

Parametrized Model for Aqueous Free Energies of Solvation Using Geometry-Dependent Atomic Surface Tensions with Implicit Electrostatics

Gregory D. Hawkins, Christopher J. Cramer,* and Donald G. Truhlar*

Department of Chemistry and Supercomputer Institute, University of Minnesota,
Minneapolis, Minnesota 55455-0431

Received: March 23, 1997; In Final Form: May 23, 1997[®]

We present a new model for predicting aqueous free energies of solvation based entirely on geometry-dependent atomic surface tensions. The model is especially suited for rapid estimations on large molecules or large sets of molecules. This method is designed to be employed with gas-phase geometries to obtain solvation free energies of organic molecules containing H, C, N, O, F, S, Cl, and Br. We parametrized the model by using a training set containing 235 neutral solutes with a variety of functional groups, and we achieve a mean unsigned error of 0.55 kcal/mol when the model is applied using gas-phase geometries calculated at the Hartree–Fock level with a heteroatom-polarized valence-double-zeta basis set (HF/MIDI!) and a mean unsigned error of 0.57 kcal/mol when it is applied using gas-phase geometries from Austin Model 1 (AM1). For a smaller set of 99 solutes, we compared the new model to two previously published models based on atomic solvation parameters, and we achieve a mean unsigned error of 0.56 kcal/mol as compared to 1.87 and 2.13 kcal/mol for the previous models. A simple extension is provided to allow treatment of certain kinds of charged groups. The model is expected to be especially useful for problems requiring high efficiency because of the size of the system, e.g., protein folding, or problems requiring rapid estimations because of the large number of calculations required, e.g., scoring of combinatorial libraries.

1. Introduction

There are many possible approaches to calculating solvation energies.^{1,2} An approach that we have espoused combines three elements: (1) explicit inclusion of electrostatics with a distributed monopole charge representation and a molecular shape approximated by an all-atom superposition of van der Waals spheres, (2) self-consistent polarization of the solute electronic charge distribution by the solvent reaction field, and (3) explicit inclusion of first-solvation-shell effects by parametrized atomic surface tensions using solvent-accessible surface areas. Elements 1 and 2 together constitute what we call the electrostatics part of the calculation, and in principle (but see next paragraph for subtleties), element 3 is nonelectrostatic.

We have developed a sequence of models that each contain in various ways all three of these elements, and we have labeled them SM1,³ SM1a,³ SM2,^{4,5} SM3,^{5,6} SM2.1,⁷ SM3.1,⁷ SM4,^{8,9} SM5.4,¹⁰ SM5.2PD,¹¹ and SM5.4PD.¹¹ Models with numbers in the range SM1 through SM4 represent our historical progress, and in addition they might still be useful for specialized purposes such as comparison to previous work or testing the sensitivity of predictions to the details of the parametrization. For the most part, though, these models are rendered obsolete by our current SMx approach, called SM5. SM5 is actually a suite of models corresponding to various levels, where—just as in *ab initio* electronic structure theory—higher levels are expected to be more reliable (in general) but also require a greater computational effort. In *ab initio* electronic structure theory, the “level” is specified primarily by the basis set, the choice of formalism for including electron correlation, and the method used to obtain geometries (e.g., optimized at the same level as energies are calculated or taken from a lower level). In the SM5 suite, the level is primarily determined by the class of partial atomic charges used, the formalism employed for the electrostatics, the

Hamiltonian used for the solute, and the method used to obtain geometries. All SM5 models involve a set of geometry-dependent atomic surface tension coefficients optimized for best performance (in predicting free energies of solvation) when used in conjunction with the other choices made for that model.

In the papers detailing the SMx models, we emphasized that the atomic surface tensions not only provide a parametrization of the nonelectrostatic aspects of first-solvation-shell effects, but they also make up for inadequacies of the theoretical solute description (such as any systematic errors in atomic partial charges or even for the very idea of replacing the continuous charge distribution of the solute by atom-centered partial charges) and for inadequacies of the theoretical model used for electrostatics (such as the intrinsic uncertainty of the boundaries between regions of unit (vacuum) dielectric constant and regions with bulk solvent dielectric constant). Taking this argument to the limit, we could ask if the surface tensions can make up for totally neglecting elements 1 and 2, i.e., not including explicit electrostatics at all. This is not only theoretically intriguing, it is of immense practical importance, since a surface-tensions-only model could lower the required computer time enormously, especially for large systems. In fact, driven by practical considerations, surface-tensions-only models have achieved a considerable prominence in the protein-folding literature,^{12–18} and we expect that they could be extremely useful for docking¹⁹ calculations as well.

In the present paper we put forth our own version of a surface-tensions-only model for aqueous solvation energies. This work differs from previous surface-tensions-only models^{12–18,20,21} in three main respects: (1) We attempt to obtain parameters that can be used for any organic solute containing H, C, N, O, F, S, Cl, or Br, rather than just for polypeptides. (2) Our model uses the exposed van der Waals surface defined by the radii proposed by Bondi,²² rather than the solvent-accessible surface area typically used in other^{14,20} models. (3) We use our experiences

[®] Abstract published in *Advance ACS Abstracts*, July 15, 1997.

with the SMx models in the way that we formulate the surface tension terms. In particular we write the surface tensions in the SM5 style that combines surface tension coefficients, which are semiempirical linear parameters, with geometry-based functional forms to recognize solute functionality.

Most of the surface tension functional forms (for the geometry dependence of the atomic surface tensions) used in the present paper were developed previously¹⁰ for the SM5.4 models. In those models, the surface tensions were combined with class IV charges, yielding the SM5.4/A and SM5.4/P models,¹⁰ and with class IV charges and a pairwise descreening algorithm, yielding models we called SM5.4PD/A and SM5.4PD/P.¹¹ In both of these cases, the surface tension functionals were used in conjunction with the solvent-accessible surface area, rather than the exposed van der Waals surface used here. Although the new model employs the exposed van der Waals surface area rather than the solvent-accessible surface area, and although it involves a fine tuning of the SM5 functional forms and additional functional forms which improve the results for halogens and certain functional groups containing nitrogen, the basic approach in making the surface tensions be functions of bond distances builds on the work of the SM5.4 models and is similar in spirit for the way group functionality is handled. Since the model presented here has no explicit charges (i.e., all partial charges are 0) and is constrained in the sense that the geometries are not allowed to relax in the presence of the solvent reaction field, it is called SM5.0R, where the 0 denotes that no explicit partial atomic charges are employed for neutral functional groups and the R denotes rigid solutes.

Where besides protein folding and docking might this cost-cutting exercise be worthwhile? When a system is so large that molecular mechanics is to be preferred over molecular orbital theory, the same size considerations will favor SM5.0R over calculations with explicit electrostatics and self-consistent geometry optimization. A second, perhaps less obvious, area of application occurs when the number of calculations is large, rather than or in addition to the size of the individual systems being large. An example of the latter would be estimating the desolvation contribution to drug–receptor binding for a whole combinatorial library of potentially biologically active agents.

2. Theory

We use a standard-state concentration of 1 mol L⁻¹ for both the gas phase and solution, assuming an ideal gas and an ideal solution as is usually done. In previous SMx models,^{3–11,23} the standard-state free energy of solvation has been expressed as

$$\Delta G_S^\circ = \Delta G_{\text{ENP}} + G_{\text{CDS}} \quad (1)$$

where ΔG_{ENP} contains the electronic and nuclear repulsion energy of the solute and the electric polarization free energy of the solvent, and G_{CDS} contains the free energy effects of the solvent cavitation, modified dispersion interactions, and structural rearrangement of the solvent that takes place in the first solvation shell. In this section we limit the discussion to neutral molecules without charged groups (i.e., zwitterions are excluded as well as cations and anions); we will return to charged groups in section 3.2. For neutral molecules without charged groups, we attempt to model aqueous solvation by ignoring all the electrostatic interactions (ΔG_{ENP} is set equal to 0) and assuming that the interaction between every atom k of a solute and the solvent is proportional to the exposed van der Waals surface area, A_k , of atom k .

Note that the exposed van der Waals surface area is different from the solvent-accessible surface area. The solvent-accessible surface area^{13,24,25} is dependent upon the geometry of a given molecule, the set of atomic radii, $\{r_k\}$, used for the solute atoms, and the effective radius, $r_{\text{H}_2\text{O}}$, used for solvent water molecules. The solvent-accessible surface area of atom k is then defined as the surface area of a sphere centered at atom k with a radius of $r_k + r_{\text{H}_2\text{O}}$ that is not contained within any of the spheres of radius $r_{k'} + r_{\text{H}_2\text{O}}$ centered at each of the other atoms k' in the molecule. In contrast, the exposed van der Waals surface area of atom k , which is used in the current work, is the surface area of a sphere centered at atom k with a radius of r_k which is not contained within any of the spheres of radii $r_{k'}$ centered at each of the other atoms in the molecule. Either of these surface areas can be (and is) calculated analytically by an algorithm described previously,⁷ and the exposed van der Waals surface area is a special case of the solvent-accessible surface area obtained by setting the solvent radius, $r_{\text{H}_2\text{O}}$, equal to 0.

In the SM5.0R model, we express the total free energy of a solute as

$$\Delta G_S^\circ = \sum_k \sigma_k A_k(\{r_k\}) \quad (2)$$

where σ_k is a proportionality constant which we call an atomic surface tension. With the exception of using the exposed van der Waals surface area rather than the solvent-accessible surface area, this quantity corresponds to our definition of G_{CDS} in previous work, with the clear difference, however, that in the present model the surface tensions must include all electrostatic effects implicitly. The atomic surface tensions depend on parameters called surface tension coefficients and on geometry.

In the SM5.0R functional forms, which are a superset of the SM5.4¹⁰ functional forms, the functional dependency of σ_k differs for each atomic number. As introduced previously,¹⁰ the functionals involve a COT (cutoff tanh) switching function which can be represented as

$$T(R_{kk'}|\bar{R},\Delta R) = \begin{cases} \exp\left[-\left(\frac{\Delta R}{\Delta R - R_{kk'} + \bar{R}}\right)\right] & R_{kk'} \leq \bar{R} + \Delta R \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where $R_{kk'}$ is the distance between two atoms k and k' and \bar{R} and ΔR are parameters. (Note that \bar{R} is the center of the switch and ΔR is the half-range over which the COT becomes engaged.) The COT function in eq 3 was, to the best of our knowledge, introduced into the molecular modeling literature by Stillinger and Weber.²⁶ It has the following properties that make it especially attractive for the present kind of application: (i) It vanishes identically for all atoms farther away than the prescribed cutoff distance of $\bar{R} + \Delta R$. (ii) It has an infinite number of continuous derivatives, even at the cutoff distance.

The functionals for the geometry dependence of σ_k when atom k is an H, S, or F atom and the portions of the functionals for C and O which do not involve nitrogen are unchanged from those developed for previously for the SM5.4/A model (except for the values of some of the \bar{R} and ΔR parameters). In other cases the functional forms were retuned for the present model. A full list of all surface tension functionals for the SM5.0R model for the supported atom types follow:

$$\sigma_{k|k=\text{H}} = \sigma_{\text{H}}^{(0)} + \sum_{k'=\text{C,O,N,S}} \sigma_{\text{HK}} T(R_{\text{HK}}|\bar{R}_{\text{HK}},W) \quad (4)$$

$$\sigma_{k|k=C} = \sigma_C^{(0)} + \sigma_{CC} \sum_{\substack{k'=C \\ k' \neq k}} T(R_{kk'}|\bar{R}_{CC}, W) + \sigma_{CC}^{(2)} \sum_{\substack{k'=C \\ k' \neq k}} T(R_{kk'}|\bar{R}_{CC}^{(2)}, W_{CC}) + \sigma_{CN} \left[\sum_{k'=N} T(R_{kk'}|\bar{R}_{CN}, W) \right]^2 \quad (5)$$

$$\sigma_{k|k=O} = \sigma_O^{(0)} + \sigma_{OC} \sum_{k'=C} T(R_{kk'}|\bar{R}_{OC}^{(2)}, W_{OC}) + \sigma_{OO} \sum_{\substack{k'=O \\ k' \neq k}} T(-T(R_{kk'}|\bar{R}_{OO}, W)|R_{TT}, W_{TT}) + \sigma_{ON} \sum_{k'=N} T(R_{kk'}|\bar{R}_{ON}, W) \quad (6)$$

$$\sigma_{k|k=N} = \sigma_N^{(0)} + \sigma_{NC} \left\{ \sum_{k'=C} T(R_{kk'}|\bar{R}_{CN}, W) \left[\sum_{\substack{k''=k \\ k'' \neq k'}} T(R_{kk''}|\bar{R}_{CK'}, W) \right]^2 \right\}^{1.3} + \sigma_{NC}^{(2)} \sum_{k'=C} [T(R_{kk'}|\bar{R}_{CN}, W) \sum_{k''=O} T(R_{kk''}|\bar{R}_{CK'}, W)] \quad (7)$$

$$\sigma_{k|k=F} = \sigma_F^{(0)} \quad (8)$$

$$\sigma_{k|k=S} = \sigma_S^{(0)} + \sigma_{SS} \sum_{\substack{k'=S \\ k' \neq k}} T(R_{kk'}|\bar{R}_{SS}, W) \quad (9)$$

$$\sigma_{k|k=Cl} = \sigma_{Cl}^{(0)} + \sigma_{ClC} \sum_{k'=C} \left[T(R_{kk'}|\bar{R}_{CCl}, W) \sum_{\substack{k'' \neq k \\ k'' \neq k'}} \frac{T(R_{kk''}|\bar{R}_{CK'}, W)}{\chi_{Z_{k'}}} \right] \quad (10)$$

and

$$\sigma_{k|k=Br} = \sigma_{Br}^{(0)} + \sigma_{BrC} \sum_{k'=C} \left[T(R_{kk'}|\bar{R}_{CBr}, W) \sum_{\substack{k'' \neq k \\ k'' \neq k'}} \frac{T(R_{kk''}|\bar{R}_{CK'}, W)}{\chi_{Z_{k'}}} \right] \quad (11)$$

where W is the standard COT width parameter of 0.30 Å, W_{CC} and W_{OC} are specialized COT width parameters of 0.07 and 0.10 Å respectively, and χ_Z is the Pauling electronegativity^{27,28} for an atom with atomic number Z .

Notice that eq 6 contains a cutoff function as an argument of another cutoff function; this was the simplest way to achieve the desired geometry dependence, which singles out geminal O—O pairs. The geometry dependence of this form cannot be obtained with a single COT because that would include bonded O—O pairs as well. Notice that there is a typo in this term in ref 10 (the arguments are interchanged).

In eq 7, the sum over k'' includes all atoms except k and k' ; i.e., if the molecule has N atoms, this is a sum over $N - 2$ atoms. However, the number of nonzero terms is restricted by the cutoff.

Note that all summations in eqs 4–11 are over all atoms of the specified atomic numbers, not just over bonded or 1–3 atom pairs. Thus the user need not assign bonds; the cutoff functions in the equations do all the work, and the results are continuous functions of geometry even along reaction paths. Furthermore all cutoff functions have an infinite number of continuous derivatives.

3. Parametrization

3.1. Neutral Molecules. *3.1.1. Training Set.* In order to parametrize the model, we first select a training set of solute

molecules for which experimental free energy of solvation data are available. The molecules in this set include a wide variety of functional groups. In addition, we are careful to select molecules that do not have complicating conformational issues. As a starting point, we selected the SM5.4 aqueous models neutral training set developed in earlier work.¹⁰ If we restrict consideration only to H, C, N, O, F, S, Cl, and Br, this set contains 205 molecules. The sources^{29–36} of the experimental data for solvation free energies of these compounds are explained in a previous paper.¹⁰ Our initial parametrization of the present model raised issues that we felt were best resolved by adding some additional nitrogen-containing and halogen-containing compounds.

To make our parametrization for nitrogen as robust as possible, we added 1,1-dimethyl-3-phenylurea, 9-methyladenine, and 1-methylthymine. For halogen compounds, we added tetrafluoromethane, hexafluoroethane, octafluoropropane, 1,1,1,2-tetrachloroethane, hexachloroethane, 2-chlorobutane, 1-chloropentane, 2-chloropentane, α -chlorotoluene, *o*-chlorotoluene, 2,2'-dichlorobiphenyl, 2,3-dichlorobiphenyl, 2,2',3-trichlorobiphenyl, 3-bromopropene, 1-bromoisobutane, α -bromotoluene, *p*-bromotoluene, dichlorodifluoromethane, trichlorofluoromethane, bromotrichloromethane, chloropentafluoroethane, *p*-bromophenol, 3,5-dibromo-4-hydroxybenzonitrile, 2,6-dichlorobenzonitrile, 2,6-dichlorothiobenzamide, and 4-amino-3,5,6-trichloropyridine-2-carboxylic acid. Experimental free energy of solvation data for these new solutes were taken from the MedChem database,³⁷ except for 9-methyladenine and 1-methylthymine, which were taken from Ferguson et al.³⁸

The MedChem database³⁷ entries that are of interest are reported as the logarithm of the partition coefficient between air and water. Many of these data points also have footnotes in the database that give additional information about the reported measurement. Any partition coefficient measurements which have any of the following footnotes were disqualified: the data was collected at a temperature outside the range 293–303 K, the partition coefficient was measured at a pH outside the range 6–8, the aqueous phase was not purely aqueous (the exception to this was a phase mixed with phosphate buffers, since results in such media tend to agree well with pure aqueous phases), the solute was not in its “true” form (e.g., dimerization occurred), salting out of the solute occurred, or the partition coefficient has a footnote which indicated that the measurement was of doubtful validity. In all cases, after the disqualified partition coefficient measurements were removed, the remaining experimental data points were within 2 standard deviations of the mean.

Our final training set for the SM5.0R model contains a total of 235 neutral solutes. Of these solutes, only 1,1-dimethyl-3-phenylurea has complicating conformational issues which required considering more than one conformation explicitly. This molecule has two low-lying conformations that differ from each other by only 1.8 kcal/mol in the gas phase. The conformations differ primarily in the orientation of the anilide, with the *Z* orientation of carbonyl oxygen to phenyl being lower in energy in the gas phase. To predict the absolute free energy of solvation for this molecule, we statistically average over both minima both in the gas phase and in solution using the following relation³⁹

$$\exp[-\Delta G_s^\circ/RT] = \sum_C P_C \exp[-\Delta G_s^\circ(C)/RT] \quad (12)$$

where P_C is the equilibrium mole fraction of conformation C in the gas phase.

3.1.2. Gas-Phase Geometries. Unlike our previous SM x models, which were parametrized with geometries that were self-consistently optimized in solution, the SM5.0R model was designed to be used with any accurate gas-phase geometry, and the structural changes that occur upon placing the molecule in solution are absorbed by the parameters of the model. To obtain good gas-phase geometries for most of the molecules of our training set, we optimized the structures by performing Hartree–Fock (HF) calculations with a heteroatom-polarized split valence basis, in particular the MIDI!⁴⁰ basis. The polarization basis functions in the MIDI! basis set were chosen⁴⁰ in previous work to economically provide accurate geometries (and reasonable charge distributions, although that feature is not used here).

HF/MIDI! geometries were used for all molecules except those containing Br. Since the MIDI! basis set was not defined for bromine at the time that we carried out this work, we used a different method for finding the gas-phase geometries for bromine-containing compounds. First, we optimized geometries for two trial bromine-containing molecules, methyl bromide and dibromomethane, by Møller–Plesset second-order perturbation theory calculations^{41–43} with the cc-pVDZ basis set.⁴⁴ We then compared these data to the predicted gas-phase geometries from AM1^{45–47} and PM3^{48,49} for the same molecules. PM3 predicted a more accurate C–Br bond length than AM1, so we chose to use PM3 to obtain the gas-phase structures for the bromine-containing molecules.

We believe that the SM5.0R model is stable with respect to small changes in the gas-phase geometry used for the calculation. To demonstrate this point, results for all solutes were also obtained using geometries calculated by the AM1 method.

3.1.3. Parameters To Be Optimized. In the SM5.0R model, the parameters can be separated into three categories: (a) nonlinear parameters which specify how the atomic surface tensions depend on the local molecular geometry; (b) the van der Waals radius for each atomic number, which is a nonlinear parameter used in calculating the exposed van der Waals surface area for a particular atom within the molecule; (c) the linear surface tension coefficients.

Parameters of type (a) were primarily taken from earlier work.¹⁰ The only changes in these \bar{R} and ΔR parameters were made to remove secondary contributions to the G_{CDS} energy. For example, in the SM5.4 functional forms, the switch point \bar{R}_{HC} for CDS contributions associated with hydrogen attached to carbon atoms is 1.85 Å, and the half-width of the switch is 0.30 Å. This implies that there is a specific contribution to the calculated CDS energy whenever the distance between a specific hydrogen and carbon is less than or equal to 2.15 Å. The contribution gradually increases as the distance between the hydrogen and carbon decreases, and for distances shorter than 1.55 Å the contribution is relatively independent of further decreases in bond length. In optimized gas-phase geometries, the distance between an alcoholic hydrogen atom and the nearest carbon atom can be as small as 1.9 Å. Thus, in the SM5.4 models there is a small H–C-type contribution to the free energy associated with hydrogen atoms in alcohol groups. Although we are not aware of any problems this has caused in the SM5.4 models, and although we were aware of them when we parametrized the SM5.4 models, in the present model we lowered the center of the switch for H–C interactions to 1.55 Å in SM5.0R to make a more distinct separation of the H–C and H–O effects. Since typical H–C bond lengths are 1.1 Å, the resulting surface tension has full contributions from H–C interactions when the hydrogen is actually bonded to the carbon and has no contributions from secondary interactions. A similar

TABLE 1: \bar{R} COT Parameters (Å) Used in SM5.0R

form ^a	k'	\bar{R}
R_{HK}	C, O, N	1.55 ^b
	S	2.14 ^c
R_{CK}	C,N	1.84 ^c
	O	1.84 ^b
	F	1.84 ^b
	S	2.20 ^b
	Cl	2.10 ^b
	Br	2.30 ^b
R_{OK}	N	1.50 ^b
	O	2.75 ^c
R_{SK}	S	2.75 ^c
$R_{\text{CK}}^{(2)}$	C	1.27 ^c
$R_{\text{OK}}^{(2)}$	C	1.33 ^c
R_{TT}		−0.4 Å ^c
W_{TT}		0.4 Å ^c

^a Parameters $R_{kk'}$ and $R_{k'k}$ have identical values. ^b Changed or new in this work. ^c Reference 10.

adjustment was made for H–O and H–N interactions. When selecting \bar{R} and ΔR parameters for the new surface tension functionals, we examined typical bond lengths for the interactions we were trying to model and selected the COT parameters which singled out the desired interactions but minimized secondary effects. The full set of \bar{R} and ΔR parameters used in the SM5.0R model are listed in Table 1.

Although we do not review here all the earlier work that led to eqs 4–11, it may be useful to point out that the exponents of 2 and 1.3 in eqs 3 and 7 are the result of a trial-and-error approach in previous work, and we found no evidence that there would be any significant benefit in changing them for the present work.

The type (b) parameters were chosen to be the atomic radii suggested by Bondi.²²

Linear parameters of type (c) were fit using the experimental free energies of solvation values for the molecules in the training set. Thus the 23 surface tension coefficients involved in eqs 3–11 were fit to minimize the mean-squared deviation of the predicted solvation free energies for the 235 molecules in our training set from experiment. The Pauling electronegativities were taken from standard sources.^{27,28}

3.1.4. Progression of the Parametrization. Initially, we started with the 107 compounds in our training set that contain at most H, C, and O. First, we fixed the nonlinear COT parameters listed in Table 1. Then we solved for the linear surface tension coefficients, σ_z , by minimizing the sum of the squares of the error between the predicted solvation free energy and the experimental solvation free energy for the 107 molecules in this portion of the test suite.

We then froze the surface tensions determined in this step and considered the 46 compounds in our training suite that contain N plus at most H, C, and O. We first developed two surface tension functionals that help distinguish certain functionalities that had systematic errors when only the original SM5.4 nitrogen-containing functionals were used. The σ_{CN} coefficient and associated functional in eq 5 help differentiate between amines with different numbers of carbon substituents. The nitrogen surface tension functional associated with the $\sigma_{\text{NC}}^{(2)}$ surface tension coefficient in eq 7 was added to identify amides and ureas which would otherwise be systematically undersolvated. Thus, we fit six surface tension coefficients to the nitrogen subset of our parametrization suite.

As discussed in section 3.1, fitting nitrogen surface tension coefficients was complicated by the need to consider more than one conformation for 1,1-dimethyl-3-phenylurea. Within the parametrization step, the conformational averaging was handled

TABLE 2: Atomic Radii (Å), Pauling Electronegativities, and Surface Tension Coefficients (cal mol⁻¹ Å⁻²) Used in SM5.0R

<i>k</i>	<i>r_k</i>	<i>χ_k</i>	<i>σ_k</i>	<i>σ_{Hk}</i>	<i>k</i>	<i>k'</i>	<i>σ_{kk'}</i>	<i>σ_{kk'}</i> ⁽²⁾
H	1.20	2.20	98.22		O	O	152.57	
C	1.70	2.55	89.21	-126.78	O	N	363.49	
N	1.55	3.04	-277.00	-212.04	O	C	171.64	
O	1.52	3.44	-376.46	-248.86	C	C	-84.71	-50.91
S	1.80	2.58	-163.04	171.42	C	N	78.99	
F	1.47	3.98	29.85		S	S	98.04	
Cl	1.75	3.16	73.52		N	C	-68.91	-529.66
Br	1.85	2.96	70.20		Cl	C	-223.18	
					Br	C	-220.05	

as follows. First we assigned arbitrary weights to each conformer, and a single data point utilizing these weights was added to the other 45 nitrogen data points. We found appropriate values for the nitrogen-based surface tension coefficients by minimizing the sum of the squares of the error for the nitrogen data points. These surface tensions were used to calculate the free energies of solvation for our two conformers, and this information was plugged into the left-hand side of eq 12. From there we determined the absolute free energy of the compound and were able to determine appropriate weights for each of the two conformers in solution. The weights were then used in the 1,1-dimethyl-3-phenylurea data point, and the nitrogen-based surface tension coefficients were found again. This loop was continued until the relative weights of the two conformers in solution stopped changing. The surface tension coefficients obtained at his point were the final nitrogen surface tension coefficients. The final weighting called for 97% of the 1,1-dimethyl-3-phenylurea molecules to be in the Z conformation in aqueous solution, as compared to 95% in the gas phase.

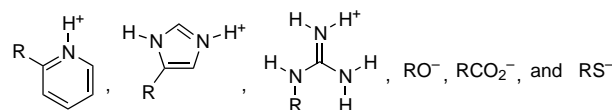
Once we set the surface tension coefficients for H, C, O, and N, we determined appropriate surface tension coefficients for the 11 sulfur-containing compounds in our test suite. As was found when parameterizing the SM5.4 models,¹⁰ the disulfides were significantly oversolvated unless the σ_{SS} surface tension coefficient and related functional from eq 9 were used. In protein modeling, it may be useful to distinguish cystine bridges from methionine and cysteine residues.

The coefficients for the halogens were determined simultaneously because a large number of the halogenated molecules in the training set contain more than one kind of halogen. Initially, we tried to parametrize the halogen-containing compounds using a single surface tension coefficient per halogen as was done for the previous SM5.4 models,¹⁰ but such an attempt resulted in a model with much larger errors than our models with explicit electrostatics. We discovered that combining a COT that differentiates between different degrees of saturation for the carbon attached to a chlorine or bromine atom with a functional form that recognizes the differences caused by different electronegativities of the attached atoms greatly improved the problem. The final SM5.0R model uses two such surface tension functionals with the corresponding coefficients, σ_{ClC} and σ_{BrC} , and these forms allowed us to halve the error in the chlorinated and brominated hydrocarbons as well as to reduce the error for the bifunctional halogenated compounds in our training set when compared to a model without the new functional forms.

The final surface tension coefficients for the SM5.0R model are listed in Table 2.

3.2. Charged Groups. The simple model presented in this paper is remarkably successful at predicting solvation energies for neutral molecules, even those with highly polar groups such as aldehydes, ketones, and carboxylic acids. For molecules with charged groups, such as cations, anions, and zwitterions, an

explicit electrostatic term should be added. In the Appendix we present such an extension for the following types of charged groups, all of which are common in biological systems: $R_nNH_{4-n}^+$ where *n* is between 0 and 3,



These functionalities include all charged groups present in naturally occurring amino acids, along with the alkoxide group, and hence allow extension of the present model to many polypeptides and proteins. In the spirit of the SM5.0R model, the solvation free energies are predicted as much as possible from local properties of the solute; in particular the charge is distributed only over the protons in the positively charged groups and only over the oxygens in the negatively charged groups. However, since electrostatic effects are intrinsically long ranged, the proposed model for including them does take account of dielectric descreening even from distant groups. Further details of the parametrization of the model for charged groups are provided in the Appendix.

When SM5.0R is extended with this electrostatic treatment, we will call it SM5.05R (because it is SM5.0R, with explicit charges of zero for neutral groups, and SM5.1R-like, with type I charges for charged groups). Of course if the charged groups play a central role in the process under consideration, one should seriously consider whether it is affordable to use a higher-level treatment of the electrostatics, such as SM5.2PD¹¹ or SM5.4.¹⁰

4. Results and Discussion

4.1. Performance of the Model. The SM5.0R model is very successful in predicting the free energies of solvation. The mean signed and unsigned errors for various types of solutes are summarized in Tables 3 and 4 (the WE and OONS columns in these tables are explained in the next subsection). The full results for all 235 solutes from which the summaries in Tables 3 and 4 were extracted are presented in the Supporting Information.

For the subset of our training suite which contains only H, C, and/or O atoms, the SM5.0R model at both the HF and AM1 geometries achieves a mean unsigned error of 0.40 kcal/mol with signed errors of -0.01 and -0.03 kcal/mol, respectively. The greatest successes of the model are for branched alkanes, alkenes, alkynes, and alcohols where the mean unsigned errors are 0.30 kcal/mol or less. The ethers and arenes proved the most challenging with mean unsigned errors of more than 0.50 kcal/mol.

As a whole, the solvation free energies predicted by SM5.0R for nitrogen-containing compounds were not as successful, with mean unsigned errors of 0.88 kcal/mol for SM5.0R/HF and 0.93 kcal/mol for SM5.0R/AM1. The model performed better for nitrohydrocarbons, aliphatic amines, and nitriles with mean unsigned errors of less than 0.70 kcal/mol. The amides and ureas and also the aromatic amines had mean unsigned errors of over 1.00 kcal/mol.

Sulfur-containing molecules were treated very successfully by the SM5.0R model with mean unsigned errors over the 11 molecules of less than 0.30 kcal/mol using either HF or AM1 geometries. The signed error for these molecules has a magnitude of less than 0.01 kcal/mol for the HF geometries and increases to 0.11 kcal/mol for the AM1 geometries.

The chloroalkanes, chloroarenes, and brominated hydrocarbons all have mean unsigned errors of less than 0.50 kcal/mol in the SM5.0R model, while the fluorinated hydrocarbons and

TABLE 3: Mean Unsigned Errors by Functional Group for All Neutral Solutes Used in the Parametrization of SM5.0R

functional group	no. of molecules	SM5.0R//HF	SM5.0R//AM1	WE//AM1	OONS//AM1	SM5.4/A	exptl dispersion ^a
Compounds Containing at Most C, H, and/or O							
unbranched alkanes	8	0.49	0.40	1.32	0.22	0.60	0.35
branched alkanes	5	0.27	0.22	1.55	0.28	0.65	0.22
cycloalkanes	5	0.46	0.42	0.33	0.81	0.21	0.34
alkenes	9	0.22	0.19	0.45	<i>b</i>	0.45	0.37
alkynes	5	0.09	0.10	0.91	<i>b</i>	0.25	0.20
arenes	8	0.54	0.63	2.61	0.85	0.15	1.16
alcohols	16	0.25	0.29	1.49	0.90	0.51	1.24
ethers	9	0.72	0.76	2.47	<i>b</i>	0.91	1.36
aldehydes	6	0.32	0.37	0.43	<i>b</i>	0.29	0.53
ketones	12	0.43	0.33	1.00	5.36	0.35	0.63
carboxylic acids	5	0.53	0.26	1.27	0.21	0.78	0.20
esters	12	0.28	0.47	0.33	<i>b</i>	0.43	0.32
bifunctional	5	0.53	0.64	1.88	<i>b</i>	0.39	3.73
water, H ₂	2	0.94	0.96	4.82 ^c	<i>b</i>	1.59	4.32
subtotal	107	0.40	0.40	1.27 ^c		0.49	3.05
Compounds Containing N							
aliphatic amines	15	0.45	0.70	2.77	3.38	0.80	1.39
aromatic amines	10	1.01	1.01	2.58	1.86	0.44	0.37
nitriles	4	0.65	0.65	2.36	<i>b</i>	0.49	0.16
nitrohydrocarbons	6	0.39	0.35	5.54	<i>b</i>	0.45	0.36
amides and ureas	4	2.80	2.26	4.60	7.96	2.64	0.86
bifunctional	5	1.08	1.21	2.20	<i>b</i>	0.79	6.51
ammonia, hydrazine	2	1.05	1.17	8.37	8.48 ^d	2.62	2.51
subtotal	46	0.88	0.93	3.40		0.89	2.69
Compounds Containing S, H, and/or C							
thiols	4	0.33	0.33	0.86	0.51	0.31	0.59
organic sulfides, H ₂ S	5	0.26	0.20	1.96	2.13 ^e	0.40	0.66
organic disulfides	2	0.07	0.29	1.28	1.63	0.30	0.10
subtotal	11	0.25	0.26	1.43	1.38 ^e	0.35	0.58
Compounds Containing Halogens							
fluorinated hydrocarbons	6	1.13	1.05	<i>b</i>	<i>b</i>	0.59	2.12
chloroalkanes	13	0.42	0.50	<i>b</i>	<i>b</i>	0.28	0.63
chloroalkenes	5	0.97	1.00	<i>b</i>	<i>b</i>	0.65	0.26
chloroarenes	8	0.27	0.29	<i>b</i>	<i>b</i>	0.23	0.61
brominated hydrocarbons	14	0.37	0.39	<i>b</i>	<i>b</i>	0.25	0.79
other halo molecules	25	0.81	0.80	<i>b</i>	<i>b</i>	0.63	4.00
subtotal	71	0.63	0.64	<i>b</i>	<i>b</i>	0.44	2.10
Overall							
OONS subset ^f	99	0.54	0.56	1.87	2.13	0.61	
WE subset ^g	163	0.51	0.54	1.85		0.59	
entire set	235	0.55	0.57			0.55	3.03

^a Experimental dispersion is defined as the root-mean-squared deviation of the members of a particular set from their mean value. ^b This model does not contain types for this functional group. ^c Since this model is a unified atom model, it has no types for hydrogen, thus H₂ was left out when calculating this number. ^d This model does not contain types sufficient to calculate the solvation free energy for ammonia, thus ammonia was left out when calculating this number. ^e This model does not contain types sufficient to calculate the solvation free energy for H₂S, thus H₂S was left out when calculating this number. ^f Subset of molecules to which the OONS model is applicable; all models are applicable to this subset. ^g Subset of molecules to which the WE model is applicable; all models except the OONS model are applicable to this subset.

chloroalkenes both have mean unsigned errors of over 1.0 kcal/mol. The most systematic error is for the chloroalkenes which are undersolvated by about 1.00 kcal/mol on average.

Overall, the SM5.0R model achieves a mean unsigned error of only 0.55 kcal/mol using the HF geometries and 0.57 kcal/mol using AM1 geometries. The mean signed error is 0.07 and 0.08 kcal/mol for the SM5.0R model with HF and AM1 geometries, respectively. The close general agreement between the results obtained using the HF and AM1 geometries indicates that the SM5.0R model is relatively stable to small perturbations of input geometry and thus could successfully be used with accurate or reasonably accurate gas-phase geometries obtained from many different sources.

4.2. Comparison to Previous Work. As mentioned in the introduction, other groups have developed nonelectrostatic solvation models using solvent-accessible surface areas and surface tension coefficients. (In these models, the surface tension coefficients are often referred to as atomic solvation parameters.) The main original intent of these models is not to

find solvation free energies for molecules such as those in our training set, but rather to empirically calculate solvation effects for much larger molecules such as proteins.

The first model is due to Ooi et al.²⁰ They used, in part, vapor-to-water transfer free energies of small solute molecules given by Cabani et al.³² to find their atomic solvation parameters. These data are a subset of the data used to parametrize the SM5.0R model. The model created by Ooi et al.²⁰ applies a united-atom approach and contains seven types of atoms or groups that allow the method to be applied to a portion of our neutral training set. Henceforth, this model will be referred to as the OONS (Ooi–Oobatake–Nemethy–Scheraga) model.

A second united-atom atomic solvation parameter model was developed by Wesson and Eisenberg¹⁴ using affinities of amino acid side chains for water. The model was parametrized using the Kyte and Doolittle⁵⁰ adjustment to the data given by Wolfenden et al.,⁵¹ and contains a single surface tension coefficient for C, O, N, and S. (Note that the radius used for oxygen is identical to the radius used for nitrogen.) Since the

TABLE 4: Mean Signed Errors by Functional Group for All Neutral Solutes Used in the Parametrization of SM5.0R

functional group	no. of molecules	SM5.0R//HF	SM5.0R//AM1	WE//AM1	OONS//AM1	SM5.4/A
Compounds Containing at Most C, H, and/or O						
unbranched alkanes	8	-0.49	-0.40	-1.32	-0.19	-0.60
branched alkanes	5	-0.19	-0.08	-1.55	-0.28	-0.65
cycloalkanes	5	-0.41	-0.31	-0.33	0.81	-0.14
alkenes	9	-0.08	0.04	-0.32	<i>a</i>	0.33
alkynes	5	0.09	0.10	0.91	<i>a</i>	-0.04
arenes	8	0.53	0.63	2.61	0.62	0.10
alcohols	16	0.19	0.22	1.49	-0.30	0.32
ethers	9	0.12	0.15	2.47	<i>a</i>	0.12
aldehydes	6	-0.22	-0.36	0.23	<i>a</i>	-0.29
ketones	12	0.33	0.21	1.00	5.36	0.27
carboxylic acids	5	0.53	0.26	-1.27	0.08	0.78
esters	12	-0.27	-0.47	0.33	<i>a</i>	-0.43
bifunctional	5	-0.52	-0.57	-0.33	<i>a</i>	0.10
water, H ₂	2	-0.94	-0.96	-4.82 ^b	<i>a</i>	-1.56
subtotal	107	-0.01	-0.03	0.48 ^b		-0.01
Compounds Containing N						
aliphatic amines	15	0.05	0.40	1.42	3.08	0.25
aromatic amines	10	0.32	0.32	2.53	1.26	0.06
nitriles	4	-0.01	-0.01	-2.36	<i>a</i>	0.29
nitrohydrocarbons	6	-0.01	-0.04	-5.54	<i>a</i>	0.23
amides and ureas	4	0.74	0.78	4.10	7.96	2.34
bifunctional	5	-0.18	-0.17	1.94	<i>a</i>	-0.25
ammonia, hydrazine	2	0.18	0.30	-8.37	-8.48 ^c	0.61
subtotal	46	0.14	0.25	0.29		0.35
Compounds Containing S, H, and/or C						
thiols	4	0.04	0.13	0.86	0.20	0.03
organic sulfides, H ₂ S	5	-0.04	0.02	1.36	2.13 ^d	0.09
organic disulfides	2	0.00	0.29	1.28	1.63	0.13
subtotal	11	0.00	0.11	1.16	0.81 ^d	0.08
Compounds Containing Halogens						
fluorinated hydrocarbons	6	0.41	0.58	<i>a</i>	<i>a</i>	0.59
chloroalkanes	13	0.15	0.02	<i>a</i>	<i>a</i>	-0.15
chloroalkenes	5	0.97	1.00	<i>a</i>	<i>a</i>	0.63
chloroarenes	8	0.20	0.22	<i>a</i>	<i>a</i>	-0.04
brominated hydrocarbons	14	0.23	0.19	<i>a</i>	<i>a</i>	-0.07
other halo molecules	25	-0.04	0.03	<i>a</i>	<i>a</i>	0.04
subtotal	71	0.18	0.20	<i>a</i>	<i>a</i>	0.06
Overall						
OONS subset ^e	99	0.20	0.27	1.10	1.73	0.22
WE subset ^f	163	0.03	0.06	0.45		0.09
entire set	235	0.08	0.10			0.10

^a This model does not contain types for this functional group. ^b Since this model is a unified atom model, it has no types for hydrogen, thus H₂ was left out when calculating this number. ^c This model does not contain types sufficient to calculate the solvation free energy for ammonia, thus ammonia was left out when calculating this number. ^d This model does not contain types sufficient to calculate the solvation free energy for H₂S, thus H₂S was left out when calculating this number. ^e Subset of molecules to which the OONS model is applicable; all models are applicable to this subset. ^f Subset of molecules to which the WE model is applicable; all models except the OONS model are applicable to this subset.

surface tension coefficients are not typed (i.e., they do not take account of bonding pattern), it is a good example of a basic surface tension model and is applicable to a larger portion of our neutral training set than is the OONS model. Henceforth, this model will be referred to as the WE (Wessen–Eisenberg) model.

The results for the SM5.0R neutral training set for SM5.0R at the HF/MIDI!-PM3 geometries will be denoted SM5.0R//HF and those calculated with AM1 geometries will be denoted SM5.0R//AM1. The WE and OONS models were also tested with AM1 geometries, and these results are denoted WE//AM1 and OONS//AM1. (As usual, “//” denotes “calculated at a geometry obtained by”.) All results are presented in Tables 3 and 4. For comparison, we have also included results from our SM5.4/A¹⁰ model which uses the CM1A charge model⁵² and involves our most complete and physical electrostatic calculation to date. In the SM5.4/A model, the geometries are minimized including the solvation effects and thus are optimized self-consistently with the model. To indicate this fact, the model does not have an R in its name and does not require a “//”.

When comparing the results obtained with the SM5.0R model to those achieved with other surface-tension-only models, the success of this approach becomes evident. As mentioned earlier, the SM5.0R//HF and SM5.0R//AM1 models achieve identical mean unsigned errors of 0.40 kcal/mol for the 107 neutral compounds containing at most C, H, and O. The WE model achieves a mean unsigned error of 1.27 kcal/mol for the C, H, and O compounds over this same set and has the most difficulty with arenes and ethers. The OONS usually does well for the bonding types for which it was parametrized. The exception to this is the ketones where the model has a mean unsigned error of 5.36 kcal/mol.

The SM5.4/A model achieves a mean unsigned error for the C, H, and O portion of our training set of 0.49 kcal/mol. Although the SM5.0R model achieved a better mean unsigned error for these molecules, this success does not necessarily indicate problems within the explicit-electrostatics treatment in SM5.4/A. The SM5.4/A model uses the solvent-accessible surface area for calculating nonelectrostatic contributions, while the present paper uses van der Waals surfaces as described in

section 2. In preliminary versions of the SM5.0R model we originally considered using the solvent-accessible surface area in conjunction with our surface tensions functionals, but we discovered that empirically the van der Waals surface worked better. Using the solvent-accessible surface area within the SM5.0R model causes the mean unsigned error for the CHO portion of our test suite to increase by nearly 0.2 kcal/mol. Thus the SM5.0R model outperforms the SM5.4/A in mean unsigned error simply because we use the exposed van der Waals surface rather than the presumably more physical solvent-accessible surface. The cycloalkanes and arenes are the only classes of neutral CHO compounds where SM5.4/A significantly outperforms the SM5.0R models.

Nitrogen-containing compounds were more challenging for our model with implicit electrostatics. The mean unsigned error for the 46 nitrogen-containing compounds in the training set is 0.88 kcal/mol for SM5.0R/HF and 0.93 kcal/mol for SM5.0R/AM1. The SM5.4/A model performed similarly in most categories of nitrogen compounds, however it had less than half the error of the SM5.0R model for aromatic amines. Aromatic ring systems in general have proved difficult for the SM5.0R model. Perhaps the reason for this problem lies in the ability of aromatic systems to redistribute charge. Since electrostatic interactions are included only implicitly in the SM5.0R model, our results tend to represent a "typical" charge distribution across multiple classes of molecules. If there are functionalities within a molecule which cause the actual charge distribution to differ from the "typical" charge distribution implicit within our parametrization, the SM5.0R model may tend to have more difficulty predicting the solvation free energy. By choosing a balanced training set which represents as many chemical functionalities as possible, we hope to have minimized this problem. Both the WE and OONS models have serious quantitative difficulties with every class of nitrogen compound.

Compounds containing sulfur are well modeled by both the SM5.0R/HF and SM5.4/A models with mean unsigned errors of 0.25 and 0.35 kcal/mol for the 11 sulfur-containing molecules in the parametrization suite. The WE and OONS models had mean unsigned errors over the sulfur-containing compounds of 1.43 and 1.38 kcal/mol, respectively.

Although the electronegativity dependence added to the chlorine and bromine parametrizations in the SM5.0R model proved somewhat successful as discussed in section 4, the SM5.0R model still has a larger error than the SM5.4/A model for the halogen-containing compounds. This demonstrates the importance of explicit electrostatics for these molecules. The fluorinated hydrocarbons and chloroalkenes have the largest errors in SM5.0R, with mean unsigned errors of 1.13 and 0.97 kcal/mol, respectively. The other classes of halogen-containing compounds all have mean unsigned errors of less than 0.50 kcal/mol with mean signed errors of less than 0.30 kcal/mol. The WE and OONS models were not parametrized for halogen-containing compounds.

It is tempting to speculate about a further consequence of the empirical finding, reported here, that the exposed van der Waals surface area correlates better with trends in solvation free energies than does the presumably more physical solvent-accessible surface area. This may account for the fact, sometimes observed,^{53,54} that the solvent-excluding surface sometimes performs better than the solvent-accessible surface because the solvent-excluding surface areas are very similar to exposed van der Waals surface areas, whereas solvent-accessible surface areas, computed with realistic effective solvent radii, are considerably larger and have a different dependence on details of molecular structure. Future research into the underly-

ing physics that makes the smaller surface areas correlate better would be very interesting.

Overall, the SM5.0R/HF method has a mean unsigned error for all 235 molecules in our test suite of only 0.55 kcal/mol, which is identical to the SM5.4/A model which explicitly includes electrostatics. When using the SM5.0R model with the AM1 geometries the error increases to only 0.57 kcal/mol, indicating that the method is fairly stable to small geometry changes and hence that we can expect reasonable results with any good gas-phase geometry calculation. When considering only the 99 molecules in our training suite which are calculable by the OONS model, the OONS model achieves a mean unsigned error of 2.13 kcal/mol when using AM1 geometries compared to 0.54 or 0.56 kcal/mol for either SM5.0R/HF or SM5.0R/AM1. For the 164 compounds of our test set which can be treated by the WE model, it achieves a mean unsigned error of 1.85 kcal/mol with AM1 geometries, compared to 0.51 and 0.54 kcal/mol for SM5.0R/HF and SM5.0R/AM1, respectively.

Starting with optimized AM1 geometries, predicting the solvation free energy with the SM5.0R/AM1 model as implemented in the AMSOL code⁵⁵ for the entire training set of 235 molecules requires a total of 6 s computing time on a Silicon Graphics Indigo-2 with an R10000 processing chip. The majority of this time is required for overhead (including reading the geometry and writing the output), which could be reduced appreciably if we cared to do so. (The WE and OONS models should theoretically take less time, but the observed difference was less than 0.005 s/trial and was therefore beneath our detection threshold.) The SM5.4/A model with explicit electrostatics requires 5325 s to complete the same 235 molecules, starting from the AM1 geometry and using version 5.4.1 of the AMSOL⁵⁵ code. (It should be noted that the SM5.4/A model uses the AM1 geometry to obtain a reference gas-phase energy for the molecule, and then optimizes the geometry including solvation effects and polarizes the CM1A charges self-consistently within the reaction field created by the solvent.) Overall, the inexpensive calculations necessary to predict solvation free energies using the SM5.0R model make it well suited for calculations on large molecules or libraries of molecules.

Acknowledgment. This work was supported in part by the National Science Foundation through Grant CHE94-23927, by the Army Research Office through Grant DAAH-04-93-G-0036, and by the National Institute of Standards and Technology through an Advanced Technology Project subcontract with Phillips Petroleum Company.

Supporting Information Available: A table of the predicted solvation free energies for all 235 molecules in the parametrization suite using the SM5.0R/MIDI! (or PM3), SM5.0R/AM1, OONS/AM1, WE/AM1, and SM5.4/A models is presented (9 pages). Ordering information is given on any current masthead page.

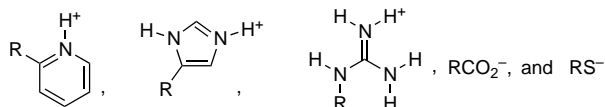
Appendix

This appendix gives the extensions of the SM5.0R model for treating charged groups. As explained in section 3.2, when the model is used with these extensions it is called SM5.05R.

The primary focus of this research was to provide a very fast method to estimate the free energy of solvation of neutral molecules, but we can envisage cases where the model will be useful even in the presence of charges. For charged groups, though, we believe that an explicit representation of the

electrostatics should be used because of the long range nature of Coulomb forces. In keeping with our goal of a very inexpensive model, we have developed a highly simplified model for the electrostatics which might serve to give the approximate contribution of a charged group to the total solvation free energy, and we present this model here, but with the caveat that we recommend using a model with a more complete electrostatic treatment for applications where ions or zwitterions are a *primary* focus.

First we identified several types of charged groups typically encountered in protein modeling, in particular: $R_n\text{NH}_{4-n}^+$ where n is between 0 and 3,



where R is a group of atoms other than a single hydrogen. In addition we added the RO^- charged group because of its general importance. These seven types of charged groups are the only ones for which we parametrized SM5.05R.

In the SM5.0R model, the ΔG_{ENP} term of eq 1 is equal to 0. In the SM5.05R model for solutes with charged groups, we employ a simplified (but nonzero) representation for the electrostatics for the charged groups. The ΔG_{ENP} term can be written as

$$\Delta G_{\text{ENP}} = \Delta E_{\text{EN}} + G_{\text{P}} \quad (13)$$

where ΔE_{EN} is the change in the electronic and nuclear energy of the solute in going from the gas phase to solution, and G_{P} is the polarization free energy. Since this is a rigid model, any change in the nuclear energy of the solute when going from the gas phase to solution must be absorbed into other parameters, and in SM5.05R we also absorb the change in internal electronic energy, so we set $\Delta E_{\text{EN}} = 0$. Using the generalized Born formula, the polarization free energy can be represented as^{3,5,7,56}

$$G_{\text{P}} = -\frac{1}{2} \left(1 - \frac{1}{\epsilon} \right) \sum_{k,k'} q_k q_{k'} \left(r_{kk'}^2 + \alpha_k \alpha_{k'} \exp \left[-\frac{r_{kk'}^2}{d_{kk'} \alpha_k \alpha_{k'}} \right] \right)^{-1/2} \quad (14)$$

where ϵ is the dielectric constant, q_k is the charge on atom k , $r_{kk'}$ is the interatomic distance between atoms k and k' , α_k is an effective Coulomb radius, and $d_{kk'}$ is an empirically optimized constant. The effective Coulomb radius has been defined previously^{3,5,7,56} in terms of a descreening algorithm and is dependent upon an intrinsic Coulomb radius, ρ_0 , for every atom k , whether or not it is charged. Although simplifications in the descreening algorithm have been proposed elsewhere,^{11,23,57,58} they are not used here.

Since we neglect partial charges in neutral molecules without charged groups, a comparable approximation for molecules with charged groups is to neglect partial charges on all the atoms except a specific small number of atoms in the identified charged groups. Thus we model the charge on nitrogen-containing charged groups as being evenly distributed among the hydrogen atoms attached to the nitrogen. Similarly the charge on the oxygen- and sulfur-containing charged groups is distributed among the oxygen or sulfur atoms. In particular, for each of the ions of type $R_n\text{NH}_{4-n}^+$, with $n = 0, 1, 2$, and 3, the charge on each hydrogen atom is taken as $1.00/(4-n)$, and the charge on the nitrogen is 0; for RO^- and RS^- the charge on O and S is -1 ; and for RCO_2^- , the charge on each O is -0.5 . In this

simplified scheme the charge for all other atoms is 0. This creates a model in which the user needs to indicate which atoms have charge and over how many atoms each charge is spread. The only nonzero terms in eq 14 are the diagonal (self-energy) terms and the interactions between charged atoms. However we need intrinsic Coulomb radii even for uncharged atoms because they appear in the descreening algorithm.

For the SM5.05R model, we optimized an intrinsic Coulomb radius for any charged H, O, or S atoms as well as the $d_{kk'}$ parameters d_{HH} and d_{OO} . For the intrinsic Coulomb radii of uncharged atoms we use the van der Waals radii determined by Bondi²² which are already a part of the SM5.0R model. See Table 2. The G_{CDs} term from eq 1 is calculated as it was for the SM5.0R model (even for charged groups) so that the SM5.05R model reduces to the SM5.0R model in the absence of any charged groups.

The training set for the solvation free energies of ions is given in Table 5. It consists of NH_4^+ , six examples of the form RNH_3^+ , eight examples each of the form $\text{RR}'\text{NH}_2^+$ and $\text{R}_2\text{R}'\text{NH}^+$, and four examples each of the forms RCO_2^- , RO^- , and RS^- . Experimental values of the free energy of solvation, when available, were taken from Pearson.⁵⁹ Unfortunately, experimental values for the free energy of solvation are not known in every case. When the experimental value is not available, we write the free energy of solvation as

$$\Delta G_{\text{S}}^{\circ}(\text{R}'\text{X}) = \Delta G_{\text{S}}^{\circ}(\text{RX}) + \Delta \Delta G_{\text{S}}^{\circ}(\text{R}'\text{X}, \text{RX}) \quad (15)$$

where

$$\Delta \Delta G_{\text{S}}^{\circ}(\text{R}'\text{X}, \text{RX}) = \Delta G_{\text{S}}^{\circ}(\text{R}'\text{X}) - \Delta G_{\text{S}}^{\circ}(\text{RX}) \quad (16)$$

Then RX is chosen as the most similar compound for which an experimental value is available. The value of $\Delta G_{\text{S}}^{\circ}(\text{RX})$ was taken from experiment and $\Delta \Delta G_{\text{S}}^{\circ}(\text{R}'\text{X}, \text{RX})$ was calculated by both SM5.4/A and SM5.4/P, and the average value was used to calculate a target value for that ion.

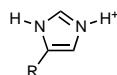
To optimize the intrinsic Coulomb radii and $d_{kk'}$ parameters, we used geometries optimized at the HF/MIDI! level, and we used the genetic algorithm implemented in GAFORTAN-version 1.6.2.⁶⁰ The method was implemented as described previously^{11,23} to maximize the negative of the mean-squared error between the predicted free energies of solvation and the target or experimental free energies of solvation. The final parameters are listed in Table 6.

The SM5.05R model achieves a mean unsigned error from the experimental or target free energies of solvation of 4.3 kcal/mol when applied using the HF/MIDI! geometries and 4.2 kcal/mol when applied using the AM1 geometries. In order to analyze the limitations of the SM5.05R model to predict solvation free energies for charged molecules, we must consider several aspects of the model. Since the SM5.05R model uses a very simplified charge distribution which does not change as a function of the substituent groups, R , attached to the charged group, we might expect difficulties when attached substituents cause the charge distribution to vary from "typical" charges. For example, considering RS^- type ions, the CM1A charge model⁵² at the SM5.4/A geometry predicts the charge on the sulfur atom to be -1.08 ± 0.01 for CH_3S^- , $(n\text{-C}_3\text{H}_7)\text{S}^-$, and $(i\text{-C}_3\text{H}_7)\text{S}^-$ while the charge on the sulfur atom in $\text{C}_6\text{H}_5\text{S}^-$ is predicted to be -0.87 . Thus, we might expect that SM5.05R would have more difficulty predicting the solvation free energy for $\text{C}_6\text{H}_5\text{S}^-$ because the charge differs from the "typical" charge distribution built into the SM5.05 model. The results in the SM5.05R model are consistent with this expectation since the predicted solvation free energies for CH_3S^- , $(n\text{-C}_3\text{H}_7)\text{S}^-$, and

TABLE 5: Calculated and Experimental Free Energies of Solvation (kcal/mol) for Selected Ionic Solutes

	SM5.05R		SM5.4		target ^a
	//AM1	//MIDI!	/A	/P	
	ΔG_S°	ΔG_S°	ΔG_S°	ΔG_S°	
NH ₄ ⁺	-64	-64	-88	-90	-79
CH ₃ NH ₃ ⁺	-62	-62	-76	-81	-70
(<i>n</i> -C ₃ H ₇)NH ₃ ⁺	-61	-61	-71	-76	(-65)
(<i>i</i> -C ₃ H ₇)NH ₃ ⁺	-60	-60	-68	-73	(-62)
(C ₆ H ₅)NH ₃ ⁺	-62	-62	-68	-73	(-62)
(<i>cyclo</i> -C ₆ H ₁₁)NH ₃ ⁺	-59	-60	-67	-69	(-59.5)
(<i>n</i> -C ₆ H ₁₃)NH ₃ ⁺	-60	-60	-71	-73	(-63.5)
(CH ₃) ₂ NH ₃ ⁺	-60	-61	-66	-70	-63
(CH ₃)(<i>n</i> -C ₃ H ₇)NH ₂ ⁺	-59	-59	-62	-66	(-59)
(CH ₃)(<i>i</i> -C ₃ H ₇)NH ₂ ⁺	-58	-58	-59	-64	(-56.5)
(CH ₃)(C ₆ H ₅)NH ₂ ⁺	-60	-60	-58	-62	(-55)
(CH ₃)H ₂ N ₂ C ₃ H ₂ ^{+b}	-54	-55	-60	-66	(-58)
(<i>n</i> -C ₃ H ₇)H ₂ N ₂ C ₃ H ₂ ^{+b}	-53	-53	-55	-62	(-53.5)
(<i>i</i> -C ₃ H ₇)H ₂ N ₂ C ₃ H ₂ ^{+b}	-53	-53	-54	-61	(-52.5)
(C ₆ H ₅)H ₂ N ₂ C ₃ H ₂ ^{+b}	-55	-55	-55	-62	(-53.5)
C ₃ H ₅ NH ^{+c}	-67	-67	-55	-60	-59
(2- <i>n</i> -C ₃ H ₇)C ₅ H ₄ NH ^{+c}	-63	-63	-49	-54	(-53)
(2- <i>i</i> -C ₃ H ₇)C ₅ H ₄ NH ^{+c}	-63	-63	-48	-53	(-52)
(2-C ₆ H ₅)C ₅ H ₄ NH ^{+c}	-65	-64	-47	-52	(-51)
(CH ₃) ₃ NH ⁺	-59	-60	-57	-60	-59
(CH ₃) ₂ (<i>n</i> -C ₃ H ₇)NH ⁺	-58	-58	-54	-57	(-56)
(CH ₃) ₂ (<i>i</i> -C ₃ H ₇)NH ⁺	-56	-56	-52	-55	(-54)
(CH ₃) ₂ (C ₆ H ₅)NH ⁺	-59	-59	-50	-52	(-51.5)
(CH ₃)H ₅ H ₃ C ^{+d}	-57	-57	-62	-68	(-65.5)
(<i>n</i> -C ₃ H ₇)H ₅ N ₃ C ^{+d}	-56	-57	-59	-66	(-63)
(<i>i</i> -C ₃ H ₇)H ₅ N ₃ C ^{+d}	-56	-56	-57	-65	(-61.5)
(C ₆ H ₅)H ₅ N ₃ C ^{+d}	-57	-57	-56	-65	(-61)
CH ₃ CO ₂ ⁻	-78	-77	-77	-79	-77
(<i>n</i> -C ₃ H ₇)CO ₂ ⁻	-76	-75	-74	-76	(-74)
(<i>i</i> -C ₃ H ₇)CO ₂ ⁻	-75	-74	-73	-76	(-73.5)
C ₆ H ₅ CO ₂ ⁻	-77	-76	-71	-74	(-71.5)
CH ₃ O ⁻	-91	-91	-85	-88	-95
(<i>n</i> -C ₃ H ₇)O ⁻	-88	-88	-79	-82	(-89)
(<i>i</i> -C ₃ H ₇)O ⁻	-85	-85	-76	-79	(-86)
C ₆ H ₅ O ⁻	-88	-88	-67	-68	-72
CH ₃ S ⁻	-78	-78	-82	-82	(-79)
(<i>n</i> -C ₃ H ₇)S ⁻	-76	-76	-79	-79	-76
(<i>i</i> -C ₃ H ₇)S ⁻	-74	-74	-76	-77	(-73.5)
C ₆ H ₅ S ⁻	-77	-77	-69	-72	-67
mean unsigned deviation ^e	4.41	4.23	3.92	5.26	

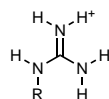
^a The values without parentheses are from ref 59, the values in parentheses are target values estimated using the trends represented in the SM5.4/A and SM5.4/P models. ^b These ions have the form



and the SM5.0R charge is spread evenly about the two hydrogen atoms.

^c Pyridinium ions, with the SM5.05R charge placed on the hydrogen.

^d These ions have the form



and the SM5.0R charge is spread evenly about the five hydrogen atoms.

^e Deviation from the target value.

(*i*-C₃H₇)S⁻ differ from the target values by at most 1 kcal/mol, but the SM5.05R model predicts the solvation free energy of C₆H₅S⁻ to be 10 kcal/mol too negative with respect to the target value.

A second limitation of the SM5.05R model is that the electronic–nuclear reorganization energy of the solute is not included explicitly within the model. Thus, variations in the electronic–nuclear reorganization energy, ΔE_{EN} , of the solute may cause errors in the model. For example, the SM5.4/A model predicts that the polarization free energy of CH₃NH₃⁺ differs by only 3.2 kcal/mol from the polarization free energy

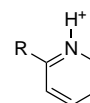
TABLE 6: Parameters Optimized for SM5.05R

type	<i>k</i>	value ^a
ρ_0 (Å)	H	2.25
	O	1.70
	S	2.07
d_{kk} (unitless)	H	1.00
	O	3.00

^a These values apply only to charge-typed atoms. Other ρ_0 and $d_{kk'}$ values are the same as in the SM5.4 models of ref 10. (The only other $d_{kk'}$ values that could possibly be needed are d_{SS} , d_{HO} , d_{HS} , and d_{OS} , and these are all 4.00 in the SM5.4 models; in the present paper we only needed d_{HH} and d_{OO} .)

of (*cyclo*-C₆H₁₁)NH₃⁺ and the predicted G_{CDS} values differ by less than 1 kcal/mol, yet the difference in the SM5.4/A predicted solvation free energies is over 9 kcal/mol. Most of the difference in the SM5.4/A predicted solvation free energies between these two molecules is due to the ΔE_{EN} effects. Since the SM5.05R model is a rigid model it does not include such differences. Thus, the SM5.05R/MIDI! model predicts a difference of only 4 kcal/mol in the solvation free energy between CH₃NH₃⁺ and (*cyclo*-C₆H₁₁)NH₃⁺.

One additional difficulty of the SM5.05R model is for



ions where our solvation free energies are too negative by about 10 kcal. Although we could correct this error by reoptimizing the intrinsic Coulomb radius for hydrogen individually for this type of ion (this would involve increasing ρ_0 from 2.25 to about 2.65 Å for this type of charged group), that would not be within the spirit of the model. We intend for the SM5.05R model to be used in cases where the site of most interest in the molecule is not a charged site. The spirit of the SM5.05R model is to give solvation free energies which are in an appropriate range so that large molecules which happen to have charged groups can be treated inexpensively.

Overall, the SM5.05R model does a fairly good job of predicting solvation free energies which are of an appropriate magnitude for the typed charge groups in the model. It could prove useful for predicting solvation free energies for molecules which have charged groups, but where the main impetus for the research is not the charged groups themselves. If one wishes to study the specific solvation effects centered at charged sites, the use of a model that treats the electrostatics without a rigid charge distribution and uses full geometry optimization in the presence of the solvent is recommended.

References and Notes

- (1) Tomasi, J.; Persico, M. *Chem. Rev.* **1994**, *94*, 2027.
- (2) Cramer, C. J.; Truhlar, D. G. In *Reviews in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH: New York, 1995; Vol. 6, p 1.
- (3) Cramer, C. J.; Truhlar, D. G. *J. Am. Chem. Soc.* **1991**, *113*, 8305.
- (4) Cramer, C. J.; Truhlar, D. G. *Science* **1992**, *256*, 213.
- (5) Cramer, C. J.; Truhlar, D. G. *J. Comput.-Aided Mol. Des.* **1992**, *6*, 629.
- (6) Cramer, C. J.; Truhlar, D. G. *J. Comput. Chem.* **1992**, *13*, 1089.
- (7) Liotard, D. A.; Hawkins, G. D.; Lynch, G. C.; Cramer, C. J.; Truhlar, D. G. *J. Comput. Chem.* **1995**, *16*, 422.
- (8) Giesen, D. J.; Storer, J. W.; Cramer, C. J.; Truhlar, D. G. *J. Am. Chem. Soc.* **1995**, *117*, 1057.
- (9) Giesen, D. J.; Cramer, C. J.; Truhlar, D. G. *J. Phys. Chem.* **1995**, *99*, 7137.

- (10) (a) Chambers, C. C.; Hawkins, G. D.; Cramer, C. J.; Truhlar, D. G. *J. Phys. Chem.* **1996**, *100*, 16385. (b) Giesen, D. J.; Gu, M. Z.; Cramer, C. J.; Truhlar, D. G. *J. Org. Chem.* **1996**, *61*, 8720. (c) Giesen, D. J.; Chambers, C. C.; Cramer, C. J.; Truhlar, D. G. *J. Phys. Chem. B* **1997**, *101*, 2061.
- (11) Hawkins, G. D.; Cramer, C. J.; Truhlar, D. G. *J. Phys. Chem.* **1996**, *100*, 19824.
- (12) Chothia, C. *Nature* **1974**, *248*, 338.
- (13) Eisenberg, D.; McLachlan, A. D. *Nature* **1986**, *319*, 199.
- (14) Wesson, L.; Eisenberg, D. *Protein Sci.* **1992**, 227.
- (15) Villa, J.; Williams, R. L.; Vasquez, M.; Scheraga, H. A. *Proteins Struct. Funct. Genet.* **1991**, *10*, 199.
- (16) Schiffer, C. A.; Caldwell, J. W.; Kollman, P. A.; Stroud, R. M. *Mol. Simul.* **1993**, *10*, 121.
- (17) Kurochkina, N.; Lee, B. *Protein Eng.* **1995**, *8*, 437.
- (18) Juffer, A. H.; Eisenhaber, F.; Hubbard, S. J.; Walther, D.; Argos, P. *Protein Sci.* **1995**, *4*, 2499.
- (19) Gschwend, D. A.; Kuntz, I. D. *J. Comput.-Aided Mol. Des.* **1996**, *10*, 330.
- (20) Ooi, T.; Oobatake, M.; Nemethy, G.; Scheraga, H. A. *Proc. Natl. Acad. Sci. U.S.A.* **1987**, *84*, 3086.
- (21) Tuñón, I.; Silla, E.; Pascual-Ahuir, J. L. *Chem. Phys. Lett.* **1993**, *203*, 289.
- (22) Bondi, A. J. *Phys. Chem.* **1964**, *68*, 441.
- (23) Hawkins, G. D.; Cramer, C. J.; Truhlar, D. G. *Chem. Phys. Lett.* **1995**, *246*, 122.
- (24) Lee, B.; Richards, F. M. *J. Mol. Biol.* **1971**, *55*, 379.
- (25) Hermann, R. B. *J. Phys. Chem.* **1972**, *76*, 2754.
- (26) Stillinger, F. H.; Weber, T. A. *Phys. Rev. A* **1983**, *28*, 2408.
- (27) Allred, A. L. *J. Inorg. Nucl. Chem.* **1961**, *17*, 215.
- (28) Allen, L. C.; Huheey, J. E. *Ibid.* **1980**, *42*, 1523.
- (29) Wauchope, R. D.; Haque, R. *Can. J. Chem.* **1972**, *50*, 133.
- (30) Hine, J.; Mookerjee, P. K. *J. Org. Chem.* **1975**, *40*, 287.
- (31) Wolfenden, R. *Biochemistry* **1978**, *17*, 201.
- (32) Cabani, S.; Gianni, P.; Mollica, V.; Lepori, L. *J. Solution Chem.* **1981**, *10*, 563.
- (33) Wagman, D. D. *J. Phys. Chem. Ref. Data* **1982**, *11*, Suppl. No. 2.
- (34) Abraham, M. H.; Whiting, G. S.; Fuchs, R.; Chambers, E. J. *J. Chem. Soc., Perkin Trans. II* **1990**, 291.
- (35) Han, P.; Bartles, D. M. *J. Phys. Chem.* **1990**, *94*, 7294.
- (36) Suleiman, D.; Eckert, C. A. *J. Chem. Eng. Data* **1994**, *39*, 692.
- (37) Leo, A. J. *Masterfile* (1994) from MedChem Software, BioByte Corp., P.O. 517, Claremont, CA 91711-0157.
- (38) Ferguson, D. M.; Pearlman, D. A.; Swope, W. C.; Kollman, P. A. *J. Comput. Chem.* **1992**, *13*, 362.
- (39) Ben-Naim, A. *Statistical Thermodynamics for Chemists and Biochemists*; Plenum: New York, 1992, p 421.
- (40) Easton, R. E.; Giesen, D. J.; Welch, A.; Cramer, C. J.; Truhlar, D. G. *Theor. Chim. Acta* **1996**, *93*, 281.
- (41) Möller, C.; Plesset, M. S. *Phys. Rev.* **1934**, *46*, 618.
- (42) Pople, J. A.; Seeger, R.; Krishnan, R. *Int. J. Quantum Chem. Symp.* **1977**, *11*, 149.
- (43) Krishnan, R.; Pople, J. A. *Int. J. Quantum Chem.* **1978**, *14*, 91.
- (44) Dunning, T. H. *J. Chem. Phys.* **1989**, *90*, 1007.
- (45) Dewar, M. J. S.; Zoebisch, E. G.; Healy, E. F.; Stewart, J. J. P. *J. Am. Chem. Soc.* **1985**, *107*, 3902.
- (46) Dewar, M. J. S.; Zoebisch, E. G. *J. Mol. Struct. (THEOCHEM)* **1988**, *180*, 1.
- (47) Dewar, M. J. S.; Yate-Ching, Y. *Inorg. Chem.* **1990**, *29*, 3881.
- (48) Stewart, J. J. P. *J. Comput. Chem.* **1989**, *10*, 209.
- (49) Stewart, J. J. P. *J. Comput. Chem.* **1989**, *10*, 221.
- (50) Kyte, J.; Doolittle, R. P. *J. Mol. Biol.* **1982**, *157*, 105.
- (51) Wolfenden, R.; Andersson, L.; Cullis, P. M.; Southgate, C. C. B. *Biochemistry* **1981**, *20*, 849.
- (52) Storer, J. W.; Giesen, D. J.; Cramer, C. J.; Truhlar, D. G. *J. Comput.-Aided Mol. Des.* **1995**, *9*, 87.
- (53) Tuñón, I.; Silla, E.; Pascual-Ahuir, J. L. *Protein Eng.* **1992**, *5*, 715.
- (54) Jackson, R. M.; Sternberg, M. J. E. *Protein Eng.* **1994**, *7*, 371.
- (55) Hawkins, G. D.; Lynch, G. C.; Giesen, D. J.; Rossi, I.; Storer, J. W.; Liotard, D. A.; Cramer, C. J.; Truhlar, D. G. *QCPE Bull.* **1996**, *16*, 11.
- (56) Still, W. C.; Tempczyk, A.; Hawley, R. C.; Hendrickson, T. *J. Am. Chem. Soc.* **1990**, *112*, 6127.
- (57) Schaefer, M.; Froemmel, C. *J. Mol. Biol.* **1990**, *216*, 1045.
- (58) Schaefer, M.; Karplus, M. *J. Phys. Chem.* **1996**, *100*, 1578.
- (59) Pearson, R. G. *J. Am. Chem. Soc.* **1986**, *108*, 6109.
- (60) Carroll, D. L. *AIAA J.* **1996**, *34*, 338.