



# An extrapolation method for computing protein solvation energies based on density fragmentation of a graphical surface tessellation

Lochana C. Menikarachchi, José A. Gascón\*

Department of Chemistry, University of Connecticut, 55 North Eagleville Rd., Unit 3060, Storrs, CT 06269, United States

## ARTICLE INFO

### Article history:

Received 2 March 2011

Received in revised form 26 May 2011

Accepted 2 June 2011

Available online 13 June 2011

### Keywords:

QM/MM

Moving-domain QM/MM

Continuum solvation

Molecular surface

## ABSTRACT

Modeling chemical events inside proteins often require the incorporation of solvent effects via continuum polarizable models. One of these approaches is based on the assumption that the interface between solute and solvent acts as a conductor. Image charges are added on the molecular surface to satisfy the appropriate conductor boundary conditions in the presence of solute charges. As in the case of other polarizable continuum models that are based on surface tessellation, the simplest implementation of this approach is often limited to several hundred atoms due to a matrix inversion, which scales as the cube of the number of tesserae. For larger systems, approaches that use iterative matrix solvers coupled to fast summation methods must be used. In the present work, we develop a self-consistent approach to obtain conductor-like screening charges suitable for applications in proteins. The approach is based on a density fragmentation of a graphical surface tessellation. This method, although approximate, provides a straightforward scheme of parallelization, which can in principle be added to existing linear scaling implementations of conductor-like models. We implement this method in conjunction with a fixed charge model for the protein, as well as with a moving domain QM/MM description of the protein. In the latter case, the overall result leads to a charge distribution within the protein determined by self-polarization and polarization due to solvent.

© 2011 Elsevier Inc. All rights reserved.

## 1. Introduction

The accurate treatment of electrostatic effects due to solvation in biological macromolecules is a common problem in all aspects of protein modeling: computational protein design [1], structure based drug design [2], and quantum mechanical (QM) treatment of enzymatic reactions [3]. Focusing on this last aspect, there are circumstances where certain chemical events require a high level of QM calculation, rendering the explicit treatment of solvent impractical. In particular, quantum mechanics/molecular mechanics (QM/MM) methods are an important tool to understand how the protein cavity affects a quantum mechanical property or a chemical reaction profile in the active site [4–6]. Explicit treatment of water molecules with a high level QM theory is often prohibitively expensive because of required averaging over times that are much longer than the dipole relaxation time of water (about 10 ps at 300 K) [7,8]. While it remains challenging to describe solvation effects for small molecules, the development of continuum solvation models and the integration of such models with quantum chemistry methods over the last 20 years, has shown remarkable success [9]. The most commonly used models, with the most prevalent

implementation in electronic structure software, are the polarizable continuum model (PCM) by Tomasi and co-workers [9,10], the conductor-like screening model (COSMO) by Klamt and co-workers [11,12], and those involving the simultaneous solution of the linearized Poisson–Boltzmann (PB) equation and its integration with ab initio quantum mechanics [13–15]. Both PCM and COSMO involve tessellation of the molecular surface, where each tessera is assigned a screening charge. The COSMO model is particularly attractive because screening charges can be obtained, in principle, by simple inversion of an  $M \times M$  matrix, where  $M$  is the number of tesserae. On the other hand, implementation of such a direct inversion method becomes rapidly intractable for protein-sized systems (for both MM and QM/MM approaches) due to the  $O(M^3)$  scalability of matrix inversion and the need of large matrix storage  $O(M^2)$  [9,16]. To overcome these limitations, linear scaling algorithms which make use of iterative matrix solvers [17,18] coupled with fast summation methods [19] have been successfully used for protein systems with several thousands of atoms. In spite of this progress, the implementation of tessellation methods such as PCM and COSMO for very large biological macromolecules remains a challenging problem, particularly in the context of QM/MM.

Efforts that focus on integrating implicit solvation effects in QM/MM methods have emerged in recent years. Dinner et al. [16] explored the use of a method derived from free energy perturbation simulations using molecular mechanics [20,21]. In this approach,

\* Corresponding author.

E-mail address: [jose.gascon@uconn.edu](mailto:jose.gascon@uconn.edu) (J.A. Gascón).

charges of exposed groups are scaled to avoid distortion of structures during the simulation. Grid-based continuum electrostatic methods [22,23] are then used to estimate the energy required to return the scaled groups to their normal state and to transfer the structure to bulk solvent. In the work by Schaefer et al. [24], a protocol was developed that extends the generalized solvent boundary potential (GSBP) [25] method to QM/MM simulations. Implementation of this method was carried out with the self-consistent-charge density-functional tight-binding method (SCC-DFTB) at the QM level with a non-polarizable force field on the MM part. More recently, Merz and co-workers have developed a method that combines a QM/MM with a Poisson–Boltzmann treatment [26]. In their approach, self-consistency due to solvent polarization was required in a small QM region while all charges outside the region are fixed MM charges. A semiempirical method with a divide and conquer algorithm was employed to solve the QM subsystem [27]. A similar technique has also been employed to integrate the implicit PCM solvation method with Morokuma's ONIOM method (ONIOM-PCM) [28] as well as with GAMESS [29]. A less computationally demanding solvation model which does not involve surface tessellation is the generalized-born surface-area (GBSA) approach, originally reported by Still et al. [30]. Integration of a GBSA model with hybrid QM/MM potentials for macromolecules has been presented by Field and Pellegrini [31]. Their effort was focused on a new approach for the calculation of Born radii which required a large number of parametrizations to include MM and QM atoms in both small and large molecules.

In the current work, the problem of integrating continuum solvation with MM and QM/MM is revisited. The motivation is to increase scalability in the computation of apparent surface charges in tessellation-based continuum models. An iterative scheme that computes the free energy of electrostatic solvation via a density-domain fragmentation (DDF) of the surface tesserae is proposed. Such a fragmentation leads to a straightforward parallelization scheme. We apply this approach to modify the COSMO model and show its feasibility to treat macromolecular systems. An ideal integration between implicit solvation and QM/MM should also allow for polarization effects on the molecular mechanics region of the protein (due to solvent and the protein itself). This is accomplished by coupling the Moving-Domain QM/MM method [32,33] with DDF-COSMO in a self-consistent manner. A new algorithm of tessellation, which was originally designed in the context of molecular surface rendering, is also implemented. This algorithm shows improved efficiency with comparison to standard methods, originally designed in the context of continuum solvation.

## 2. Theory and implementation

In conductor-like screening models [11,12], the solute-solvent boundary is assumed to behave like a conductor surface. The reaction field provided by this conducting surface is described by means of apparent polarization charges distributed on a set of mosaics (tesserae), placed on the surface of the solute molecule. The advantage of this model, as opposed to dielectric-type continuum models, is that the required boundary condition (a surface of zero electrostatic potential) is easily met by placing image charges, and therefore does not require the explicit solution of the Poisson equation [34]. Instead, the solution of these apparent polarization charges is analytical. However, as it is shown below, the most time consuming part of this calculation involves inversion of a matrix of dimension  $M_0 \times M_0$ , where  $M_0$  is the number of tesserae. Typically,  $M_0$  is between 100 and 500 for molecules with a few atoms up to 20 atoms. This model has been successfully applied to relatively small systems providing hydration energies in very good agreement with experimental results. Implementing this implicit

solvation model for a large size protein would require the inversion and storage of a matrix based on a surface as large as a few hundred thousand tesserae. Inversion or storage of a matrix with such dimensions is computationally intractable. However, this problem has been solved by the use of iterative matrix solvers coupled to fast summation methods allowing calculations up to several thousands of atoms [17,19,35]. The iterative solution proposed here, although approximate, aims to be a complement of such linear scaling algorithms. An additional layer of parallelization is also possible on top of already parallelizable iterative solvers and fast summation methods.

The starting point in the conductor-like solvation model is to write the electrostatic component of the screening energy of a solute (represented by the protein partial charges) as:

$$E = \sum_i \int z_i \frac{1}{|r_i - r_q|} q dS + \frac{1}{2} \iint q \frac{1}{|r_q - r_{q'}|} q' dS dS' \quad (1)$$

where the  $z_i$  are the partial atomic charges in the solute molecule, at position  $r_i$ , and  $q$  represents the surface charge in an element of surface area  $dS$  centered at position  $r_q$ . By discretizing the molecular surface, Eq. (1) can be conveniently written in matrix notation as:

$$E(\mathbf{Q}) = \mathbf{Z}^T \mathbf{B} \mathbf{Q} + \frac{1}{2} \mathbf{Q}^T \mathbf{A} \mathbf{Q} \quad (2)$$

where  $\mathbf{Z}$  is an  $N$ -dimensional vector containing the solute charges (upper script  $T$  indicates transpose),  $\mathbf{B}$  is the coulomb matrix describing the interaction between solute charges and solvent charges ( $B_{ij} = 1/|r_i - r_j|$ ),  $\mathbf{Q}$  is an  $M_0$ -dimensional vector containing the solvent charges (to be determined), and  $\mathbf{A}$  is the coulomb matrix describing the interaction between solvent charges:  $A_{ij} = 1.07 \sqrt{4\pi/a_i}$ ,  $A_{ij} = 1/|r_i - r_j|$ , where  $a_i$  is the area of tessera  $i$ .

For small systems,  $\mathbf{Q}$  can be easily obtained by imposing the condition:

$$\nabla E_Q = 0 \quad (3)$$

This condition leads to the linear equation:

$$\mathbf{A} \mathbf{Q} = -\mathbf{B}^T \mathbf{Z} \quad (4)$$

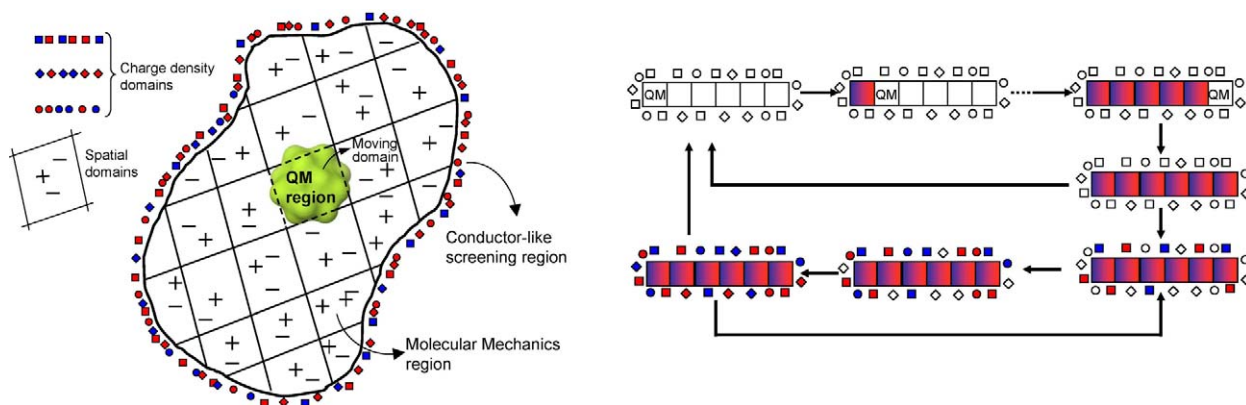
whose solution can be obtained by direct inversion of  $\mathbf{A}$  or by iterative numerical procedures. The solution  $\mathbf{Q}$  is then scaled by the correction factor  $f(\epsilon) = (\epsilon - 1)/(\epsilon + 1/2)$  to extend the theory to finite values of a dielectric with an error of less than  $1/(2\epsilon)$  [12]. Finally, the electrostatic solvation energy is computed as:

$$E(\mathbf{Q}) = \frac{1}{2} \mathbf{Z}^T \mathbf{B} \mathbf{Q} \quad (5)$$

We now consider the case in which the surface is large enough so that the solution of linear equations is impractical even with numerical iterative solvers. The surface tesserae are partitioned into  $L$  domains of dimension equal to  $M$  ( $M_0 = M \times L$ ). Such partitioning is done not in spatial domains, but rather in the density domain (see Fig. 1, left). Therefore, the charges within each domain are distributed along the entire surface, using a subset of the entire tesserae. Provided that  $L$  is large enough, the linear equations can be easily solved, even by direct inversion. Thus, the charge density of a particular domain is  $L$ -times smaller than the density of all screening charges before the partitioning.

Charges are then obtained via an iteration procedure where at each step condition [3] is imposed on only one domain, domain  $k$ . Thus, at any given step the screening energy can be written by explicitly separating the contribution of charges from domain  $k$  ( $\mathbf{q}^k$ ) and charges from the rest of the domains ( $\mathbf{q}^0$ ):

$$E(\mathbf{q}^k, \mathbf{q}^0) = \mathbf{Z}^T \mathbf{B}^0 \mathbf{q}^0 + \mathbf{Z}^T \mathbf{B}^k \mathbf{q}^k + \frac{1}{2} (\mathbf{q}^0)^T \mathbf{A}^0 \mathbf{q}^0 + \frac{1}{2} (\mathbf{q}^k)^T \mathbf{A}^k \mathbf{q}^k + \frac{1}{2} (\mathbf{q}^0)^T \mathbf{A}^{0k} \mathbf{q}^k \quad (6)$$



**Fig. 1.** Essential parts in the MOD-QM/MM method and its integration with the DDF-COSMO model. On the right, each square in a bin represents a spatial domain in a protein. Coloring of the various elements represents update in charge.

where  $\mathbf{B}^0$  is the coulomb matrix between solute charges and the solvent charges for all domains other than domain  $k$ ,  $\mathbf{B}^k$  is the coulomb matrix between solute charges and solvent charges in domain  $k$ ,  $\mathbf{A}^0$  is the coulomb matrix between solute charges within the domains other than  $k$ ,  $\mathbf{A}^k$  contains the interactions within charges of domain  $k$ , and  $\mathbf{A}^{0k}$  contains the interactions between charges in domain  $k$  and the rest of the domains. Solving the equation  $\nabla E|_{\mathbf{q}^k} = 0$  we obtain:

$$\mathbf{q}^k = -(\mathbf{A}^k)^{-1} \left( \mathbf{B}^k \mathbf{Z} + \frac{1}{2} (\mathbf{A}^{0k})^T \mathbf{q}^0 \right) \quad (7)$$

By substituting Eq. (7) in Eq. (6) we obtain:

$$E(\mathbf{q}^k, \mathbf{q}^0) = -\frac{1}{2} (\mathbf{q}^k)^T \mathbf{A}^k \mathbf{q}^k + \frac{1}{2} (\mathbf{q}^0)^T \mathbf{A}^0 \mathbf{q}^0 + \mathbf{Z}^T \mathbf{B}^0 \mathbf{q}^0 \quad (8)$$

The process is then repeated sequentially for all other domains until convergence is achieved either by monitoring the total energy or the charge distribution. In summary, the procedure to obtain an approximation to  $\mathbf{Q}$ , after tessellation has been carried out, follows these steps:

- Partition the tesserae into  $L$  density domains. All screening charges are initially zero.
- From  $k = 1$  to  $k = L$ , solve for  $\mathbf{q}^k$  in Eq. (7), updating at each step  $\mathbf{q}^0$  using an iterative matrix solver such as the biconjugate gradient stabilized method.
- Evaluate  $\mathbf{Q} = (\mathbf{q}^1, \mathbf{q}^2, \dots, \mathbf{q}^k, \dots, \mathbf{q}^M)$ .
- Compute the solvation energy according to Eq. (8).
- If the solvation energy is converged, stop the algorithm; otherwise repeat from step ii.

**Integration with MM and QM/MM:** The procedure described above can be integrated with molecular mechanics because point charges are readily available from any standard force field. To integrate the solvation model with a QM/MM approach we use solute charges taken from the force field for the MM region and derive electrostatic potential charges (ESP) for the QM region assuming the protein is in vacuum. In response to these solute charges, the screening charges obtained by the DDF-COSMO method are then put back into the QM/MM Hamiltonian as external charges. This procedure is iterated until convergence in the solvation energy is obtained, which usually requires two to three cycles. Notice that this iteration is in addition to the iteration described in the previous section regarding self-consistency within the DDF method.

The treatment of implicit solvation via the proposed density-domain fragmentation approach is intrinsically related to the space-domain decomposition schemes such as the moving-domain QM/MM (MOD-QM/MM) protocol developed recently [33,36]. The

latter method aims to obtain a polarized MM region by moving the QM region (as part of a QM/MM calculation) throughout the protein and recomputing charges at each step, leading to a QM-based self-consistent charge distribution for the whole macromolecular system. DDF-COSMO can be naturally integrated with MOD-QM/MM by a double iteration in both the spatial domains and the density domains (Fig. 1 right depicts this integration). In Section 4 we present calculations of the electrostatic component of solvation energy in which the protein charges are either obtained from an MM force field, QM, or MOD-QM/MM calculations.

### 3. Computational details and graphical surface tessellation

All QM calculations (either in QM/MM or MOD-QM/MM) were carried out at the Hartree-Fock level with the basis set 6–31 g\*. Because our interest is centered around the performance and feasibility of the DDF implicit solvation model, this theory level provides enough accuracy and efficiency to carry out the desired test calculations. To benchmark the method, the sole focus was the electrostatic contribution to the solvation energy. Hydrophobic terms could eventually be added via the solvent accessible surface area (SASA) approach [37]. Charges for the MM region were taken from the Amber force field [38] for the case in which the MM region was not polarized. For cases where molecular regions were treated quantum mechanically, charges on those regions were obtained via the electrostatic potential charges (ESP) method [39]. The molecular mechanics package Tinker [40] was used for the assignment of force field parameters in the MM region. The program Gaussian 09 [41] was used for the QM/MM calculations, with the ONIOM [42–44] method and electronic embedding. MOD-QM/MM calculations were carried out with our program MODQ3M [36]. The coupling between MODQ3M and DDF-COSMO was performed with in-house scripts. Van der Waals radii were taken from the values used in the COSMO implementation of Klamt et al. [45]. The biconjugate gradient stabilized method as implemented in the iterative matrix solver package SPARSKIT2 [46] was used for solving linear equations. Because the purpose is only to test the solvation method, calculations were done as single point energies using X-ray configurations of the polypeptides without additional minimization.

Tessellation of the molecular surface was performed with the MSMS algorithm (also a program) by M. Sanner [47] with a density of 2.0 vertex/Å<sup>2</sup> and a probe radius of 1.4 Å. Tessellation was created on the solvent excluded surface. This algorithm has proven robust in handling singularities for protein-sized systems and to scale as  $O(N \log(N))$  for  $N$  atoms. Although the MSMS algorithm was originally conceived for graphical rendering of molecular surfaces, we have found that the this method performs equally as well or better

than algorithms designed exclusively for implicit solvent calculations, such as GEPOL93 [48]. Indeed, the MSMS algorithm requires consistently fewer tesseræ than GEPOL93 while yielding the same level of accuracy. Moreover, for some systems the SPARSKIT algorithm requires considerably fewer iterations, while maintaining equivalent accuracy, to solve the linear equations using GEPOL93 compared to MSMS (see [supporting information](#)).

As shown in the next section, electrostatic solvation energies obtained via DDF with  $L = 10$  correlate almost linearly with a full COSMO calculation. The value of 10 was used to reduce 10-fold the task of solving the linear equations in the biconjugate method. More precisely, by timing the total time needed to obtain the solvation energy with DDF and with full COSMO, we find that the gain is in fact between 6 and 10 fold. Although one would expect even larger gains with larger values of  $L$ , as  $L$  increases there is a slight deterioration on the quality of the correlation. Thus, we have also chosen  $L$  as a compromise between efficiency and accuracy. Additional details on the effect of  $L$  on the accuracy of the DDF-COSMO energy is presented in the supplementary information. Encouraged by the excellent correlation, a set of linear and multi-linear regression parameters was generated to obtain an approximation to the full COSMO electrostatic solvation energy (i.e.,  $L = 1$ ). The quality of the multi-linear regression is assessed through a K-fold cross validation of the data set. The electrostatic solvation energies calculated with COSMO and DDF are also compared with Poisson Boltzmann electrostatic solvation energies. The Poisson Boltzmann (PB) calculations were done for a sub set of proteins with the software APBS 1.3 [49].

## 4. Results

Table 1 presents the electrostatic component of the solvation free energy for various proteins and protein fragments. The third column specifies the set of amino-acids with charges computed self-consistently with DDF-COSMO either via QM/MM or MOD-QM/MM. Blank boxes indicate charges obtained directly from Amber, without self-consistent polarization of the protein or fragment. The purpose of the selected set is to produce solvation energies more or less equally distributed along a large range of solvation energies ( $\sim 2750$  kcal/mol). With this set of data we determined the correlation between a full COSMO calculation (1 domain) and the DDF-COSMO calculation (10 domains). Such a correlation can ultimately be used to extrapolate the full COSMO solvation energy from DDF.

Energies computed with DDF with 10 domains have a remarkable correlation with a full COSMO calculation (see Fig. 2, left panel). This is a very encouraging result, as it suggests that for systems with thousand of atoms, solvation energy could be easily computed with the appropriate fragmentation and be corrected with predetermined linear or multi-linear regression parameters. Considering that DDF-COSMO could be eventually implemented with solutes treated as MM, QM/MM or MOD-QM/MM, we have included, as part of the training set, all these types of solute charges. In fact, we observe that before and after regression is applied, all sets fall under the same correlation (no multimodal distributions are apparent, see Fig. 2). This shows that the correlation between the fragmentation method and the standard COSMO method is independent of the level of theory used to treat the solute charges. Although the correlation parameters ( $a$ ,  $b$ ,  $c$  and  $d$ ) depend on  $L$ , we envision that  $L = 10$  would be a standard value, providing an order of magnitude reduction in the cost employing the iterative matrix solvers and three orders of magnitude using direct matrix inversion. The multiple linear regression parameters obtained for the full data set are presented in Table 2.

### 4.1. Convergence analysis

The cycle connecting MOD-QM/MM with DDF-COSMO converges in two or three iterations, a property that is independent of the size of the system. Here, we consider in more detail the convergence properties of the self-consistent polarization cycle within the DDF approach. To test how the computed solvation energy converges, we considered three polypeptides (listed by PDB entry) of three different sizes: 3ftk (100 atoms), 1q8h (538 atoms), and 1dsl (1473 atoms). As shown in Fig. 3, the DDF-COSMO energy converges by the second or third cycle. We have further proved that this result holds for even larger molecules. This means that the solution of linear equations is performed at most  $3 \times L$  times, acting only as a prefactor in the scaling properties.

Because of the fragmentation approach it would seem natural to implement a parallelization procedure in which each processor solves Eq. (7) for a different domain. However, solution of Eq. (7) for a given domain depends on the screening charges of all other domains. Therefore, the parallel procedure should, in principle, take longer to converge than the serial version, since it uses older versions of surface charges. Remarkably, the parallel implementation of DDF-COSMO takes one or, at most, two more cycles than the serial version to converge. Considering the fact that the calculations on all domains are done at the same time in a multi-core computer this would lead to an enormous reduction in computational time. Fig. 4 illustrates the convergence properties of serial and parallel algorithms for 1dsl (19,180 tesseræ).

### 4.2. Cross validation

A rigorous K-fold cross validation method was employed to assess the quality of the multi-linear regression within the range of protein sizes considered (see [supporting information](#) for further details). The data points used for the multi-linear regression were randomly partitioned into 10 sub samples each containing 8 data points. A single sub sample is retained for validating the regression while the rest is used as the training data set. The procedure is repeated for all sub samples. The root mean square error (RMSE) of the electrostatic solvation energies predicted with smaller training data sets with respect to the regression with the full data set is taken as a measure of predictive accuracy ( $\sim 4$  kcal/mol). Therefore, as far as predicting absolute solvation free energies of biomacromolecules, this value would typically represent less than 1% error. On the other hand, if the computation of relative energies was desired (e.g., in the evaluation of protein–protein association within similar complexes), uncertainty of 4 kcal/mol would be certainly critical. In such cases, the production of training sets within a particular family of complexes would be necessary, so that the magnitudes of solvation energies remain within a smaller range. Under these conditions the standard error will be considerably smaller. This is in fact the philosophy of linear response free energy methods, which use multi-linear regression to predict ligand affinity using training sets within a certain family of proteins [50]. As further validation of the present approach, DDF-COSMO was compared with the Poisson–Boltzmann (PB) method, as a way of comparison with models that are not based on a conductor-like assumption and do not require surface tessellation. DDF-COSMO correlates well with PB for a representative sample of the proteins studied (Fig. 5).

### 4.3. Influence of solute characteristics on DDF-COSMO energy

As shown in Table 1, the DDF solvation energy, prior to any extrapolation, gives a consistent underestimation of the full COSMO solvation energy. The ratio between these two energies ranges from 0.5 to 0.85. In order to discern what protein properties give rise to



**Table 1**  
Electrostatic solvation free energy in kcal/mol using COSMO and DDF-COSMO.

Protein <sup>a</sup>	No. of tesserae	QM region <sup>b</sup>	DDF COSMO (no. of domains = 10)	Multiple regression	Full COSMO
2onx	2112	Residue 1 All atoms	−103.95	−175.07	−172.53
3dg1	2114		−93.28	−153.49	−153.13
1itt	2194		−102.23	−153.86	−158.92
3fpo	2520		−200.09	−312.74	−293.33
3fpo	2520		−195.76	−306.07	−286.62
3fpo	2520		−191.38	−300.67	−285.86
1gcn (9–14)	2610		−147.61	−243.04	−237.69
1yjo	2632		−111.80	−181.44	−186.48
3ftk	2662		−109.12	−171.80	−180.87
3fva	2862		−114.12	−185.74	−188.58
1etl	3134	Residue 1 All atoms	−115.27	−178.61	−183.72
1hje	3516		−126.68	−207.24	−196.60
1not	3734		−243.11	−385.44	−363.39
1not	3734		−227.59	−358.98	−339.76
1not	3734		−200.63	−326.09	−370.43
1vrz	4264		−78.98	−142.13	−138.97
1pen	4286		−227.71	−334.56	−329.67
1akg	4318		−230.47	−346.89	−338.91
1rpb	4912		−67.04	−113.52	−114.61
1gcn (15–29)	6294		−264.87	−417.02	−417.87
2bf9 (1–20)	6338	Residue 1 All atoms	−423.10	−567.51	−573.72
3e4h	6566		−92.45	−169.44	−178.11
1kyc	6600		−288.77	−481.90	−464.60
1pef	6802		−244.02	−401.52	−380.55
1edn	6936		−256.61	−405.37	−394.03
1jo8 (1–25)	7688		−1126.64	−1350.77	−1343.44
2bf9 (1–25)	7980		−525.69	−715.60	−700.20
2bf9 (2–28)	8564		−480.42	−649.83	−654.93
2bf9 (1–30)	9194		−485.50	−677.97	−669.47
3e21	9196		−1032.21	−1278.59	−1257.64
3e7u	9472	Residue 1 All atoms	−477.34	−750.12	−723.98
1amc	9608		−547.58	−855.56	−829.25
1jo8 (4–32)	9794		−1829.56	−2181.11	−2154.20
1jo8 (4–32)	9794		−1817.33	−2163.29	−2137.07
1jo8 (4–32)	9794		−1737.34	−2072.44	−2046.00
1crn	9846		−157.79	−278.29	−286.62
1cbn	9974		−152.85	−259.38	−278.79
2bf9 (5–35)	10,010		−376.51	−566.82	−588.67
1q8h	10,032		−518.44	−758.89	−759.17
1cnr	10,076		−154.92	−273.75	−280.93
1q9b	10,094	Residue 1 All atoms	−473.35	−725.03	−724.07
1ab1	10,126		−157.89	−273.39	−287.15
1ab1	10,126		−145.32	−241.70	−267.13
1ab1	10,126		−146.98	−256.23	−283.92
2erl	10,216		−720.32	−919.19	−904.32
1jo8 (5–35)	10,550		−2104.46	−2436.39	−2418.12
1gcn	10,946		−372.50	−601.43	−602.50
1gcn	10,946		−367.27	−590.90	−594.30
1gcn	10,946		−370.53	−595.94	−607.92
2bf9	10,986		−447.07	−684.27	−689.73
2bf9	10,986	Residue 1 All atoms	−441.46	−675.09	−681.32
2bf9	10,986		−374.11	−603.78	−611.54
1jo8 (1–35)	11,154		−2171.95	−2480.43	−2505.64
1jo8(8–48)	11,186		−1999.33	−2337.31	−2312.52
2gkr	11,904		−443.32	−719.50	−716.51
1aie	12,096		−553.42	−826.12	−827.76
1jo8 (1–42)	12,266		−2285.71	−2606.68	−2632.21
1zlm	12,270		−507.30	−833.65	−932.89
1jo8	12,656		−2315.02	−2673.04	−2687.29
2qmt	12,930	Residue 1 All atoms	−592.44	−859.10	−871.58
2erw	12,982		−496.27	−732.65	−732.56
1yp5	13,026		−481.36	−776.91	−771.24
1oot	13,032		−670.69	−1053.96	−1033.10
2vvk	13,138		−598.52	−876.47	−894.78
2ovo	13,420		−432.98	−714.80	−715.19
1fas	13,450		−612.07	−874.66	−881.51
1bpi	13,458		−876.68	−1173.58	−1172.79
3ca7	13,520		−575.75	−857.37	−867.76
3ca7	13,520		−554.73	−844.98	−839.13
3ca7	13,520	Residue 1 All atoms	−521.62	−819.17	−826.64
1vb0	13,638		−438.85	−719.86	−711.58
1mhn	13,698		−839.11	−1193.09	−1182.88
1fan	13,852		−874.44	−1180.05	−1168.46
4pti	13,890		−887.56	−1189.41	−1179.72
2fma	13,952		−549.99	−875.04	−872.68
1enh	14,066		−1405.62	−1767.67	−1783.79

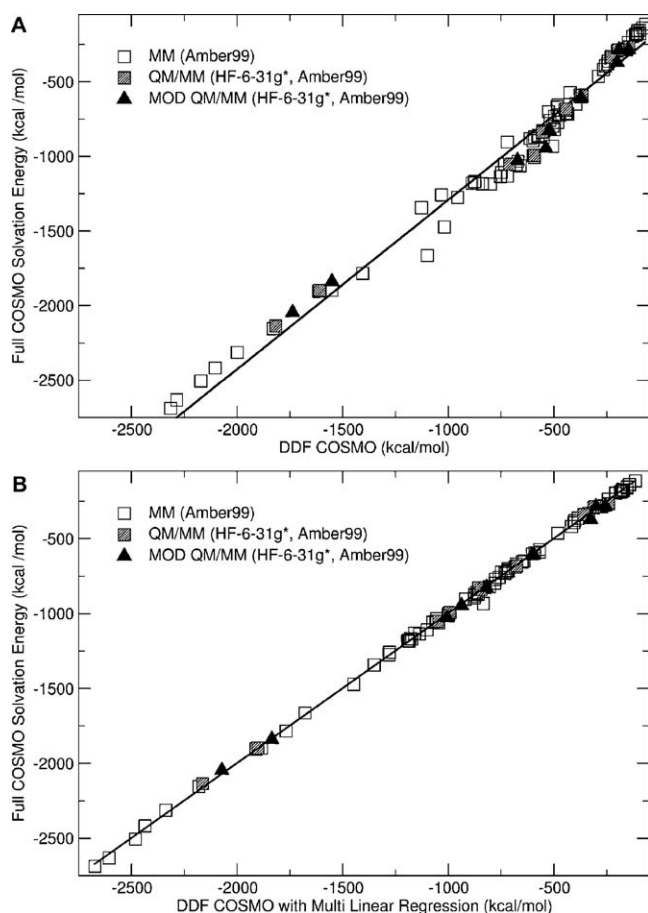
Table 1 (Continued)

Protein <sup>a</sup>	No. of tesserae	QM region <sup>b</sup>	DDF COSMO (no. of domains = 10)	Multiple regression	Full COSMO
1bti	14,178		−955.83	−1279.51	−1275.48
1guu	14,264		−801.38	−1185.25	−1185.41
1ctf	14,346		−719.56	−1159.41	−1132.46
1gvd	14,372		−1550.61	−1883.41	−1897.59
1mjc	14,472		−678.65	−1073.45	−1058.63
1npi	14,496		−732.96	−1058.48	−1052.25
3ebx	14,556		−394.87	−642.08	−648.94
2sn3	14,574		−497.24	−804.83	−819.80
3a03	14,698		−1612.90	−1910.31	−1903.78
3a03	14,698	Residue 1	−1609.46	−1901.30	−1897.64
3a03	14,698	All atoms	−1550.64	−1835.62	−1840.62
3i8z	15,724		−748.54	−1100.40	−1108.77
3i8z	15,724	Residue 1	−709.64	−1046.47	−1052.69
3i8z	15,724	All atoms	−671.96	−1005.20	−1027.04
2c0x	16,146		−573.88	−866.13	−859.88
1ubq	16,484		−597.88	−996.73	−997.59
1ubq	16,484	Residue 1	−594.39	−991.19	−992.40
1ubq	16,484	All atoms	−538.03	−937.48	−945.91
1hyp	16,578		−505.01	−779.74	−797.72
1vcc	17,568		−659.38	−1044.55	−1063.29
1pgx	17,974		−752.87	−1137.72	−1135.05
1dsl	19,182		−591.88	−999.87	−1008.09
1tif	19,204		−1018.13	−1448.73	−1472.50
1l2p	22,294		−1099.94	−1679.26	−1664.84

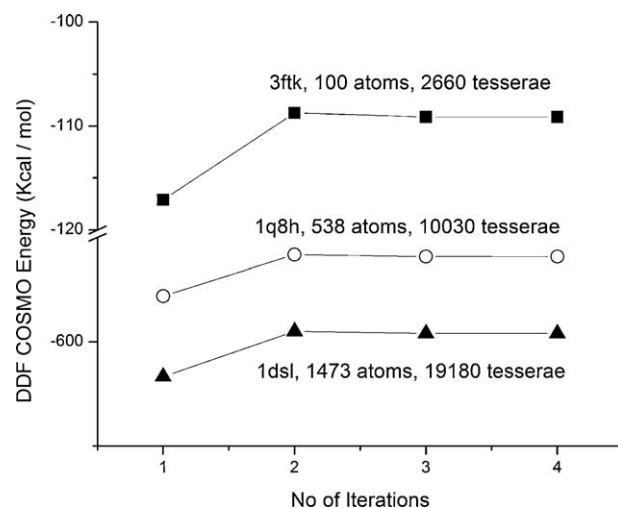
<sup>a</sup> 'Residue 1' means that the first residue in the protein sequence was considered as QM, while the rest was treated via MM. 'All atoms' means that charges are derived for the specified protein sequence via MOD-QM/MM.

<sup>a</sup> Values in parenthesis mean that only that particular residue sequence was considered as the solute.

<sup>b</sup> A blank box under "QM region" means that protein charges are taken directly from the Amber force field.



**Fig. 2.** Correlation between a full COSMO calculation (1 density domain) and the DDF-COSMO approximation using 10 density domains. (A) Uncorrected correlation. (B) Corrected correlation by Multi-linear regression.



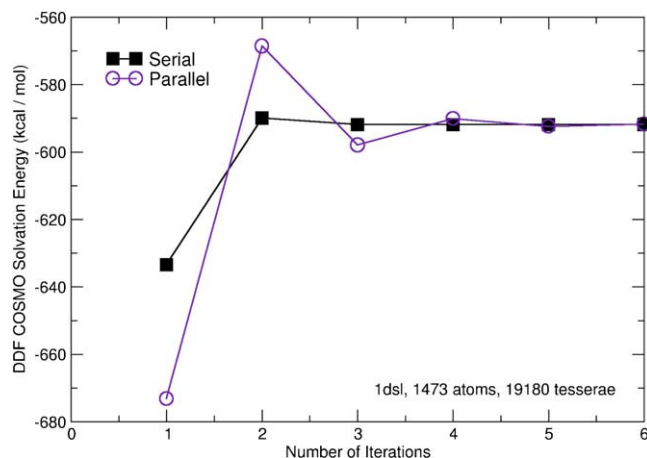
**Fig. 3.** DDF-COSMO energy as a function of the iteration cycle. Number of density domains is  $L = 10$ .

a larger or smaller ratio, a comprehensive analysis on the effect of various molecular properties on the DDF energy was carried out. The influence of total charge, total surface area, protein compactness, surface polarity was considered. Protein compactness was measured by evaluating the normalized radius of gyration [51]. Surface polarity was measured by the ratio between polar and non-polar surface area. No correlation with surface polarity and surface area was found. On the other hand, highly charged proteins ( $|q| \sim 10$ ) gave consistently high ratios (DDF-COSMO reproduces around 85% of the full COSMO energy). This observation made us look into DNA and RNA systems which normally present even larger total charge (Table 3). Indeed, it was found that as much as 95% can be recovered by DDF-COSMO, suggesting that for such systems our approach can be particularly efficient. To a lesser extent, it was found that protein compactness also correlated with this energy

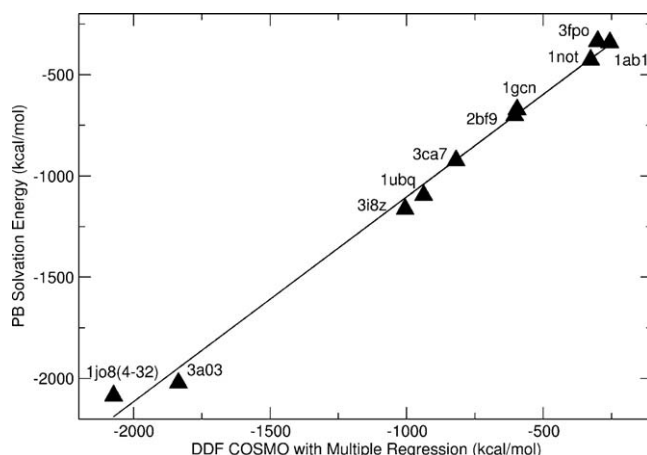
**Table 2**  
Multi linear regression parameters.

Energy term	$a$ (Intercept)	$b$ $-(1/2)(q^k)^T A^k q^k$	$c$ $(1/2)(q^0)^T A^0 q^0$	$d$ $Z^T B^0 q^0$	RMSE (kcal/mol)
Full data set <sup>a</sup>	2.4654	290.9696	4.9290	1.0383	4.15

<sup>a</sup> Multi-linear regression from all sets in Table 1. RMSE was obtained by cross validation (see supporting information for details).



**Fig. 4.** DDF-COSMO energy convergence for serial and parallel implementations. Number of density domains is  $L=10$ . Energies in the serial implementation converges in 3 cycles whereas the parallel implementation energy converges in 5 cycles.



**Fig. 5.** Correlation between PB electrostatic solvation energies and DDF-COSMO electrostatic solvation energies.

**Table 3**  
DDF-COSMO and COSMO solvation energies (in kcal/mol) for highly charged DNA/RNA systems.

System	No. of tesserae	Total charge	DDF-COSMO	COSMO	Energy ratio
1JRN	12,530	-22	-6229.1	-6530.8	0.954
3QK4	17,700	-26	-8476.5	-8936.2	0.953
3PBX	12,972	-18	-4201.2	-4516.5	0.930
3MJ3	24,930	-41	-18074.8	-18865.8	0.958

ratio (i.e., the less compact the protein is the smaller the energy recovered by DDF-COSMO).

## 5. Conclusions

We have presented an approximate method to compute the electrostatic solvation energy of a macromolecular system by integrating a conductor-like solvation model with MM, QM/MM and

MOD-QM/MM. One of the main motivations is to further expand the scope of state-of-the-art scalable conductor-like models to handle even larger systems or to achieve a large-scale calculation with restricted computational resources. We have termed this method Density Domain Fragmentation COSMO (DDF-COSMO). The main feature of the method is the fragmentation of screening charges in the density domain. Iterative solution of the electrostatic equations on each domain one at a time, assuming that they behave individually as a conductor, leads to a solvation energy that can be scaled up to one-density domain energy by a linear or multi-linear regression. We have shown that for a large test set, the DDF-COSMO does correlate linearly with the full COSMO energy. Furthermore, the DDF-COSMO algorithm can be easily parallelized saving a considerable amount of computational time. We have presented the necessary modifications of the basic equations and algorithms used in the standard COSMO approach.

## Acknowledgement

José A. Gascón acknowledges financial support from the Camille and Henry Dreyfus New Faculty Award, the NSF Career Award (CHE-0847340), the Hewlett-Packard Junior Faculty Award, and start-up package funds from University of Connecticut. The authors thank Dr. Christian Bruckner and Daniel Sandberg for helpful input and proof reading the manuscript.

## Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at [doi:10.1016/j.jmgm.2011.06.001](https://doi.org/10.1016/j.jmgm.2011.06.001).

## References

- [1] A. Jaramillo, S.J. Wodak, Computational protein design is a challenge for implicit solvation models, *Biophys. J.* 88 (2005) 156–171.
- [2] S. Wong, F. Lightstone, Accounting for water molecules in drug design, *Expert Opin. Drug Discov.* 6 (2011) 65–74.
- [3] O. Acevedo, W.L. Jorgensen, Cope elimination: elucidation of solvent effects from QM/MM simulations, *J. Am. Chem. Soc.* 128 (2006) 6141–6146.
- [4] A. Warshel, M. Levitt, Theoretical studies of enzymic reactions: dielectric, electrostatic and steric stabilization of the carbonium ion in the reaction of lysozyme, *J. Mol. Biol.* 103 (1976) 227–249.
- [5] A. Warshel, *Computer Modeling of Chemical Reactions in Enzymes and Solutions*, John Wiley & Sons, NY, 1991.
- [6] R.A. Friesner, M.D. Beachy, Quantum mechanical calculations on biological systems, *Curr. Opin. Struct. Biol.* 8 (1998) 257–262.
- [7] Z. Kurtovi, M. Marchi, D. Chandler, Umbrella sampling molecular dynamics study of the dielectric constant of water, *Mol. Phys.* 78 (1993) 1155–1165.
- [8] M. Neumann, The dielectric constant of water. Computer simulations with the MCY potential, *J. Chem. Phys.* 82 (1985) 5663–5672.
- [9] J. Tomasi, M. Persico, Molecular interactions in solution: an overview of methods based on continuous distributions of the solvent, *Chem. Rev.* 94 (1994) 2027–2094.
- [10] S. Miertus, E. Scrocco, J. Tomasi, Electrostatic interaction of a solute with a continuum – a direct utilization of ab initio molecular potentials for the prevision of solvent effects, *Chem. Phys.* 55 (1981) 117–129.
- [11] A. Klamt, Conductor-like screening model for real solvents: a new approach to the quantitative calculation of solvation phenomena, *J. Phys. Chem.* 99 (1995) 2224–2235.
- [12] A. Klamt, G. Schüürmann, COSMO: a new approach to dielectric screening in solvents with explicit expressions for the screening energy and its gradient, *J. Chem. Soc.: Perkin Trans. 2* (1993) 799–805.
- [13] P. Bandyopadhyay, M.S. Gordon, B. Mennucci, J. Tomasi, An integrated effective fragment-polarizable continuum approach to solvation: theory and application to glycine, *J. Chem. Phys.* 116 (2002) 5023–5032.

- [14] R.B. Murphy, D.M. Philipp, R.A. Friesner, A mixed quantum mechanics/molecular mechanics (QM/MM) method for large-scale modeling of chemistry in protein environments, *J. Comput. Chem.* 21 (2000) 1442–1457.
- [15] B. Marten, K. Kim, C. Cortis, R.A. Friesner, R.B. Murphy, M.N. Ringnalda, D. Sitkoff, B. Honig, New model for calculation of solvation free energies: correction of self-consistent reaction field continuum dielectric theory for short-range hydrogen-bonding effects, *J. Phys. Chem.* 100 (1996) 11775–11788.
- [16] A.R. Dinner, X. Lopez, M. Karplus, A charge-scaling method to treat solvent in QM/MM simulations, *Theor. Chem. Acc.* 109 (2003) 118–124.
- [17] C.S. Pomelli, J. Tomasi, V. Barone, An improved iterative solution to solve the electrostatic problem in the polarizable continuum model, *Theor. Chem. Acc.* 105 (2001) 446–451.
- [18] S.H. Li, T. Fang, An efficient fragment-based approach for predicting the ground-state energies and structures of large molecules, *J. Am. Chem. Soc.* 127 (2005) 7215–7226.
- [19] D.M. York, T.-S. Lee, W. Yang, Parameterization and efficient implementation of a solvent model for linear-scaling semiempirical quantum mechanical calculations of biological macromolecules, *Chem. Phys. Lett.* 263 (1996) 297–304.
- [20] T. Simonson, G. Archontis, M. Karplus, A Poisson–Boltzmann study of charge insertion in an enzyme active site: the effect of dielectric relaxation, *J. Phys. Chem. B* 103 (1999) 6142–6156.
- [21] T. Simonson, G. Archontis, M. Karplus, Continuum treatment of long-range interactions in macromolecular free energy calculations. Application to protein–ligand binding, *J. Phys. Chem. B* 101 (1997) 8349–8362.
- [22] M.S.K. Gilson, B. Honig, Calculation of the total electrostatic energy of a macromolecular system: solvation energies, binding energies, and conformational analysis, *Proteins: Struct. Funct. Genet.* 4 (1988) 7–18.
- [23] A. Nicholls, B. Honig, A rapid finite difference algorithm utilizing successive over-relaxation to solve the Poisson–Boltzmann equation, *J. Comput. Chem.* 12 (1991) 435–445.
- [24] P. Schaefer, D. Riccardi, Q. Cui, Reliable treatment of electrostatics in combined QM/MM simulation of macromolecules, *J. Chem. Phys.* 123 (2005).
- [25] W. Im, S. Berneche, B. Roux, Generalized solvent boundary potential for computer simulations, *J. Chem. Phys.* 114 (2001) 2924–2937.
- [26] S.A. Hayik, N. Liao, K.M. Merz, A combined QM/MM Poisson–Boltzmann approach, *J. Chem. Theory Comput.* 4 (2008) 1200–1207.
- [27] V. Gogonea, K.M. Merz, Fully quantum mechanical description of proteins in solution. Combining linear scaling quantum mechanical methodologies with the Poisson–Boltzmann equation, *J. Phys. Chem. A* 103 (1999) 5171–5188.
- [28] S.J. Mo, T. Vreven, B. Mennucci, K. Morokuma, J. Tomasi, Theoretical study of the S(N)2 reaction of Cl–(H<sub>2</sub>O) + CH<sub>3</sub>Cl using our own N-layered integrated molecular orbital and molecular mechanics polarizable continuum model method (ONIOM-PCM), *Theor. Chem. Acc.* 111 (2004) 154–161.
- [29] Q. Cui, Combining implicit solvation models with hybrid quantum mechanical/molecular mechanical methods: a critical test with glycine, *J. Chem. Phys.* 117 (2002) 4720–4728.
- [30] W.C. Still, A. Tempczyk, R.C. Hawley, T. Hendrickson, Semianalytical treatment of solvation for molecular mechanics and dynamics, *J. Am. Chem. Soc.* 112 (2002) 6127–6129.
- [31] E. Pellegrini, M.J. Field, A generalized-born solvation model for macromolecular hybrid-potential calculations, *J. Phys. Chem. A* 106 (2002) 1316–1326.
- [32] J.A. Gascon, S.S.F. Leung, E.R. Batista, V.S. Batista, A self-consistent space-domain decomposition method for QM/MM computations of protein electrostatic potentials, *J. Chem. Theory Comput.* 2 (2006) 175–186.
- [33] L. Menikarachchi, J. Gascón, Optimization of cutting schemes for the evaluation of molecular electrostatic potentials in proteins via Moving-Domain QM/MM, *J. Mol. Model.* 14 (2008) 479–487.
- [34] J.D. Jackson, *Classical Electrodynamics*, Wiley, NY, 1962.
- [35] H. Li, C.S. Pomelli, J.H. Jensen, Continuum solvation of large molecules described by QM/MM: a semi-iterative implementation of the PCM/EFP interface, *Theor. Chem. Acc.* 109 (2003) 71–84.
- [36] J.A. Gascon, *MODQ3M V 1.1*, University of Connecticut, Storrs, CT, 2006.
- [37] M.L. Connolly, Solvent-accessible surfaces of proteins and nucleic acids, *Science* 221 (1983) 709–713.
- [38] W.D. Cornell, P. Cieplak, C.I. Bayly, P.A. Kollman, Application of RESP charges to calculate conformational energies, hydrogen-bond energies, and free energies of solvation, *J. Am. Chem. Soc.* 115 (1993) 9620–9631.
- [39] W.D. Cornell, P. Cieplak, C.I. Bayly, I.R. Gould, K.M. Merz, D.M. Ferguson, D.C. Spellmeyer, T. Fox, J.W. Caldwell, P.A. Kollman, A 2nd generation force-field for the simulation of proteins, nucleic-acids, and organic-molecules, *J. Am. Chem. Soc.* 117 (1995) 5179.
- [40] J.W. Ponder, *Tinker V 4.2*, Washington University School of Medicine, St. Louis, Missouri, 2004.
- [41] M.J. Frisch, G.W. Trucks, H.B. Schlegel, G.E. Scuseria, M.A. Robb, J.R. Cheeseman, G. Scalmani, V. Barone, B. Mennucci, G.A. Petersson, H. Nakatsuji, M. Caricato, X. Li, H.P. Hratchian, A.F. Izmaylov, J. Bloino, G. Zheng, J.L. Sonnenberg, M. Hada, M. Ehara, K. Toyota, R. Fukuda, J. Hasegawa, M. Ishida, T. Nakajima, Y. Honda, O. Kitao, H. Nakai, T. Vreven, J.A. Montgomery Jr., J.E. Peralta, F. Ogliaro, M. Bearpark, J.J. Heyd, E. Brothers, K.N. Kudin, V.N. Staroverov, R. Kobayashi, J. Normand, K. Raghavachari, A. Rendell, J.C. Burant, S.S. Iyengar, J. Tomasi, M. Cossi, N.J. Millam, M. Klene, J.E. Knox, J.B. Cross, V. Bakken, C. Adamo, J. Jaramillo, R.E. Gomperts, O. Stratmann, A.J. Yazyev, R. Austin, C. Cammi, J.W. Pomelli, R. Ochterski, R.L. Martin, K. Morokuma, V.G. Zakrzewski, G.A. Voth, P. Salvador, J.J. Dannenberg, S. Dapprich, A.D. Daniels, O. Farkas, J.B. Foresman, J.V. Ortiz, J. Cioslowski, D.J. Fox, *Gaussian 09 Rev. A.1*, Gaussian Inc., Wallingford, CT, 2009.
- [42] S. Dapprich, I. Komaromi, K.S. Byun, K. Morokuma, M.J. Frisch, A new ONIOM implementation in Gaussian98. Part I. The calculation of energies, gradients, vibrational frequencies and electric field derivatives, *J. Mol. Struct.-Theochem.* 462 (1999) 1–21.
- [43] M. Maseras, K. Morokuma, IMOMM – A new integrated ab-initio plus molecular mechanics geometry optimization scheme of equilibrium structures and transition-states, *J. Comput. Chem.* 16 (1995) 1170–1179.
- [44] M. Svensson, S. Humbel, R.D.J. Froese, T. Matsubara, S. Sieber, K. Morokuma, ONIOM: A multilayered integrated MO + MM method for geometry optimizations and single point energy predictions. A test for Diels-Alder reactions and Pt(P(t-Bu)(3))(2) + H-2 oxidative addition, *J. Phys. Chem.* 100 (1996) 19357–19363.
- [45] A. Klamt, V. Jonas, T. Burger, J.C.W. Lohrenz, Refinement and parametrization of COSMO-RS, *J. Phys. Chem. A* 102 (1998) 5074–5085.
- [46] Y. Saad, *Sparskit: A Basic Tool Kit for Sparse Matrix Computations*, Research Institute for Advanced Computer Science, NASA Ames Research Center, Moffett Field, CA, 1990.
- [47] M.F. Sanner, A.J. Olson, J.C. Spehner, Reduced surface: an efficient way to compute molecular surfaces, *Biopolymers* 38 (1996) 305–320.
- [48] J.L. Pascual-ahuir, E. Silla, I. Tuñón, GEPOL: an improved description of molecular surfaces. III. A new algorithm for the computation of a solvent-excluding surface, *J. Comput. Chem.* 15 (1994) 1127–1138.
- [49] N.A. Baker, D. Sept, S. Joseph, M.J. Holst, J.A. McCammon, Electrostatics of nanosystems: application to microtubules and the ribosome, *Proc. Natl. Acad. Sci. U.S.A.* 98 (2001) 10037–10041.
- [50] A. Khandelwal, S. Balaz, QM/MM linear response method distinguishes ligand affinities for closely related metalloproteins, *Proteins: Struct. Funct. Bioinf.* 69 (2007) 326–339.
- [51] M. Lobanov, N. Bogatyreva, O. Galzitskaya, Radius of gyration as an indicator of protein structure compactness, *Mol. Biol.* 42 (2008) 623–628.