

Testing the semi-explicit assembly solvation model in the SAMPL3 community blind test

Charles W. Kehoe · Christopher J. Fennell ·
Ken A. Dill

Received: 24 October 2011 / Accepted: 14 December 2011 / Published online: 29 December 2011
© Springer Science+Business Media B.V. 2011

Abstract We report here a test of the Semi-Explicit Assembly (SEA) model in the solvation free energy category of the SAMPL3 blind prediction event (summer 2011). We tested how dependent the SEA results are on the chosen force field by performing calculations with both the General Amber and OPLS force fields. We compared our SEA results with full molecular dynamics simulations in explicit solvent. Of the 20 submissions, our SEA/OPLS results gave the second smallest RMS errors in free energies compared to experiments. SEA gives results that are very similar to those of its underlying force field and explicit solvent model. Hence, while the SEA water modeling approach is much faster than explicit solvent simulations, its predictions appear to be just as accurate.

Keywords Solvation · Implicit · Explicit · Semi-explicit · Force Field

Introduction

Motivated by the need for improved computational models of water and aqueous solutions, we recently developed a solvation model called Semi-Explicit Assembly (SEA) water [1, 2].

Electronic supplementary material The online version of this article (doi:10.1007/s10822-011-9536-8) contains supplementary material, which is available to authorized users.

C. W. Kehoe
Graduate Group in Bioinformatics, University of California
at San Francisco, San Francisco, CA 94158, USA

C. J. Fennell · K. A. Dill (✉)
Laufer Center for Physical and Quantitative Biology, Stony
Brook University, Stony Brook, NY 11794, USA
e-mail: dill@laufercenter.org

In this paper, we report a test of the SEA model through our participation in the solvation-free-energy category of the SAMPL3 blind prediction challenge.

Description of the SEA-water method

SEA has been described in detail elsewhere [1, 2], so here we give only a brief summary. As shown in Fig. 1, SEA divides the calculation of the solvation free energy of a solute into two parts. First, there is a pre-computation stage which samples the interaction of various test spheres (having different radii, charges and van der Waals properties) with a chosen explicit-solvent model of water, such as TIP3P. Among other properties, these pre-simulations provide the average axial dipole moment of first-shell waters as a function of the local electric field around the different solute spheres. Second, at runtime, SEA models the solute under study as an assembly of those pre-computed spheres. We can then rapidly compute a solvation free energy as a sum over the properties of all the water molecules around the solute. In this way, the SEA model captures many of the physical and structural properties of the explicit-water model on which it was parameterized, and also captures the particulate and electrostatic properties of first-shell waters. Yet the calculations at runtime are nearly as fast to compute as those of implicit-solvent models.

SAMPL3 small molecules: a test of chlorinated molecule solvation

The SAMPL event, developed and run by OpenEye Software, is a community-wide blind test of computational chemistry prediction methods, including the prediction of solvation free energies [3–5]. At the start, participants are provided with a set of varied solutes for which they do not

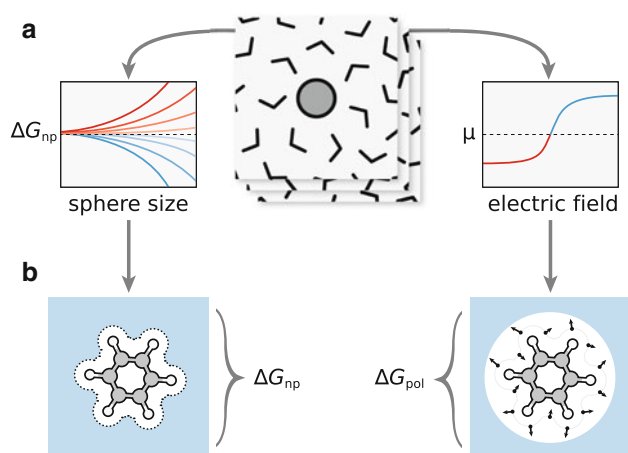


Fig. 1 **a** The SEA solvation model involves pre-simulations of neutral and charged *spheres* in explicit water to generate maps of the ΔG_{np} versus sphere size and solvent dipole response versus surface electric field. **b** After these one-time calculations, we can rapidly calculate the ΔG_{np} and ΔG_{pol} terms of the total solvation free energy for an arbitrary solute molecule

know the solvation free energies. Each group then uses its particular methodology to compute solvation free energies for these compounds. After these predictions are submitted, the SAMPL organizers then provide experimental solvation free energies to compare against. This year's event, SAMPL3, included 36 different solute molecules, roughly evenly divided among ethanes, biphenyls, and dibenzo-*p*-dioxins. The molecules in each class varied only in the number and location of substituted chlorines. This systematic approach allowed for an analysis of trends across these classes of molecules, but it also had the potential to exaggerate errors for any methods that are challenged by these particular molecule types or by chlorinated solutes. In this work, we assess our predictions for these 36 molecules, and explore what this and other SAMPL datasets tell us about SEA, explicit, and implicit solvation models.

Simulation methods

The SEA method relies on underlying pre-simulations using an atomically detailed force field. For comparisons of computational solvation methods, we performed SEA and explicit-solvent calculations with the General Amber Force Field (GAFF) and the OPLS force field [6, 7]. We used AM1-BCC partial charges [8] for the GAFF calculations. As an additional implicit solvent comparison, we performed Poisson–Boltzmann (PB) calculations with the resultant GAFF solute topologies. Our GAFF parameters and AM1-BCC charges came from the ANTECHAMBER program in AmberTools 1.4 [9], and we used the TIP3P water model [10] to solvate these structures, following the

practices of Mobley et al [11]. OPLS parameters were from Desmond 2.4.2.1 [12], and we solvated these structures using the SPC water model [13], as was done in a published study of OPLS small-molecule solvation free energy calculations [14]. For each force field, we also submitted two sets of results: one set with minimally relaxed, single conformations, and another with multiple conformations per molecule, sampled from explicit water simulations.

We submitted the SEA and PB calculations to the SAMPL3 event as blind predictions, and performed the explicit calculations for the present study. We also compared results from solutes put into a single (dominant) configuration against results that sampled some of the solute's conformational options.

Setting up molecules with Amber-based and OPLS-based force fields

We performed all molecular dynamics (MD) calculations with GROMACS 4.0.4 [15, 16]. We switched off Lennard–Jones interactions between 8 and 9 Å and applied long-ranged energy and pressure corrections. We used smooth particle-mesh Ewald for long-ranged electrostatics accumulation [17], this with a real-space cutoff of 10 Å, a spline order of 6, fourier spacing of 1 Å, and a real-space energy tolerance parameter of 10^{-6} kJ/mol. We used the SETTLE algorithm to constrain the geometry of water molecules in explicit-solvent simulations [18], and the LINCS algorithm to constrain covalent bonds to hydrogen atoms on solute molecules [19].

We solvated the target solutes in a rhombic dodecahedral box with a 12 Å buffer of explicit water between the atoms of the solute and the edges of the box. We relaxed these systems in two steps, first with 1,000 steps of steepest descent minimization followed by 10 ps of constant energy MD with a 1 fs timestep. All timesteps after this relaxation phase were 2 fs. For the single-conformation calculations with SEA, we used the final solute structure following 100 ps of Langevin Dynamics (LD) equilibration of these relaxed structures in explicit solvent. For the multi-conformation sets, we followed relaxation with 100 ps of constant temperature LD (300 K), 100 ps of constant pressure (1 atm) equilibration using the Berendsen thermostat, and 700 ps of constant pressure dynamics using the Parrinello–Rahman barostat. We rescaled the simulation box to the average volume from last 500 ps of this constant pressure trajectory, and equilibrated this system with 500 ps of constant volume simulation before our 1 ns LD production run. We selected 10 conformations at equal time intervals from the production run calculations with SEA. Our final multi-conformation results are arithmetic averages of the SEA calculations on these 10 solute conformers.

SEA solvation free energy calculations

The SEA solvation method uses an extensive set of pre-simulations that were calculated using Lorentz–Berthelot mixing rules for Lennard–Jones (LJ) interactions. These resulting property tables are compatible with Amber force fields, but not directly compatible with OPLS force fields because OPLS uses a geometric mean for determining the LJ size (σ) parameters between dissimilar atom types. To perform OPLS force field calculations with SEA, we converted the OPLS LJ parameters to an arithmetic mean equivalent,

$$\sigma_{\text{new}} = 2(\sqrt{\sigma_{\text{OPLS}} \cdot \sigma_{\text{wat}}}) - \sigma_{\text{wat}} \quad (1)$$

where σ_{wat} is the LJ σ value for the explicit water model of interest, in this case SPC water [13]. Note that converting atom parameters in this way will result in differences in solute intramolecular LJ interactions. This does not directly affect the SEA calculations as there is no solute intramolecular interaction contribution in its estimation of the solvation free energy.

PB solvation free energy calculations

We performed our PB tests by combining a polar term and a non-polar term. For the polar term, we solved the linearized Poisson equation provided in APBS [20]. We set the dielectric boundary as the molecular surface and used a solvent dielectric constant of 78. For the non-polar term, we used the standard expression

$$\Delta G_{\text{np}} = 0.00542 \times \text{SASA} + 0.92, \quad (2)$$

where SASA is the solvent-accessible surface area in \AA^2 to give ΔG_{np} values in kcal/mol [21]. Our PB calculations were performed on the multi-conformer GAFF structures.

Explicit solvent solvation free energy calculations

For our explicit solvent comparisons, we performed solvation free energy calculations as others have done in previous studies [22, 23]. Here, we used thermodynamic integration to transform the solute molecules between the relevant state-points. For the ΔG_{chg} term, we turned off the partial charges of the solute in TIP3P or SPC water (for the GAFF and OPLS calculations respectively) over λ windows of {0, 0.2, 0.4, 0.6, 0.8, 1.0}. Each of these windows involved a 5 ns simulation with GROMACS 4.0.4 using the same simulation protocols described previously for the conformation sampling calculations. In order to determine the absolute ΔG_{pol} term for transfer from air to water, we subtracted out the internal Coulombic and conformation distribution contributions to ΔG_{chg} . We calculated these by determining the ΔG_{chg} term in vacuum, using λ window

values of {0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0}. For the ΔG_{np} term calculation, we converted the uncharged solute to a soft-core representation[24] over the set of λ window values {0, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95, 1.0}. We used the standard trapezoid rule for integration, and accumulated errors via the limiting value of block averages [25].

Results and discussion

Figure 2 shows an overview of our SAMPL3 submissions. For each method, we show the bootstrapped RMS difference in free energy between the predicted solvation free energy and the experimental value (as reported by the SAMPL organizers). After we parameterized the molecules and sampled their conformations, SEA generated each prediction in <1 s on a 3 GHz Intel Core 2 E8400 processor. As shown, SEA was one of the more accurate predictors in the event, despite its modest computational cost. In terms of absolute performance relative to the rest of the predictions, it should be noted that this is a small set of very related solutes, so it is difficult to make definitive assessments. We perform additional comparisons below to explore these SEA results in more detail.

As an implicit solvent reference, we submitted PB results using the GAFF and TIP3P LJ parameters to define the molecular surface and solvent accessible surface area along with the same AM1-BCC partial charges used in the SEA calculations. This combination of parameters did not perform particularly well, about 1.7 kcal/mol higher RMSE than the analogous SEA submission. Most of this difference comes from the use of Eq. 2 as opposed to a more microscopically physical approach to ΔG_{np} . If the SEA

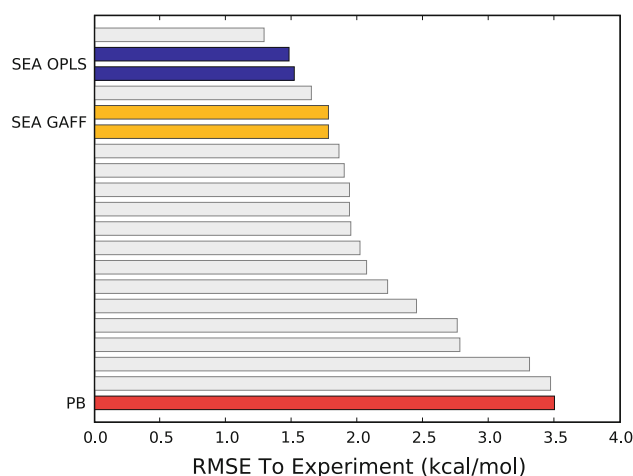


Fig. 2 Overall performance of the 20 submissions to SAMPL3. Bootstrapped RMSE of each method; best predictions are at the *top*. Our submissions are *highlighted*

ΔG_{np} term is used in place of Eq. 2, this application of PB would have given an only 0.2 kcal/mol higher RMSE than SEA.

Multiple conformations do not alter the solvation free energies

The 36 compounds in the SAMPL3 test were all fairly rigid. Not surprisingly, we found that the free energies of the conformational ensembles of the solutes are accurately approximated, in this case, by the single dominant conformer. The differences are only a few hundredths of a kcal/mol, similar to the deviations between runs on identical input. This is markedly different than for previous SAMPL events, in which some solute molecules (like the sugars in SAMPL2) involved multiple energetically important conformations.

After the SAMPL meeting, it came to our attention that three of the biphenyl molecules, 1,2,4,5-tetrachloro-3-(3,4-dichlorophenyl) benzene, 1,2,3,4-tetrachloro-5-(3,4-dichlorophenyl) benzene, and 1,2,3,4-tetrachloro-5-(3,4,5-trichlorophenyl) benzene, have a large dihedral barrier between two stereochemically unique conformations that might affect the resulting solvation free energies. While the multi-conformation calculations failed to show any added uncertainty for these molecules, no barrier crossings were identified over the course of a 10 ns MD conformation calculation.

To test if these unsampled conformations would play any role in the resulting solvation free energies, we performed multi-conformer SEA calculations about relaxed structures on either side of this dihedral barrier. Table S1 shows the results of these calculations. The differences between these conformations, though not completely negligible, are in most cases much less than the difference between the force field's predictions and the experimental results.

The SEA model is an accurate mimic of the underlying explicit water model

Figure 3 shows a comparison of SEA results with experimental numbers, and the same comparison with explicit solvent free energy calculations. The plots show how the SEA model captures quite accurately the much more expensive explicit force field simulations on which it was based. That is, SEA/GAFF gives very similar results to GAFF explicit solvent calculations, and SEA/OPLS gives very similar results to OPLS explicit solvent calculations. In fact, in this limited data set, the SEA results happen to be more faithful to the experimental numbers than do the explicit solvent results. The SEA methods are also six orders of magnitude faster to compute.

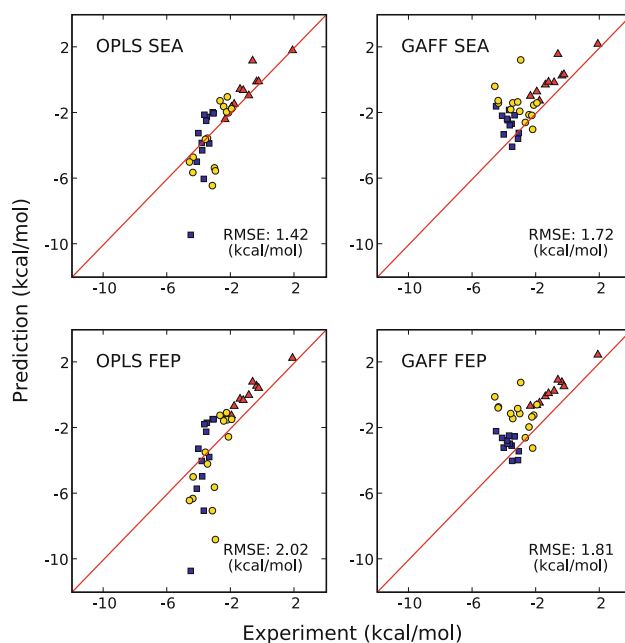


Fig. 3 The performance of SEA/GAFF, SEA/OPLS, GAFF explicit and OPLS explicit versus experimental data for all the solutes in SAMPL3. Ethanes appear as *triangles*, biphenyls as *circles*, and dioxins as *squares*

SEA inherits not only the strengths, but also weaknesses of the underlying force fields fed into it. For both force fields we tested, the most heavily chlorinated biphenyl and dioxin solutes were the most problematic, and these problems were reflected equally in both the SEA and explicit simulations. With the GAFF topologies, SEA and explicit solvent systematically *under-solvated* these molecules. With the OPLS topologies, SEA and explicit solvent systematically *over-solvated* these molecules. The reason for these errors, at least in the dioxins, is shown in Fig. 4. While the LJ parameters are similar for the aromatic carbon and substituted chlorine atoms, the partial charges are roughly 7 times greater with OPLS than with GAFF/AM1-BCC. The resulting difference in ring electric fields is why we see an 8 kcal/mol difference between GAFF in OPLS in both the SEA and explicit calculations. It is unclear which set of partial charges gives a better representation of the actual electric field, since the GAFF calculations are closer to experiment in this case, but not in others. However, SEA and explicit solvent respond to the topologies in the same way.

We also performed retrospective tests on past SAMPL solutes using SEA

As a further test, we performed SEA calculations on all the previous SAMPL solute sets; our results are shown in Fig. 5. We used the same simulation protocol as that used for the SAMPL3 calculations. We compared the multi-

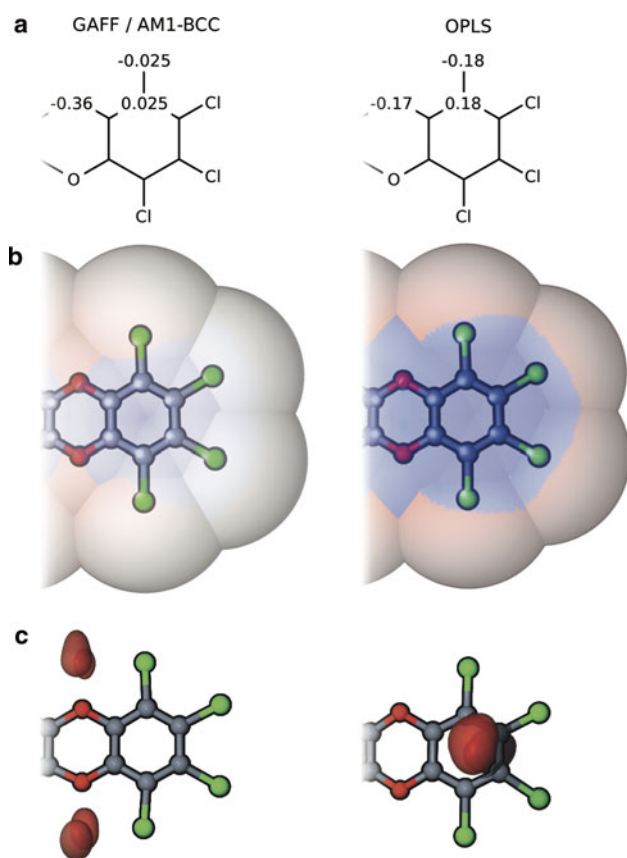


Fig. 4 A comparison of the aromatic ring region of octachlorodibenzo-*p*-dioxin with the GAFF/AM1-BCC and OPLS topologies. The **a** partial charges are larger in OPLS, resulting in a **b** stronger field seen by SEA at the solvent-accessible surface. This difference in electric fields also affects explicit solvent, with **c** a greater water occupancy probability interacting with the ring with the OPLS partial charges

conformer results because previous SAMPL events had some solutes with more flexible dihedrals.

OPLS performed poorly in SAMPL2, while GAFF performed poorly in SAMPL1. In SAMPL2, OPLS performed poorly on cyclic nitrogen compounds. When these compounds are removed from the dataset, OPLS performance improves markedly, while GAFF performance degrades slightly. Similarly in SAMPL1, GAFF (and OPLS) performs poorly on sulfur compounds. Removing the sulfur compounds reduces the problem. Results from both forcefields also improve when we skip the amides from SAMPL0.

These results illustrate a key challenge of force field development: parameters that work better for some molecules may be worse for others. SEA provides a tool for diagnosing these problems. It is notable that for most SAMPLs, one of either SEA/GAFF or SEA/OPLS gives predictions that are among the best predictions of that year.

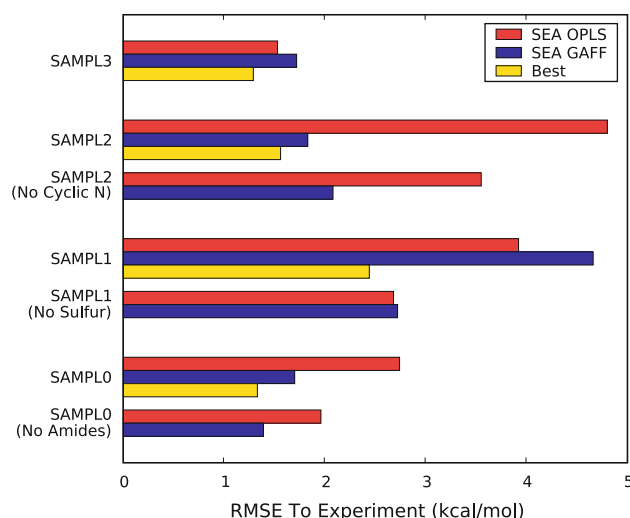


Fig. 5 The performance of SEA/GAFF and SEA/OPLS on all SAMPL events. Which force field performs better depends heavily on the set of molecules in question. In most cases, however, SEA is able to produce results rivaling the *top* performer from that year, given a sufficiently accurate force field. We are also able to identify some force field weaknesses: both forcefields have trouble with the sulfur compounds in SAMPL1, and the amides in SAMPL0. OPLS also improves markedly if we skip the cyclic nitrogen compounds from SAMPL2. SEA's accuracy and speed can greatly facilitate force field tuning in cases like these

Conclusions

We have tested the SEA water method of computational solvation in the blind test SAMPL3 event on 36 solute molecules. We submitted SEA calculation results using both GAFF and the OPLS force field. SEA/OPLS performed second best among 20 submissions. For these rigid solutes, we confirmed that sampling solute conformations added no further value. We compared our SEA approach with explicit solvent free-energy calculations, and found that both methods give quantitatively similar solvation free energies when based on the same underlying force field. The advantage of the SEA method over explicit solvation simulations is that it is orders of magnitude faster to compute. Hence SEA water promises to be faster than explicit solvation and more accurate than implicit solvation.

Acknowledgments The authors thank Professor David Mobley for helpful discussions. The authors appreciate the support from NIH Grant GM063592.

References

- Fennell CJ, Kehoe C, Dill KA (2010) Oil/water transfer is partly driven by molecular shape, not just size. *J Am Chem Soc* 132:234–240

2. Fennell CJ, Kehoe CW, Dill KA (2011) Modeling aqueous solvation with semi-explicit assembly. *Proc Natl Acad Sci USA* 108:3234–3239
3. Nicholls A, Mobley DL, Guthrie JP, Chodera JD, Bayly CI, Cooper MD, Pande VS (2008) Predicting small-molecule solvation free energies: an informal blind test for computational chemistry. *J Med Chem* 51:769–779
4. Guthrie JP (2009) A blind challenge for computational solvation free energies: introduction and overview. *J Phys Chem B* 113:4501–4507
5. Geballe MT, Skillman AG, Nicholls A, Guthrie JP, Taylor PJ (2010) The SAMPL2 blind prediction challenge: introduction and overview. *J Comput Aided Mol Des* 24:259–279
6. Wang J, Wolf RM, Caldwell JW, Kollman PA, Case DA (2004) Development and testing of a general amber force field. *J Comput Chem* 25:1157–1174
7. Jorgensen WL, Maxwell DS, Tirado-Rives J (1996) Development and testing of the opls all-atom force field on conformational energetics and properties of organic liquids. *J Am Chem Soc* 118:11225–11236
8. Jakalian A, Bush BL, Jack DB, Bayly CI (2000) Fast, efficient generation of high-quality atomic charges. AM1-BCC model: I. method. *J Comput Chem* 21(2):132–146
9. Wang J, Wang W, Kollman P, Case D (2006) Automatic atom type and bond type perception in molecular mechanical calculations. *J Mol Graph Model* 25:247–260
10. Jorgensen WL, Chandrasekhar J, Madura JD, Impey RW, Klein ML (1983) Comparison of simple potential functions for simulating liquid water. *J Chem Phys* 79:926–935
11. Mobley DL, Liu S, Cerutti DS, Swope WC, Rice JE (2011) Alchemical prediction of hydration free energies for SAMPL. *J Comput Aided Mol Des*. doi:[10.1007/s10822-011-9528-8](https://doi.org/10.1007/s10822-011-9528-8)
12. Bowers KJ et al (2006) Scalable algorithms for molecular dynamics simulations on commodity clusters. In: *Proceedings of the 2006 ACM/IEEE conference on supercomputing*, New York, SC '06, ACM
13. Berendsen HJC, Postma JPM, van Gunsteren WF, Hermans J (1981) Simple point charge water. In: Pullmann B (ed) *Inter-molecular forces*. Reidel, Dordrecht pp 331–342
14. Shivakumar D, Williams J, Wu Y, Damm W, Shelley J, Sherman W (2010) Prediction of absolute solvation free energies using molecular dynamics free energy perturbation and the opls force field. *J Chem Theory Comput* 6:1509–1519
15. Berendsen HJC, van der Spoel D, van Drunen R (1995) GRO-MACS: a message-passing parallel molecular dynamics implementation. *Comput Phys Comm* 91:43–56
16. Hess B, Kutzner C, van der Spoel D, Lindahl E (2008) GROMACS 4: algorithms for highly efficient, load-balanced, and scalable molecular simulation. *J Chem Theory Comput* 4:435–447
17. Essman U, Perela L, Berkowitz ML, Darden T, Lee H, Pedersen LG (1995) A smooth particle mesh ewald method. *J Chem Phys* 103:8577–8592
18. Miyamoto S, Kollman PA (1992) SETTLE: an analytical version of the SHAKE and RATTLE algorithms for rigid water models. *J Comput Chem* 13:952–962
19. Hess B, Bekker H, Berendsen HJC, Fraaije JGEM (1997) LINCS: a linear constraint solver for molecular simulations. *J Comput Chem* 18:1463–1472
20. Baker NA, Sept D, Joseph S, Holst MJ, McCammon JA (2001) Electrostatics of nanosystems: application to microtubules and the ribosome. *Proc Natl Acad Sci USA* 98:10037–10041
21. Rizzo RC, Aynechi T, Case DA, Kuntz ID (2006) Estimation of absolute free energies of hydration using continuum methods: accuracy of partial charge models and optimization of nonpolar contributions. *J Chem Theory Comput* 2:128–139
22. Mobley DL, Dumont E, Chodera JD, Dill KA (2007) Comparison of charge models for fixed-charge force fields: small-molecule hydration free energies in explicit solvent. *J Phys Chem B* 111: 2242–2254
23. Mobley DL, Bayly CI, Cooper MD, Shirts MR, Dill KA (2009) Small molecule hydration free energies in explicit solvent: an extensive test of fixed-charge atomistic simulations. *J Chem Theory Comput* 5:350–358
24. Steinbrecher T, Mobley DL, Case DA (2007) Nonlinear scaling schemes for Lennard–Jones interactions in free energy calculations. *J Chem Phys* 127:214108
25. Hess B (2002) Determining the shear viscosity of model liquids from molecular dynamics simulations. *J Chem Phys* 116:209–217