on their conformations in solution. The significant contribution of self-solvation effects in solvation free energies of dipeptides is thus consistent with the nonlinearity in protein-solvent interactions that reflects the significant role of intramolecular interactions in protein solvation.

Despite the improved accuracy in estimating the solvation free energies of dipeptides, some problems still remain to be solved for the present solvation model to be extended to cope with the protein-solvent interaction problems. First, it is difficult to determine the atomic parameters of some atom types such as C.cat (CZ atom of arginine) and H.4 (HG atom of cysteine) with accuracy due to the rarity in amino acids. Those parameters may be improved further in polypeptides because training and test sets can be constructed in such a way to include a large number of polypeptides involving Arg and Cys residues. Second, conformational diversity of proteins should be considered in the parameterization because the volumes of solvent-exposed and buried regions can vary with the conformational change of proteins. For this purpose, molecular dynamics or Monte Carlo simulations can be applied prior to the parameterization to collect various local structural minima of proteins. Finally, the solvation free energy function needs to be decomposed into enthalpy and entropy terms. Because both thermodynamic quantities are experimentally accessible, the potential parameters in the enthalpic and entropic terms can be optimized independently using their respective corresponding experimental data. Apparently, this dual parameterization warrants the better correlation between the experimental and computational solvation free energies than the single parameterization because more diverse experimental data can be included in reference dataset. Because the sign of solvation free energy is determined by the combination of enthalpic and entropic contributions, the decomposition analysis of solvation free energy can also provide thermodynamic insight into the mechanism of protein solvation. Future studies for protein solvation will focus on further improvement in the accuracy of solvation free energy function with the three above-mentioned points kept in mind.

4. Conclusions

We have shown the outperformance of a solvent-contact model modified with the self-solvation term in predicting the solvation free energies of dipeptides. Total 64 atomic parameters for 16 atom types could be optimized with the standard genetic algorithm using 3-D atomic coordinates of 401 dipeptides and their solvation free energies obtained from quantum chemical calculations at B3LYP/6-31G* level of theory with PCM solvation model. As a consequence of the addition of the self-solvation term to the original solvation free energy function, the r^2 values between the solvation free energies estimated with the potential function and those obtained from PCM calculations increase from 0.932 and 0.922 to 0.967 and 0.953 for training and test sets, respectively. Such an enhancement in accuracy confirms the significance of intramolecular interactions in the stability of dipeptides in solution, which were also implicated in the experimental studies on protein-solvent interactions. Considering the simplicities in energy calculation and model refinement,

we expect that the present solvation model can be extended to cope with protein solvation problems.

Acknowledgments

This work was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (2012-0008440).

Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at <http://dx.doi.org/10.1016/j.jmgm.2013.02.006>.

References

- [1] B. Rost, C. Sander, Prediction of protein secondary structure at better than 70% accuracy, *Journal of Molecular Biology* 232 (1993) 584–599.
- [2] V.A. Parsegian, R.P. Rand, N.L. Fuller, D.C. Rau, Osmotic stress for the direct measurement of intermolecular forces, *Methods in Enzymology* 127 (1986) 400–416.
- [3] K. Liu, J. Cruzan, R. Saykally, Water clusters, *Science* 271 (1986) 929–933.
- [4] C.Y. Hu, H. Kokubo, G.C. Lynch, D.W. Bolen, B.M. Pettitt, Backbone additivity in the transfer model of protein solvation, *Protein Science* 19 (2010) 1011–1022.
- [5] Y.N. Vorobjev, J.A. Vila, H.A. Scheraga, FAMBE-pH: a fast and accurate method to compute the total solvation free energies of proteins, *The Journal of Physical Chemistry B* 112 (2008) 11122–11136.
- [6] W.L. Jorgensen, E.M. Duffy, Prediction of drug solubility from structure, *Advanced Drug Delivery Reviews* 54 (2002) 355–366.
- [7] D. Eisenberg, A.D. McLachlan, Solvation energy in protein folding and binding, *Nature* 319 (1986) 199–203.
- [8] T. Ooi, M. Oobatake, Prediction of the thermodynamics of protein unfolding: the helix-coil transition of poly(L-alanine), *Proceedings of the National Academy of Sciences of the United States of America* 88 (1991) 2859–2863.
- [9] C.A. Schiffer, J.W. Caldwell, B.M. Stroud, P.A. Kollman, Inclusion of solvation free energy with molecular mechanics energy: alanyl dipeptide as a test case, *Protein Science* 1 (1992) 396–400.
- [10] Y.K. Kang, G. Nemethy, H.A. Scheraga, Free energies of hydration of solute molecules. 1. Improvement of the hydration shell model by exact computations of overlapping volumes, *The Journal of Physical Chemistry* 91 (1987) 4105–4109.
- [11] J.D. Thompson, C.J. Cramer, D.G. Truhlar, New universal solvation model and comparison of the accuracy of the SM5.42R, SM5.43R, C-PCM, D-PCM, and IEF-PCM continuum solvation models for aqueous and organic solvation free energies and for vapor pressures, *The Journal of Physical Chemistry A* 108 (2004) 6532–6542.
- [12] I. Klapper, R. Hagstrom, R. Fine, K. Sharp, B. Honig, Focusing of electric fields in the active site of Cu-Zn superoxide dismutase: effects of ionic strength and amino-acid modification, *Proteins* 1 (1986) 47–59.
- [13] S. Garde, G. Hummer, A. Garcia, L. Pratt, M. Paulaitis, Hydrophobic hydration: Inhomogeneous water structure near nonpolar molecular solutes, *Physical Review E* 53 (1996) R4310–R4313.
- [14] F. Colonna-Cesari, C. Sander, Excluded volume approximation to protein-solvent interaction: the solvent contact model, *Biophysical Journal* 57 (1990) 1103–1107.
- [15] P.F.W. Stouten, C. Frömmel, H. Nakamura, C. Sander, An effective solvation term based on atomic occupancies for use in protein simulations, *Molecular Simulation* 10 (1993) 97–120.
- [16] H. Kang, H. Choi, H. Park, Prediction of molecular solvation free energy based on the optimization of atomic solvation parameters with genetic algorithm, *Journal of Chemical Information Model* 47 (2007) 509–514.
- [17] J.H. Park, J.W. Lee, H. Park, Computational prediction of solvation free energies of amino acids with genetic algorithm, *Bulletin of the Korean Chemical Society* 31 (2010) 1247–1251.
- [18] W.C. Wimley, T.P. Creamer, S.H. White, Solvation energies of amino acid side chains and backbone in a family of host-guest pentapeptides, *Biochemistry* 35 (1996) 5109–5124.

- [19] G. König, S.J. Boresch, Hydration free energies of amino acids: why side chain analog data are not enough, *The Journal of Physical Chemistry B* 113 (2009) 8967–8974.
- [20] J. Chang, A.M. Lenhoff, S.I. Sandler, Solvation free energy of amino acids and side-chain analogues, *The Journal of Physical Chemistry B* 111 (2007) 2098–2106.
- [21] R. Roth, Y. Harano, M. Kinoshita, Morphometric approach to the solvation free energy of complex molecules, *Physical Review Letters* 97 (2006) 078101.
- [22] M.H. Abraham, J. Andonian-Haftvan, G.S. Whiting, A. Leo, R.S. Taft, Hydrogen bonding. Part 34. The factors that influence the solubility of gases and vapours in water at 298 K, and a new method for its determination, *Journal of the Chemical Society, Perkin Transactions 2* (1994) 1777–1791.
- [23] J. Tomasi, B. Mennucci, R. Cammi, Quantum mechanical continuum solvation models, *Chemical Reviews* 105 (2005) 2999–3094.
- [24] C. Cappelli, B. Mennucci, Modeling the solvation of peptides. The case of (s)-N-acetylproline amide in liquid water, *The Journal of Physical Chemistry B* 112 (2008) 3441–3450.
- [25] M. Gupta, E.F. da Silva, H.F. Svendsen, Modeling temperature dependency of amine basicity using PCM and SM8T implicit solvation models, *The Journal of Physical Chemistry B* 116 (2012) 1865–1875.
- [26] R. Wolfenden, L. Andersson, P.M. Cullis, C.C.B. Southgate, Affinities of amino acid side chains for solvent water, *Biochemistry* 20 (1981) 849–855.
- [27] A. Villa, A.E. Mark, Calculation of the free energy of solvation for neutral analogs of amino acid side chains, *Journal of Computational Chemistry* 23 (2002) 548–553.
- [28] A. Cheng, K.M. Merz Jr., Prediction of aqueous solubility of a diverse set of compounds using quantitative structure–property relationships, *Journal of Medicinal Chemistry* 46 (2003) 3572–3580.
- [29] R. Liu, S.S. So, Development of quantitative structure–property relationship models for early ADME evaluation in drug discovery. 1. Aqueous solubility, *Journal of Chemical Information and Computer Sciences* 41 (2001) 1633–1639.
- [30] M.M. Gromiha, S. Selvaraj, Inter-residue interactions in protein folding and stability, *Progress in Biophysics and Molecular Biology* 86 (2004) 235–277.
- [31] R.L. Baldwin, Desolvation penalty for burying hydrogen-bonded peptide groups in protein folding, *The Journal of Physical Chemistry B* 114 (2010) 16223–16227.
- [32] C.N. Pace, Polar group burial contributes more to protein stability than nonpolar group burial, *Biochemistry* 40 (2001) 310–313.
- [33] A. Klamt, F. Eckert, M. Diedenhofen, Prediction of the free energy of hydration of a challenging set of pesticide-like compounds, *The Journal of Physical Chemistry B* 113 (2009) 4508–4510.
- [34] D. Shivakumar, J. Williams, Y. Wu, W. Damm, J. Shelley, W. Sherman, Prediction of absolute solvation free energies using molecular dynamics free energy perturbation and the OPLS force field, *Journal of Chemical Theory and Computation* 6 (2010) 1509–1519.
- [35] D. Shivakumar, E. Harder, W. Damm, R.A. Friesner, W. Sherman, Improving the prediction of absolute solvation free energies using the next generation OPLS force field, *Journal of Chemical Theory and Computation* 8 (2012) 2553–2558.
- [36] D.S. Palmer, A.I. Frolov, E.L. Ratkova, M.V. Fedorov, Toward a universal model to calculate the solvation thermodynamics of druglike molecules: the importance of new experimental databases, *Molecular Pharmaceutics* 8 (2011) 1423–1429.
- [37] E.L. Ratkova, M.V. Fedorov, Combination of RISM and cheminformatics for efficient predictions of hydration free energy of polyfragment molecules: application to a set of organic pollutants, *Journal of Chemical Theory and Computation* 7 (2011) 1450–1457.
- [38] A.M. Gronenborn, D.R. Filpula, N.Z. Essig, A. Achari, M. Whitlow, P.T. Wingfield, G.M. Clore, A novel, highly stable fold of the immunoglobulin binding domain of streptococcal protein G, *Science* 253 (1991) 657–661.
- [39] T. Imai, M. Kinoshita, A. Kovalenko, F. Hirata, Theoretical analysis on changes in thermodynamic quantities upon protein folding: essential role of hydration, *Journal of Chemical Physics* 126 (2007) 225102.