



Cite this: *Phys. Chem. Chem. Phys.*, 2015, 17, 12441

# Predicting accurate absolute binding energies in aqueous solution: thermodynamic considerations for electronic structure methods

Jan H. Jensen

Recent predictions of absolute binding free energies of host–guest complexes in aqueous solution using electronic structure theory have been encouraging for some systems, while other systems remain problematic. In this paper I summarize some of the many factors that could easily contribute 1–3 kcal mol<sup>-1</sup> errors at 298 K: three-body dispersion effects, molecular symmetry, anharmonicity, spurious imaginary frequencies, insufficient conformational sampling, wrong or changing ionization states, errors in the solvation free energy of ions, and explicit solvent (and ion) effects that are not well-represented by continuum models. While I focus on binding free energies in aqueous solution the approach also applies (with minor adjustments) to any free energy difference such as conformational or reaction free energy differences or activation free energies in any solvent.

Received 31st January 2015,  
Accepted 7th April 2015

DOI: 10.1039/c5cp00628g

www.rsc.org/pccp

## Introduction

The prediction of accurate absolute binding energies in aqueous solution is one of the holy grails of computational chemistry,

Department of Chemistry, University of Copenhagen, Universitetsparken 5, 2100 Copenhagen, Denmark. E-mail: jhjensen@chem.ku.dk; Web: www.twitter.com/janhjensen



Jan H. Jensen

Jan H. Jensen obtained his PhD in theoretical chemistry in 1995 from Iowa State University working with Mark Gordon, where he continued as a postdoc until he joined the faculty at the University of Iowa in 1997. In 2006 he moved to the University of Copenhagen, where he is now professor of bio-computational chemistry. In addition to his research interests in quantum biochemistry, he is interested in blended learning (receiving the

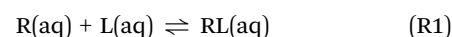
University of Copenhagen teacher of the year award in 2013), open access publishing and open science. He is an active blogger (on *Molecular Modeling Basics and Proteins and Wave Functions*), active on Twitter and Google+ and initiated the overlay journal *Computational Chemistry Highlights* and the aggregator *Computational Chemistry Daily*. He welcomes further discussion on this Perspective on PubPeer.

mainly because of the potential use in rational drug design. “Accurate” is typically taken to be 1 kcal mol<sup>-1</sup>, which roughly corresponds to predicting the binding constant within an order of magnitude at room temperature and it is understood that the method must be generally applicable. The recent blind prediction challenge SAMPL4 has shown that this goal has yet to be met even for host–guest complexes that are significantly smaller than proteins.<sup>1</sup> Interestingly, the entry that arguably performed best for one of the hosts (curcurbit[7]uril or CB7) was, for the first time, based on the rigid rotor-harmonic oscillator (RRHO) approximation and electronic structure theory and involved no direct parameterization against experimental binding free energies.<sup>2</sup> This method reproduced 14 experimental CB7–guest binding free energies with a mean absolute deviation of 2.02 ± 0.46 kcal mol<sup>-1</sup> suggesting that, perhaps, the holy grail is within reach. However, the mean absolute error was significantly larger for another host–guest system indicating that there remains some work to be done.

In this paper I summarize why electronic structure/RRHO-based approaches are starting to yield accurate binding free energies. I also discuss many of the possible sources of error when computing aqueous binding free energies with electronic structure theory and how to correct for them.

## General approach

The general approach for predicting the standard free energy of binding ( $\Delta G_{b,aq}^\circ$ ) of a receptor (R or host) and ligand (L or guest) molecule in aqueous (aq) solution



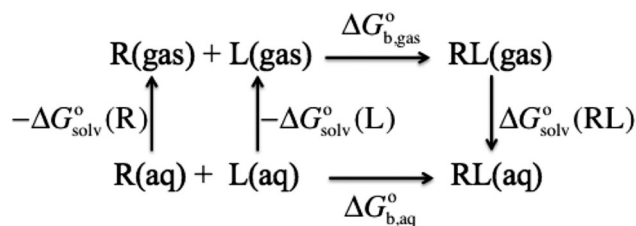


Fig. 1 Thermodynamic cycle for computing the binding free energy in aqueous solution for a ligand (L) binding to a receptor (R) to form a complex (RL).

using electronic structure theories is through a thermodynamic cycle (Fig. 1)

$$\Delta G_{b, aq}^{\circ} = G_{aq}^{\circ}(\text{RL}) - G_{aq}^{\circ}(\text{R}) - G_{aq}^{\circ}(\text{L}) \quad (1)$$

where

$$\begin{aligned} G_{aq}^{\circ}(\text{X}) &= E_{\text{gas}}(\text{X}) + G_{\text{gas, RRHO}}^{\circ}(\text{X}) + \Delta G_{\text{solv}}^{\circ}(\text{X}) \\ &= G_{\text{gas}}^{\circ}(\text{X}) + \Delta G_{\text{solv}}^{\circ}(\text{X}) \end{aligned} \quad (2)$$

$E_{\text{gas}}(\text{X})$ ,  $G_{\text{gas, RRHO}}^{\circ}(\text{X})$ , and  $\Delta G_{\text{solv}}^{\circ}(\text{X})$  is the electronic energy, rigid rotor-harmonic oscillator (RRHO), and solvation free energy, respectively, of molecule X. Note that  $G_{\text{gas, RRHO}}^{\circ}(\text{X})$  contains the zero point energy. The standard state (denoted by “ $\circ$ ”) throughout this paper is 1 M, unless otherwise noted. The solvation free energy is typically computed using a continuum solvation model as described in detail below.

## The electronic energy

One of the reasons electronic structure-based approaches are starting to yield accurate binding free energies is the use of dispersion corrections<sup>3</sup> in the evaluation of the electronic energy and the structure (as well as the vibrational frequencies as discussed below). Grimme<sup>4</sup> has shown that dispersion typically makes a very big ( $>10 \text{ kcal mol}^{-1}$ ) contribution to binding free energies of host-guest complexes. Dispersion corrections are therefore a must if DFT is used to compute the electronic binding energy. Furthermore, Grimme has shown that three-body dispersion makes a non-negligible ( $2\text{--}3 \text{ kcal mol}^{-1}$ ) contribution to the electronic binding energy. For convergent methods this effect is only included in rather expensive methods that involve triple-excitations such as MP4 and CCSD(T).

Interestingly, it has been found that dispersion corrected, and short-range corrected, semiempirical methods such as DFTB or PM6, yield binding energies with accuracies similar to conventional DFT calculations with large basis sets. For example, Muddana and Gilson<sup>5</sup> used PM6-DH+ to compute reasonably accurate relative binding energies for CB7-ligand complexes. On the other hand, Yilmazer and Korth<sup>6</sup> found significant deviations for PM6-DH+ and similar methods when applied to larger protein-ligand models. Whether these minimal basis set-based methods are sufficiently flexible to handle large many-body polarization effects involving many charged groups remains to be determined. In any case, Grimme and co-workers

have computed  $E_{\text{gas}}(\text{X})$  at the PW6B95-D3(BJ)/def2-QZVP//TPSS27-D3(BJ)/def2-TZVP level of theory with good results.<sup>2</sup>

## Molecular thermodynamics

The translational, rotational and vibrational thermodynamic contribution to the binding free energy is very large ( $>10 \text{ kcal mol}^{-1}$ ) and must be included for accurate results. Some years ago there was a bit of confusion in the literature about whether the RRHO approach was appropriate for condensed phase systems, but Zhou and Gilson<sup>7</sup> have clarified this beautifully. The accuracy of the dispersion and hydrogen bond-corrected semi-empirical methods mentioned above has now made it feasible to compute the vibrational frequencies for typical host-guest complexes and this is another reason why electronic structure-based approaches are starting to yield accurate binding free energies. (They appear to be a qualitative step forward in accuracy compared to standard force fields in this regard.) For example, Grimme has computed  $G_{\text{gas, RRHO}}^{\circ}(\text{X})$  with PM6-D3H<sup>4</sup> and HF-3c<sup>2</sup> with good results.

### The standard state

Most electronic structure codes compute the RRHO energy corrections for an ideal gas, where the standard state is a pressure of 1 bar. As I'll discuss further below the solvation free energies are computed for a 1 M standard state so the gas phase free energy must be corrected accordingly

$$G_{\text{gas, RRHO}}^{\circ}(\text{X}) = G_{\text{gas, RRHO}}^{\circ}(1 \text{ bar})(\text{X}) - RT \ln(V^{-1}) \quad (3)$$

where  $V$  is the volume of an ideal gas at a temperature  $T$  and  $R$  is the ideal gas constant. At 298 K this correction increases the free energy by  $1.90 \text{ kcal mol}^{-1}$ .

It is tempting to argue that since the volume change in solution is negligible one should use the Helmholtz free energy  $A_{\text{gas, RRHO}}^{\circ}(\text{X})$  instead of the Gibbs free energy. However, as I discuss below, the solvation free energy corrects for the change in volume on going from the gas phase to solution, so the Gibbs free energy change should be used throughout.

### The vibrational enthalpy for NDDO based semiempirical methods

NDDO based semiempirical methods such as PM6 are parameterized against experimental standard enthalpies of formation ( $\Delta H_{f, \text{gas}}^{\circ}$ ). However, in the case of intermolecular interactions such as hydrogen binding the parameterization was done by fitting  $\Delta \Delta H_{f, \text{gas}}^{\circ}$  to  $\Delta E_{\text{gas}}$  values computed using electronic structure theory (Stewart 2007).<sup>8</sup> The same is true for dispersion and hydrogen bond corrected PM6 methods. Thus, if a PM6 based method is used to compute the interaction energy the RRHO enthalpy corrections should still be included, *i.e.*

$$G_{aq}^{\circ}(\text{X}) = \Delta H_{f, \text{gas}}^{\circ}(\text{X}) + G_{\text{gas, RRHO}}^{\circ}(\text{X}) + \Delta G_{\text{solv}}^{\circ}(\text{X}) \quad (4)$$

### Molecular symmetry

Many host molecules and some guest molecules are symmetric and this affects the rigid-rotor rotational entropy ( $S_{\text{RR}}$ ) through



the symmetry number ( $\sigma$ ), which is a function of the molecular point group.

$$S_{\text{RR}} = R \ln \left( \frac{8\pi^2}{\sigma} \left( \frac{2\pi ekT}{h^2} \right)^{3/2} \sqrt{I_1 I_2 I_3} \right) \quad (5)$$

Here  $h$  and  $k$  are Planck's and Boltzmann's constant, respectively and  $I_x$  is the moment of inertia for principal axis  $x$ . In practice it can be very difficult to build large molecules with the correct point group and most studies use  $C_1$  symmetry. In this case the effect of symmetry must be added manually to the free energy

$$G_{\text{gas,RRHO}}^{\circ}(\text{X}) = G_{\text{gas,RRHO}}^{\circ(C_1)}(\text{X}) + RT \ln(\sigma_{\text{X}}) \quad (6)$$

As an example, CB7 has  $D_{7h}$  symmetry and a corresponding  $\sigma$  value of 14, in which case the correction contributes  $1.56 \text{ kcal mol}^{-1}$  to the free energy at 298 K.

### Anharmonicity and low frequency modes

Host-guest complexes can exhibit very low frequency vibrations on the order of  $50 \text{ cm}^{-1}$  or less, which tend to dominate the vibrational entropy contribution.<sup>4</sup> Many researchers have questioned whether the harmonic approximation is valid for such low frequency modes and this is an open research question. The main problem is that it is very difficult to compute the vibrational entropy exactly. Most methods for computing anharmonic effects are developed to obtain the 1 or 2 lowest energy states, but for very low frequency modes 10–20 states are likely significantly populated at room temperature and therefore contribute to the entropy.

In the absence of theoretical benchmarks, comparison to experiment can prove constructive. Kjærgaard and co-workers<sup>9,10</sup> have recently measured standard binding free energies for small gas phase compounds and compared them to CCSD(T)/aug-cc-pVT+d calculations. For example, in the case of acetonitrile-HCl the measured binding free energy at 295 K is between 1.2 and 1.9  $\text{kcal mol}^{-1}$ , while the predicted value is 1.9  $\text{kcal mol}^{-1}$  using the harmonic approximation.<sup>10</sup> Since the errors in  $\Delta E$  and the rigid-rotor approximation presumably are quite low, this suggests an error in the vibrational free energy of at most 0.7  $\text{kcal mol}^{-1}$ , despite the fact that the lowest vibrational frequency is only about  $30 \text{ cm}^{-1}$ . Furthermore, the error can be reduced by 0.4  $\text{kcal mol}^{-1}$  by scaling the harmonic frequencies by anharmonic scaling factors suggested by Shields and co-workers.<sup>11,12</sup> Similar results were found for dimethylsulfide-HCl.<sup>9</sup> So there are some indications that the harmonic approximation yields free energy corrections that are reasonable and possibly can be improved upon by relatively minor corrections.

On the other hand in a recent study Piccini and Sauer<sup>13</sup> show that anharmonic effects need to be included to obtain agreement with the experimental binding free energy of methane to H-CHA zeolite. Specifically, they compute the vibrational binding free energy by computing the 1-dimensional potential energy surface for each low frequency mode and compute the vibrational energy levels and corresponding partition function numerically (as opposed to using the anharmonic fundamental frequency together

with the harmonic oscillator partition function). This decreases the binding free energy by 2.5  $\text{kcal mol}^{-1}$  compared to the standard harmonic oscillator treatment.

Grimme<sup>4</sup> has taken a different approach by arguing that low-frequency modes resemble free rotations and using the corresponding entropy term for low frequency modes. This changes the RRHO free energy correction by 0.5–4  $\text{kcal mol}^{-1}$ , depending on the system.

Low frequencies are especially susceptible to numerical error and it is not unusual to see 1 or 2 imaginary frequencies of low magnitude in a vibrational analysis of a host-guest complex. Since imaginary frequencies are excluded from the vibrational free energy this effectively removes 1 or 2 low frequency contributions to the vibrational free energy. For example, a  $30 \text{ cm}^{-1}$  frequency contributes about 1.7  $\text{kcal mol}^{-1}$  to the free energy at 298 K.

Imaginary frequencies resulting from a flat PES and numerical errors can often be removed by making the convergence criteria for the geometry optimization and electronic energy minimization more stringent and making the grid size finer in the case of DFT calculations. If the Hessian is computed using finite difference it is important to use central-differencing. If all else fails, it is probably better to pretend that the imaginary frequency is real and add the corresponding vibrational free energy contribution. However, this needs to be systematically tested.

### Conformations

One of the main problems in computing accurate binding free energies is to identify the structures of the host, guest and (especially) the host-guest complex with the lowest free energy. Because both the RRHO and solvation energy contributions contribute greatly to the binding free energy change, simply finding the structure with the lowest electronic energy and computing the free energy only for that conformation is unlikely to result in the global free energy minimum.

For a molecule (X) with  $N_{\text{conf}}$  conformations the standard free energy is

$$G_{\text{aq}}^{\circ}(\text{X}) = G_{\text{aq}}^{\circ}(\text{X}_{\text{ref}}) - RT \ln \left( 1 + \sum_{\substack{i=1 \\ i \neq \text{ref}}}^{N_{\text{conf}}-1} e^{-\Delta G_{\text{aq}}^{\circ}(\text{X}_i)/RT} \right) \quad (7)$$

where

$$\Delta G_{\text{aq}}^{\circ}(\text{X}_i) = G_{\text{aq}}^{\circ}(\text{X}_i) - G_{\text{aq}}^{\circ}(\text{X}_{\text{ref}}) \quad (8)$$

and where  $\text{X}_{\text{ref}}$  is some arbitrarily chosen reference geometry – for example the global minimum. With that choice for  $\text{X}_{\text{ref}}$ , conformations with free energies higher than 1.36  $\text{kcal mol}^{-1}$  contribute less than 0.1 to the sum at 298 K. So a significant number of very low free energy structures are needed to make even a 0.5  $\text{kcal mol}^{-1}$  contribution to the free energy. Conformations related by symmetry should not be included here as their effects are accounted for in the rotational entropy (see above). Note that if the binding measurements are done for racemic mixtures then all stereoisomers must be included in the sum.



## Molecular charge and pH

Virtually all binding measurements in aqueous solution are performed in a buffer with a constant pH and many ligands and/or receptors contain one or more ionizable groups. The charge ( $q$ ) of an ionizable (acid/base) group in aqueous solution depends on its  $pK_a$  and the pH:

$$q = \frac{1}{1 + 10^{pH-pK_a}} - \delta \quad (9)$$

where  $\delta$  is 1 for an acid and 0 for a base. This is an average charge for all the molecules in solution and will not be an integer. This section describes how to handle charges that differ significantly from an integer value and/or change as a result of binding. The  $pK_a$  can be computed using electronic structure theory or empirically using software such as Marvin.<sup>14</sup> However, if the  $pK_a$  value is perturbed by the binding the situation may be complicated further. Here I illustrate this point for a simple example where the ligand has a basic group that is neutral when deprotonated and the receptor is non-ionizable.



The apparent equilibrium constant is then (throughout this paper I assume ideal solutions where the activity is equal to the concentration)

$$K' = \frac{[RL] + [RLH^+]}{[R]([L] + [LH^+])} \quad (10)$$

and the corresponding binding free energy is

$$\begin{aligned} \Delta G'_{aq} &= \Delta G_{aq}^{\circ}(+) - RT \ln \left( \frac{1 + 10^{pH-pK_a^c}}{1 + 10^{pH-pK_a^f}} \right) \\ &= \Delta G_{aq}^{\circ}(0) - RT \ln \left( \frac{1 + 10^{pK_a^c-pH}}{1 + 10^{pK_a^f-pH}} \right) \end{aligned} \quad (11)$$

where  $\Delta G_{aq}^{\circ}(+)$  and  $\Delta G_{aq}^{\circ}(0)$  is the binding free energy computed using the charged (protonated) and neutral form of the ligand and  $pK_a^c$  and  $pK_a^f$  are the  $pK_a$  values the ligand bound to the receptor and the free ligand, respectively.

For example, Koner *et al.*<sup>15,16</sup> have shown that binding of benzimidazole and derivatives to CB7 can increase the  $pK_a$  of the ligand by as much as 4 pH units (from  $pK_a^f = 4.6$  to  $pK_a^c = 8.6$ ) which results in a 3.3 kcal mol<sup>-1</sup> pH-dependent correction to the binding free energy at pH 7. Put another way, using  $pK_a^f$  to determine the protonation state of the bound ligand would result in an 3.3 kcal mol<sup>-1</sup> error in the binding free energy.

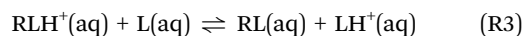
For many ligands of interest the  $pK_a^c$  can be estimated fairly accurately in a matter of second using programs such as Marvin. The effect of binding on  $pK_a^c$  can often be estimated by chemical intuition since hydrogen bonds to charged acid and basic groups tend to, respectively, lower or raise the  $pK_a$  even further. For example, if an amine with  $pK_a^f = 9$  binds to the receptor *via* hydrogen bonding, then  $pK_a^c$  is likely higher than 9 and  $\Delta G'_{aq} \approx \Delta G_{aq}^{\circ}(+)$  is a good approximation. However, if  $pK_a^f$  is close to 7 then  $pK_a^c$  should be computed. Also, it is possible

for charged ligands to change to their neutral state if they bind to hydrophobic or similarly charged receptors.

If  $pK_a^f$  is known with some degree of confidence (*e.g.* from experiment or Marvin) then  $pK_a^c$  can be estimated by

$$pK_a^c = pK_a^f - \frac{\Delta G_{\Delta pK_a, aq}^{\circ}}{RT \ln(10)} \quad (12)$$

where  $\Delta G_{\Delta pK_a, aq}^{\circ}$  is the free energy change for this reaction<sup>17</sup>



However, if one suspects that empirical  $pK_a$  predictors such as Marvin give inaccurate results for  $pK_a^f$  then this value can be computed using quantum chemistry. Ho and Coote<sup>18</sup> have written a very useful summary of different approaches to such predictions. The accuracy for phenol and carboxyl  $pK_a$  values are as low as 1 pH units (unfortunately they did not give a value for amines). However, if the  $pK_a$  value is close to the pH of interest a 1 pH unit-error can lead to prediction of the wrong protonation and result in errors in the binding free energy on the order of 1–3 kcal mol<sup>-1</sup>.

If there are several ( $N_{ionz}$ ) ionizable groups then eqn (11) generalizes to

$$\Delta G'_{aq} = \Delta G_{aq}^{\circ}(-/+) - RT \ln \left( \frac{\sum_{i=1}^{N_{ionz}} \frac{1 + 10^{s_i(pH-pK_{a,i}^c)}}{1 + 10^{s_i(pH-pK_{a,i}^f)}}}{1 + 10^{s_i(pH-pK_{a,i}^f)}} \right) \quad (13)$$

where  $\Delta G_{aq}^{\circ}(-/+)$  is the binding free energy when all acids and bases are deprotonated and protonated, respectively, the sum runs over all ionizable groups and  $s_i$  is 1 and  $-1$  if  $i$  is a base or acid, respectively.

However, this assumes that the ionizable groups titrate independently of one another, *i.e.* that the  $pK_a$  value of one group is independent of the protonation states of all other ionizable groups. If that is not the case then it is difficult to give a general expression for the pH-dependent free energy correction in terms of  $pK_a$  values (though it can be derived for a specific case). Next I present an alternative approach, but note that in practice because one can obtain more accurate *relative*  $pK_a$  values (using eqn (12)) or similar<sup>18</sup> than absolute  $pK_a$  values it may be worth the extra effort to derive the pH-dependent free energy correction in terms of  $pK_a$  values.

### Legendre transformed free energies

Instead a general expression can be written in terms of Legendre transformed free energies as suggested by Alberty<sup>19,20</sup> and modified here to electronic structure calculations:<sup>21</sup>

$$G'_{aq}(\bar{X}) = -RT \ln \left( \sum_{i=1}^{2^{N_{ionz}}} e^{-G'_{aq}(X_i)/RT} \right) \quad (14)$$

where  $\bar{X}$  denotes an average over several protonation states of X,  $2^{N_{ionz}}$  is the number of possible protonation states given  $N_{ionz}$  sites and

$$G'_{aq}(X_i) = G_{aq}^{\circ}(X_i) - n_i(H^+) (\Delta G_{sol}^{\circ}(H^+) - RT \ln(10)pH) \quad (15)$$



**Table 1** Common continuum solvation models used with electronic structure theory, the level of theory used for parameterization and the solvation energy of the proton used as a reference for the experimental solvation energies of ions used in the parameterization. Adapted from Ho<sup>27</sup>

Method	Level of theory used for parameterization	Solvation energy of proton used as reference for ions
IEFPCM-MST <sup>a</sup>	HF/6-31+G(d)	−264.0 kcal mol <sup>−1</sup>
DPCM-UAHF <sup>b</sup>	HF/6-31(+G(d)) <sup>c</sup>	−261.4 kcal mol <sup>−1</sup>
PCM-UAKS <sup>d</sup>	PBE1PBE/6-31G(d)	Unknown
IEFPCM-SMD <sup>e,f</sup>	M05-2X98/MIDI!6D M05-2X/6-31G* M05-2X/6-31+G** M05-2X/cc-pVTZ B3LYP/6-31G* HF/6-31G*	−265.9 kcal mol <sup>−1</sup>
COSMO-RS <sup>g</sup> SM8 <sup>h</sup>	BP/TZVP Independent of level of theory	Not specifically parameterized for ions −265.9 kcal mol <sup>−1</sup>

<sup>a</sup> Ref. 28 IEF and CPCM give virtually identical results for water. <sup>b</sup> Ref. 29 UAHF spheres have been used with CPCM with good results. <sup>c</sup> Diffuse functions are used only for anions. <sup>d</sup> This parameterization has not been published and the information is taken from the Gaussian09 manual. The method has been benchmarked for CPCM by Takano and Houk.<sup>30</sup> <sup>e</sup> Ref. 31. <sup>f</sup> The parameterization was performed by minimizing the error for all six methods simultaneously and any of the six methods can be used with the same parameter set. <sup>g</sup> Ref. 32. <sup>h</sup> Ref. 33.

where  $n_i(\text{H}^+)$  is the number of ionizable protons in protonation state  $i$ , and  $\Delta G_{\text{sol}}^{\circ}(\text{H}^+)$  is the solvation free energy of the proton. So in the case of ligand L considered above,  $n_i(\text{H}^+)$  is 0 and 1 for L and  $\text{LH}^+$ , respectively.

$\Delta G_{\text{sol}}^{\circ}(\text{H}^+)$  is usually taken from the literature where estimates vary between −264 and −266 kcal mol<sup>−1</sup>,<sup>22</sup> which can add to the uncertainty in the predicted binding free energy change. There are at least two ways of reducing the error. One way is to maximize error cancelation by computing  $\Delta G_{\text{sol}}^{\circ}(\text{H}^+)$  (using explicit solvent molecules as discussed below) using the same level of theory method use to compute  $\Delta G_{\text{b, aq}}^{\circ}$ . The other way is to choose the value of  $\Delta G_{\text{sol}}^{\circ}(\text{H}^+)$  used as reference for the experimental solvation free energies of ions that are used to parameterize the continuum solvation model you use (Table 1). The first way is best if explicit solvent molecules are used to compute the solvation free energies of ions in the binding study and otherwise the second method is best.

Using Legendre transformed free energies, eqn (1) can be rewritten as

$$\Delta G_{\text{b, aq}}^{\circ} = G_{\text{aq}}^{\circ}(\text{R}\bar{\text{L}}) - G_{\text{aq}}^{\circ}(\text{R}) - G_{\text{aq}}^{\circ}(\bar{\text{L}}) \quad (16)$$

Since the electronic energy contribution to the standard free energy can be very large in magnitude this form is more easily evaluated

$$G_{\text{aq}}^{\circ}(\bar{\text{X}}) = G_{\text{aq}}^{\circ}(\text{X}_{\text{ref}}) - RT \ln \left( 1 + \sum_{\substack{i=1 \\ i \neq \text{ref}}}^{2N_{\text{ionz}}-1} e^{-\Delta G_{\text{aq}}^{\circ}(\text{X}_i/RT)} \right) \quad (17)$$

where

$$\Delta G_{\text{aq}}^{\circ}(\text{X}_i) = G_{\text{aq}}^{\circ}(\text{X}_i) - G_{\text{aq}}^{\circ}(\text{X}_{\text{ref}}) \quad (18)$$

and where  $\text{X}_{\text{ref}}$  is some arbitrarily chosen reference protonation state, for example that for which  $n_i(\text{H}^+) = 0$ . The sum can be combined with that over different conformations [eqn (7)] as discussed below.

## Other ions and ionic strength

If the ligand and/or hosts contain ionizable groups then the binding measurements were likely performed in a buffer, with a certain ionic strength, to regulate pH. It is possible to include this effect in continuum solvation models such as the PCM method.<sup>23</sup> However, given the relatively low (10–100 mM) concentrations usually used in the experiments this will only have a noticeable (>0.5 kcal mol<sup>−1</sup>) effect on the energetics involving multiply charged ions. As discussed below, the error in the computed solvation energy for such ions is already large and it is not clear whether it is worth including non-specific ionic strength effects in the computations. At high ion concentrations, it is possible that these ions bind at certain sites in the ligand, receptor, or ligand–receptor complex with sufficient probability that they must be included in the thermodynamics. If so the exact same equations and considerations outlined above for  $\text{H}^+$  also apply to, e.g.  $\text{Cl}^-$  and  $\text{pCl}^-$  (computed from the specified buffer concentration) is used instead of pH.

## Solvation thermodynamics

### Background

Most continuum models (CMs) of solvation compute the solvation free energy as the difference between the free energy in solution ( $G_{\text{soln}, E}^{\circ, \text{CM}}(\text{X})$ ) and the gas phase electronic energy ( $E_{\text{gas}}(\text{X})$ )

$$\Delta G_{\text{sol}}^{\circ}(\text{X}) = G_{\text{soln}, E}^{\circ, \text{CM}}(\text{X}) - E_{\text{gas}}(\text{X}) \quad (19)$$

$G_{\text{soln}, E}^{\circ, \text{CM}}(\text{X})$  typically contains energy terms describing the electrostatic interaction of the molecule and the continuum as well as the van der Waals interactions with the solvent and free energy required to create the molecular cavity in the solvent (cavitation). The electrostatic interaction with the solvent alters the molecular wavefunction and is computed self-consistently. Usually the gas phase structure of X is used for the computation of  $G_{\text{soln}, E}^{\circ, \text{CM}}(\text{X})$ , though for COSMO-RS the structure is optimized in solution.



There is typically no explicit RRHO contribution for  $G_{\text{soln},E}^{\circ,\text{CM}}(\text{X})$  so the computational cost is comparable to that for  $E_{\text{gas}}(\text{X})$ .

Some software packages automatically compute  $\Delta G_{\text{soln},E}^{\circ,\text{CM}}(\text{X})$  and  $E_{\text{gas}}(\text{X})$  in one run, while other packages only compute  $G_{\text{soln},E}^{\circ,\text{CM}}(\text{X})$ . Also, some programs just compute the electrostatic component of  $G_{\text{soln},E}^{\circ,\text{CM}}(\text{X})$  by default. However, the van der Waals and, especially, the cavitation component can make sizable contributions to the binding free energy and must be included for accurate results. It is worth noting that any hydrophobic contribution to binding will derive primarily from the change in cavitation energy.<sup>24</sup>  $G_{\text{soln},E}^{\circ,\text{CM}}(\text{X})$  contains parameters (e.g. atomic radii) that are adjusted to reproduce experimentally measured solvation free energies

$$\Delta G_{\text{solv}}^{\circ,\text{exp}}(\text{X}) = G_{\text{soln}}^{\circ,\text{exp}}(\text{X}) - G_{\text{gas}}^{\circ,\text{exp}}(\text{X}) \quad (20)$$

The standard state for both  $G_{\text{soln}}^{\circ,\text{exp}}(\text{X})$  and  $G_{\text{gas}}^{\circ,\text{exp}}(\text{X})$  is generally chosen to be 1 M.<sup>25,26</sup> The latter is the reason a 1 M reference state also must be used when computing  $G_{\text{gas,RRHO}}^{\circ}(\text{X})$ .

Notice that the volume on going from the gas phase to solution is included in the solvation free energy

$$\Delta G_{\text{solv}}^{\circ}(\text{X}) = \Delta A_{\text{solv}}^{\circ,\text{exp}}(\text{X}) + p^{\circ}(\Delta V_{\text{solv}} - V_{\text{gas}}) \quad (21)$$

where  $\Delta V_{\text{solv}}$  is the volume change in solution due to addition of the solute X to the neat solvent. For an ideal gas ( $p^{\circ}V_{\text{gas}} = RT$ ) it follows that

$$\Delta \Delta G_{\text{solv}}^{\circ} = \Delta \Delta A_{\text{solv}}^{\circ} + p^{\circ} \Delta \Delta V_{\text{solv}} - RT \quad (22)$$

and

$$\Delta G_{\text{b,aq}}^{\circ} = \Delta A_{\text{b,aq}}^{\circ} + p^{\circ} \Delta V_{\text{soln}} \quad (23)$$

because the  $-RT$  term is cancelled by a corresponding term in the translational enthalpy contribution to  $\Delta G_{\text{gas,RRHO}}^{\circ}$ .  $\Delta V_{\text{soln}} = \Delta \Delta V_{\text{solv}}$  is the change in the volume of the solution on upon binding.

### Atomic radii

The solvation energy is computed using a set of atomic radii that define the solute–solvent boundary surface. These radii are usually obtained by fitting to experimentally measured solvation energies. Accurate solvation energies should not be expected from methods that use iso-electron density surfaces or van der Waals radii without additional empirical fitting. When using fitted radii one should use the same level of theory for the solute as was used in the parameterization (Table 1).

### Ions

For neutral molecules solvation free energies can be measured with an accuracy of roughly 0.2 kcal mol<sup>-1</sup> and reproduced theoretically to within roughly 0.5–1.0 kcal mol<sup>-1</sup>, depending on the method. However, the solvation energies of ions cannot be directly measured and must be indirectly inferred relative to a standard (usually the solvation energy of the proton). The experimentally obtained solvation energies are typically accurate to within 3 kcal mol<sup>-1</sup> and can be reproduced computationally

with roughly the same accuracy.<sup>22</sup> The solvation energy of ions are therefore an especially likely source of error in binding free energies – especially if the ionic regions of the molecules become significantly desolvated due to binding.

### Gas phase vs. solution optimization

The fitting of the radii described above is usually done using gas phase optimized structures only, *i.e.* any change in structure and corresponding rotational and vibrational effects are “included” in the radii *via* the parameterization. However, for ionic species gas phase optimization can lead to significantly distorted structures or even proton transfer and in these cases solution phase optimizations and, hence, vibrational frequency calculations, tend to be used. However, numerical noise in the continuum models can make it necessary to increase (*i.e.* make less stringent) the geometry convergence criteria and can lead to more imaginary frequencies than in the gas phase. One option is to compute the vibrational contribution to  $\Delta G_{\text{gas,RRHO}}^{\circ}$  using gas phase optimized structures as Sure *et al.* have done.<sup>2</sup>

When using solution phase geometries the gas phase single point energies needed to evaluate  $\Delta G_{\text{solv}}^{\circ}(\text{X})$  represent added computational expense one option is to use solution phase free energies to evaluate the binding free energies

$$\Delta G_{\text{b,aq}}^{\circ} = \Delta G_{\text{b,soln},E}^{\circ,\text{CM}} + \Delta G_{\text{b,soln,RRHO}}^{\circ,\text{CM}} \quad (24)$$

One problem with this approach is that  $\Delta G_{\text{b,soln},E}^{\circ,\text{CM}}$ , unlike  $\Delta E_{\text{gas}}$ , is not systematically improveable due to the empirical parameterization. For a more thorough discussion of this issue see Ho *et al.*,<sup>34</sup> Ribeiro *et al.*,<sup>35</sup> and Ho.<sup>27</sup>

### Cavities

The atomic radii and corresponding cavity generation algorithms are parameterized for small molecules. For more complex molecules such as receptors this can lead to continuum solvation of regions of molecules, *e.g.* deep in the binding pocket, that are not accessible to the molecular solvent. Furthermore, any solvent molecule inside such pocket is likely to be quite “un-bulk-like” and not well-represented by the bulk solvent or fixed by the underlying parametrization. However, how big an error this may introduce to the binding free energy is not really known, but certain models for the cavitation energy have been shown to give unrealistically large contributions to the binding free energy.<sup>36,37</sup>

### Explicit water molecules

Adding explicit solvent molecules to the receptor and/or ligand can potentially lead to more accurate results. For example, including explicit water molecules around ionic sites reduces the strong dependence of the solvation energy on the corresponding atomic radii. Also, “un-bulk-like” water molecules now are treated more naturally and the risk of solvating non-solvent-accessible regions is reduced somewhat. However, adding explicit solvent molecules increases the computational cost by increasing the CPU time needed to compute energies, perform conformational searches, and compute vibrational frequencies.



There are several approaches to include the effect of explicit solvent molecules in the binding free energy. Bryantsev *et al.*<sup>38</sup> suggest computing the solvation energy by

$$G_{\text{aq},n}^{\circ}(\text{X}) = G_{\text{gas}}^{\circ}(\text{X}) + \Delta G_{\text{solv},n}^{\circ}(\text{X}) \quad (25)$$

where

$$\Delta G_{\text{solv},n}^{\circ}(\text{X}) = \Delta G_{\text{gas}}^{\circ}(\text{X}(\text{H}_2\text{O})_n) + \Delta G_{\text{solv}}^{\circ}(\text{X}(\text{H}_2\text{O})_n) - \Delta G_{\text{solv}}^{\circ}((\text{H}_2\text{O})_n) \quad (26)$$

(note that  $\Delta G_{\text{solv},0}^{\circ}(\text{X}) = \Delta G_{\text{solv}}^{\circ}(\text{X})$ ) and

$$\Delta G_{\text{gas}}^{\circ}(\text{X}(\text{H}_2\text{O})_n) = G_{\text{gas}}^{\circ}(\text{X}(\text{H}_2\text{O})_n) - G_{\text{gas}}^{\circ}(\text{X}) - G_{\text{gas}}^{\circ}((\text{H}_2\text{O})_n) \quad (27)$$

and

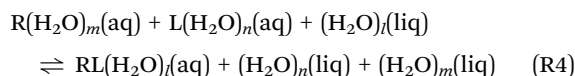
$$\Delta G_{\text{solv}}^{\circ}(\text{liq})((\text{H}_2\text{O})_n) = \Delta G_{\text{solv}}^{\circ}((\text{H}_2\text{O})_n) + RT \ln([\text{H}_2\text{O}]/n) \quad (28)$$

with “(liq)” referring to a standard state of 55.34 M (the concentration of liquid water at 298 K), respectively. The term  $RT \ln([\text{H}_2\text{O}]/n)$  is the free energy required to change the standard state of  $(\text{H}_2\text{O})_n$  from 1 M to 55.34/n M.

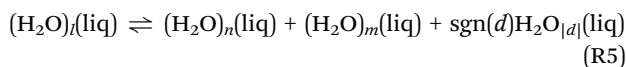
Bryantsev *et al.* have shown that using this water cluster approach leads to a smooth convergence of the solvation free energy with respect to the cluster size  $n$ . The optimum choice of  $n$  is the one where an additional water molecule changes the solvation energy by less than a certain amount defined by the user. One can thereby compute the optimum number of water molecules for the receptor ( $n$ ), ligand ( $m$ ) and receptor–ligand complex ( $l$ ) and then compute the change in solvation free energy as

$$\Delta \Delta G_{\text{b,solv},x}^{\circ} = \Delta G_{\text{solv},l}^{\circ}(\text{RL}) - \Delta G_{\text{solv},n}^{\circ}(\text{L}) - \Delta G_{\text{solv},m}^{\circ}(\text{R}) \quad (29)$$

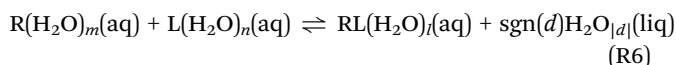
and computing  $\Delta E_{\text{gas}}$  and  $\Delta G_{\text{gas,RRHO}}^{\circ}$  as before. One can show that this corresponds to the free energy change for this reaction



In principle, the free energy is zero for



where  $d = l - m - n$  and  $\text{sgn}(d)$  returns the sign of  $d$ . So the free energy change for reaction (4) can also be computed as the free energy change for

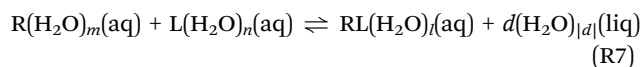


However, this is only approximately true in practice due to errors in the computed gas phase and solvation free energies. Furthermore, reaction (6) does not really lead to any significant reduction in CPU time because the water cluster free energies only have to be computed once. However, if reaction (6) is

used then one must add an additional term correcting for the indistinguishability of water molecules

$$G_{\text{gas,RRHO}}^{\circ}(\text{X}(\text{H}_2\text{O})_n) \rightarrow G_{\text{gas,RRHO}}^{\circ}(\text{X}(\text{H}_2\text{O})_n) - RT \ln(n!) \quad (30)$$

and similarly for the water clusters. Using reaction (4) leads to a cancellation of this term and also maximizes error cancellation in the other energy terms. Similar considerations apply when using individual water molecules to balance the reaction instead of water clusters



One of the main reasons reaction (4) maximizes error cancellation is that the number and type of hydrogen bonds involving water molecules are very similar on each side of the equilibrium. This can also be achieved when using reaction (6) or (7) by ensuring that  $l = m + n$ , in which case the error cancellation may be comparable and will depend on the nature of the ligand, host, and water arrangement. However, eqn (30) must still be used when using reaction (6) or (7) in this way.

When using many explicit water molecules the error in the continuum solvation energies can be reduced by ensuring that the continuum solvation energy of a single water molecule matches the experimental value of  $-6.32 \text{ kcal mol}^{-1}$  at 298.15 K as close as possible.

## Enthalpy and entropy contributions to the binding free energy

It is often instructive to decompose the binding free energy into enthalpy and entropy contributions. The standard enthalpy and entropy of molecule X in aqueous solution is

$$H_{\text{aq}}^{\circ}(\text{X}) = E_{\text{gas}}^{\circ}(\text{X}) + H_{\text{gas,RRHO}}^{\circ}(\text{X}) + \Delta H_{\text{solv}}^{\circ}(\text{X}) \quad (31)$$

and

$$S_{\text{aq}}^{\circ}(\text{X}) = S_{\text{gas,RRHO}}^{\circ}(\text{X}) + \Delta S_{\text{solv}}^{\circ}(\text{X}) \quad (32)$$

where the standard state eqn (3) and symmetry correction [eqn (6)] is applied to the entropy term. Thus, in order to compute these quantities one must compute the enthalpy and entropy of solvation, which can be done by the COSMO-RS<sup>32</sup> and SM8T<sup>39</sup> solvation methods. Chamberlin *et al.*<sup>39</sup> have noted that most of the temperature dependence of the aqueous solvation free energy comes from the non-polar term so simply including the effect of temperature on the dielectric constant is unlikely to give accurate results. Plata and Singleton<sup>40</sup> have recently shown that  $\Delta S_{\text{solv}}^{\circ}(\text{X})$  can make an appreciable contribution to the energy change for reaction energies.

For a molecule (X) with  $N_{\text{conf}}$  conformations the standard enthalpy and entropy is

$$H_{\text{aq}}^{\circ}(\text{X}) = \sum_{i=1}^{N_{\text{conf}}} H_{\text{aq}}^{\circ}(\text{X}_i) p(\text{X}_i) \quad (33)$$



and

$$S_{\text{aq}}^{\circ}(\mathbf{X}) = \sum_{i=1}^{N_{\text{conf}}} S_{\text{aq}}^{\circ}(\mathbf{X}_i) p(\mathbf{X}_i) - R \sum_{i=1}^{N_{\text{conf}}} p(\mathbf{X}_i) \ln(p(\mathbf{X}_i)) \quad (34)$$

where

$$p(\mathbf{X}_i) = \frac{e^{-\Delta G_{\text{aq}}^{\circ}(\mathbf{X}_i)/RT}}{\sum_{i=1}^{N_{\text{conf}}} e^{-\Delta G_{\text{aq}}^{\circ}(\mathbf{X}_i)/RT}} \quad (35)$$

and  $\Delta G_{\text{aq}}^{\circ}(\mathbf{X}_i)$  is computed relative to the conformation with the lowest free energy.

The Legendre transformed entropy and enthalpy is

$$\begin{aligned} S_{\text{aq}}^{\prime\circ}(\mathbf{X}_i) &= - \left( \frac{\partial G_{\text{aq}}^{\prime\circ}(\mathbf{X}_i)}{\partial T} \right)_{p,\text{pH}} \\ &= S_{\text{aq}}^{\circ}(\mathbf{X}_i) - n_i(\text{H}^+) (\Delta S_{\text{solv}}^{\circ}(\text{H}^+) + R \ln(10) \text{pH}) \end{aligned} \quad (36)$$

and

$$H_{\text{aq}}^{\prime\circ}(\mathbf{X}_i) = H_{\text{aq}}^{\circ}(\mathbf{X}_i) - n_i(\text{H}^+) \Delta H_{\text{solv}}^{\circ}(\text{H}^+) \quad (37)$$

When comparing computed enthalpy and entropy changes to experimental measurements on systems with ionizable groups note that the observed values will depend on the buffer used *if* protonation states change upon binding (see *e.g.* ref. 41). Unless the experimental study has corrected for this effect by repeating the measurements in different buffers, this effect can contribute to the difference between the computed and experimental values.

## A concrete example

In this section I apply the key equations discussed above to a specific example: *p*-xylylenediamine (L, Fig. 2) binding to CB7 (R) for which a binding free energy of  $-9.9 \pm 0.1$  kcal mol<sup>-1</sup> has been measured at pH 7.4 and 298 K.<sup>1</sup> The conformations and other details such as the number of water molecules are just selected and constructed *for illustration purposes only* using the Avogadro program<sup>42</sup> and the MMFF force field and should not be considered accurate.

CB7 has one conformation with  $D_{7h}$  symmetry and no ionizable groups. It is assumed that the solvation energy can be computed accurately without explicit water molecules. Thus, the free energy in aqueous solution is

$$G_{\text{aq}}^{\circ}(\mathbf{R}) = G_{\text{aq},0}^{\circ}(\mathbf{R}) = G_{\text{gas}}^{\circ}(\mathbf{R}) + \Delta G_{\text{solv}}^{\circ}(\mathbf{R}) + RT \ln(14) \quad (38)$$

where  $G_{\text{gas}}^{\circ}(\mathbf{R})$  is computed in  $C_1$  symmetry and 14 is the symmetry number ( $\sigma$ ) corresponding to the  $D_{7h}$  point group.

Ligand L has two basic groups and is assumed to have two conformations *a* and *b* for each protonation state. The  $\text{p}K_{\text{a}}$  values for the basic groups are 9.2 and 9.8 according to Marvin, so both groups are likely 100% protonated at pH 7. However, for illustration purposes I will include all three protonation states in the computation of the free energy. Furthermore, I will assume that each charged amine group is microsolvated by three explicit water molecules.

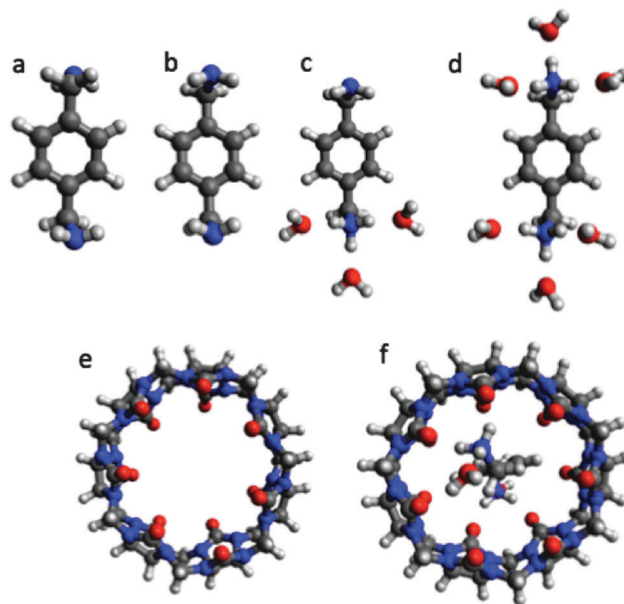


Fig. 2 Representative conformations of ligand L (*p*-xylylenediamine), receptor R (CB7), and a receptor–ligand complex RL used to illustrate the use of the equations presented in this paper. (a)  $L_a$ , (b)  $L_b$ , (c)  $L\text{H}^+b$ , (d)  $L\text{H}_2^{2+}b$ , (e) R, and (f)  $RL\text{H}_2^{2+}a$ . The coordinates for the structures are available here: <http://dx.doi.org/10.6084/m9.figshare.1290639>.

The free energy of conformer *a* of the doubly protonated state ( $L\text{H}_2^{2+}$ ) is thus

$$\begin{aligned} G_{\text{aq},6}^{\circ}(L\text{H}_2^{2+}a) &= G_{\text{gas}}^{\circ}(L\text{H}_2^{2+}(\text{H}_2\text{O})_6a) + \Delta G_{\text{solv}}^{\circ}(L\text{H}_2^{2+}(\text{H}_2\text{O})_6a) \\ &\quad - G_{\text{gas}}^{\circ}((\text{H}_2\text{O})_6) - \Delta G_{\text{solv}}^{\circ}((\text{H}_2\text{O})_6) \\ &\quad - RT \ln([\text{H}_2\text{O}]/6) - RT \ln(2) \end{aligned} \quad (39)$$

where the gas phase energy is computed in  $C_1$  symmetry and 2 is the symmetry number of the  $C_2$  point group. The lowest energy structure of  $(\text{H}_2\text{O})_6$  suggested by Bransyev *et al.* can be used for compute  $G_{\text{aq},n}^{\circ}((\text{H}_2\text{O})_6)$ , or the effect of additional conformations can be included using eqn (7). Finally, the Legendre transformed free energy [eqn (15)] at pH 7 is computed by

$$G_{\text{aq},6}^{\prime\circ}(L\text{H}_2^{2+}a) = G_{\text{aq},6}^{\circ}(L\text{H}_2^{2+}a) - 2(\Delta G_{\text{solv}}^{\circ}(\text{H}^+) - RT \ln(10) \text{pH}) \quad (40)$$

The corresponding free energy of conformer *b*,  $G_{\text{aq},6}^{\prime\circ}(L\text{H}_2^{2+}b)$ , which has  $C_{2v}$  symmetry and for which  $\sigma$  is also 2, is computed in the same way. Notice that each conformation in principle can have different numbers of water associated with them. Similarly, the free energies of the singly protonated and neutral ligand (with  $C_1$  and  $C_2$  symmetry) is computed by

$$\begin{aligned} G_{\text{aq},3}^{\prime\circ}(L\text{H}^+a) &= G_{\text{gas}}^{\circ}(L\text{H}^+(\text{H}_2\text{O})_3a) + \Delta G_{\text{solv}}^{\circ}(L\text{H}^+(\text{H}_2\text{O})_3a) \\ &\quad - G_{\text{gas}}^{\circ}((\text{H}_2\text{O})_3) - \Delta G_{\text{solv}}^{\circ}((\text{H}_2\text{O})_3) \\ &\quad - RT \ln([\text{H}_2\text{O}]/3) \\ &\quad - (\delta G_{\text{solv}}^{\circ}(\text{H}^+) - RT \ln(10) \text{pH}) \end{aligned} \quad (41)$$





and

$$G_{\text{aq},0}^{\circ}(\text{LHa}) = G_{\text{gas}}^{\circ}(\text{La}) + \Delta G_{\text{solv}}^{\circ}(\text{La}) + RT \ln(2) \quad (42)$$

(here for conformer *a* and similarly for conformer *b*). Finally, the free energy of L averaged over conformations and protonation states is

$$\begin{aligned} G_{\text{aq},x}^{\circ}(\bar{\text{L}}) = & G_{\text{aq},0}^{\circ}(\text{La}) \\ & - RT \ln \left( 1 + e^{-\Delta G_{\text{aq},0}^{\circ}(\text{Lb})/RT} + e^{-\Delta G_{\text{aq},3}^{\circ}(\text{LH}^+a)/RT} \right. \\ & + e^{-\Delta G_{\text{aq},3}^{\circ}(\text{LH}^+b)/RT} + e^{-\Delta G_{\text{aq},6}^{\circ}(\text{LH}_2^{2+}a)/RT} \\ & \left. + e^{-\Delta G_{\text{aq},6}^{\circ}(\text{LH}_2^{2+}b)/RT} \right) \end{aligned} \quad (43)$$

where

$$\Delta G_{\text{aq},0}^{\circ}(\text{Lb}) = G_{\text{aq},0}^{\circ}(\text{Lb}) - G_{\text{aq},0}^{\circ}(\text{La}) \quad (44)$$

and similarly for the remaining terms in the sum. Notice that for each conformation there are three protonation states rather than (2<sup>2</sup>) because the two singly protonated structures are equivalent.

For the host-guest complex I have assumed that each conformation can bind CB7 in only one way and that two explicit water molecules per protonated group is lost upon binding, so that

$$\begin{aligned} G_{\text{aq},x}^{\circ}(\text{RL}) = & G_{\text{aq},0}^{\circ}(\text{RLa}) \\ & - RT \ln \left( 1 + e^{-\Delta G_{\text{aq},0}^{\circ}(\text{RLb})/RT} + e^{-\Delta G_{\text{aq},1}^{\circ}(\text{RLH}^+a)/RT} \right. \\ & + e^{-\Delta G_{\text{aq},1}^{\circ}(\text{RLH}^+b)/RT} + e^{-\Delta G_{\text{aq},2}^{\circ}(\text{RLH}_2^{2+}a)/RT} \\ & \left. + e^{-\Delta G_{\text{aq},2}^{\circ}(\text{RLH}_2^{2+}b)/RT} \right) \end{aligned} \quad (45)$$

Note that the effect of the 28 equivalent binding modes to other oxygen atoms for *e.g.* LH<sub>2</sub><sup>2+</sup>*a* (Fig. 2f) is accounted for by the symmetry factors. Finally, the binding free energy is computed using eqn (16).

## Protein–ligand binding

In order for the electronic structure approach to be used in drug design corresponding calculation have to be carried out on proteins, which are significantly larger than the hosts that have been used to benchmark the approach so far. QM/MM is of course the obvious choice for computing the geometries and gas phase energies, although linear scaling all QM methods such as the FMO<sup>43</sup> method is also possible. Furthermore, continuum methods such as PCM have been adapted for large systems and interfaced to both QM/MM<sup>44</sup> and the FMO method.<sup>45</sup> Of course as the system size increases conformational sampling will become a bigger practical issue.

The main issue is the computation of vibrational frequencies for the protein and protein–ligand complex. The fast semi-empirical methods currently used for computing the vibrational frequencies (dispersion and hydrogen bond-corrected PM6 and DFTB as well as HF-3c) must be interfaced with QM/MM codes and/or be implemented in a linear scaling approach that allow for

frequency calculations. Dispersion-corrected PM6 and DFTB are already implemented in AMBER, a FMO implementation of DFTB has recently been added to GAMESS<sup>46</sup> and a similar HF-3c/FMO implementation is forthcoming from my lab.

Most QM/MM studies of enzyme catalysis constrain the geometry of a significant portion of the system to avoid spurious structural fluctuation far away from the active site contributing to the barrier. This may well be necessary for binding free energy calculations as well, in which case the effect of the constraints on the vibrational frequencies must be accounted for.<sup>47</sup> Alternatively, only the Hessian of the un-constrained region can be computed.<sup>48</sup>

So while there is some code-adjustment to be done it may well be that the promising developments in electronic structure-based prediction of aqueous binding free energies may also be brought to bear on drug design within the next few years.

## Summary and outlook

Recent predictions of absolute binding free energies of host-guest complexes in aqueous solution using electronic structure theory have been encouraging for some systems. It is interesting to consider the underlying innovations that have led to the recent increase in accuracy in predicted binding free energies. Advances in computer hardware and coupled cluster algorithms made it possible to construct benchmark sets of accurate electronic binding energies for a diverse set of molecules. These benchmark sets were then used to develop the dispersion corrections needed for accurate DFT-based electronic binding energies and the short-range (hydrogen bond) corrections to the semi-empirical methods needed to compute accurate vibrational frequencies for the RRHO free energy corrections. In fact methods like HF-3c,<sup>49</sup> while containing empirical corrections, was developed without reference to any experimental data. Another interesting observation is that the dispersion and RRHO free energy contributions to the binding free energy have roughly the same magnitude, but opposite signs. So including just one of the corrections is likely to significantly increase the error relative to experiment and lead to the wrong conclusions regarding their importance.

While there have been reasonably accurate predictions for some host-guest systems, other systems remain problematic. In this paper I summarize some of the many factors that could easily contribute 1–3 kcal mol<sup>-1</sup> at 298 K: three-body dispersion effects, molecular symmetry, anharmonicity, spurious imaginary frequencies, insufficient conformational sampling, wrong or changing ionization states, errors in the solvation free energy of ions, and explicit solvent (and ion) effects that are not well-represented by continuum models.

While I focus on binding free energies in aqueous solution it is worth noting that the approach also applies to any free energy difference in solution, such as conformational and reaction free energy differences or activation free energies. Furthermore, the equations apply to solvents other than water as long as the concentration of liquid water, the solvation free energy of the proton changed, and the parameterization of the continuum solvation model is changed to match the solvent of interest.



Furthermore, while the recent successes with electronic structure-based approaches have been for host-guest complexes they can be extended to protein-ligand complexes with a few methodological improvements (mainly related to the computation of vibrational frequencies). Thus, it may well be that the promising developments in electronic structure-based prediction of aqueous binding free energies may also be brought to bear on drug design within the next few years.

## Acknowledgements

I thank the following people for very illuminating discussions and/or comments on the manuscript: Vyacheslav Bryantsev, Chris Cramer, Mike Gilson, Stefan Grimme, Fahmi Himo, Junming Ho, Adrian Jinich, Andreas Klamt, Pedro Silva, Dan Singleton and Casper Steinmann.

## References

- H. S. Muddana, A. T. Fenley, D. L. Mobley and M. K. Gilson, The SAMPL4 host-guest blind prediction challenge: an overview, *J. Comput.-Aided Mol. Des.*, 2014, 1–13.
- R. Sure, J. Antony and S. Grimme, Blind Prediction of Binding Affinities for Charged Supramolecular Host-Guest Systems: Achievements and Shortcomings of DFT-D3, *J. Phys. Chem. B*, 2014, **118**(12), 3431–3440.
- S. Grimme, J. Antony, S. Ehrlich and H. Krieg, A consistent and accurate *ab initio* parametrization of density functional dispersion correction (DFT-D) for the 94 elements H-Pu, *J. Chem. Phys.*, 2010, **132**(15), 154104.
- S. Grimme, Supramolecular binding thermodynamics by dispersion-corrected density functional theory, *Chem. – Eur. J.*, 2012, **18**(32), 9955–9964.
- H. S. Muddana and M. K. Gilson, Calculation of host-guest binding affinities using a quantum-mechanical energy model, *J. Chem. Theory Comput.*, 2012, **8**(6), 2023–2033.
- N. D. Yilmazer and M. Korth, Comparison of molecular mechanics, semi-empirical quantum mechanical, and density functional theory methods for scoring protein-ligand interactions, *J. Phys. Chem. B*, 2013, **117**(27), 8075–8084.
- H. X. Zhou and M. K. Gilson, Theory of free energy and entropy in noncovalent binding, *Chem. Rev.*, 2009, **109**(9), 4092–4107.
- J. J. P. Stewart, Optimization of parameters for semiempirical methods V: modification of NDDO approximations and application to 70 elements, *J. Mol. Model.*, 2007, **13**(12), 1173–1213.
- N. Bork, L. Du and H. G. Kjaergaard, Identification and Characterization of the HCl-DMS Gas Phase Molecular Complex *via* Infrared Spectroscopy and Electronic Structure Calculations, *J. Phys. Chem. A*, 2014a, **118**(8), 1384–1389.
- N. Bork, L. Du, H. Reiman, T. Kurtén and H. G. Kjaergaard, Benchmarking *Ab Initio* Binding Energies of Hydrogen-Bonded Molecular Clusters Based on FTIR Spectroscopy, *J. Phys. Chem. A*, 2014b, **118**(28), 5316–5322.
- B. Temelso, K. A. Archer and G. C. Shields, Benchmark structures and binding energies of small water clusters with anharmonicity corrections, *J. Phys. Chem. A*, 2011, **115**(43), 12034–12046.
- B. Temelso and G. C. Shields, The role of anharmonicity in hydrogen-bonded systems: the case of water clusters, *J. Chem. Theory Comput.*, 2011, **7**(9), 2804–2817.
- G. Piccini and J. Sauer, Effect of Anharmonicity on Adsorption Thermodynamics, *J. Chem. Theory Comput.*, 2014, **10**(6), 2479–2487.
- Marvin, <http://www.chemaxon.com/marvin/sketch/index.php>, 2014.
- A. L. Koner, I. Ghosh, N. I. Saleh and W. M. Nau, Supramolecular encapsulation of benzimidazole-derived drugs by cucurbit[7]uril, *Can. J. Chem.*, 2011, **89**(2), 139–147.
- M. O. Kim, P. G. Blachly, J. W. Kaus and J. A. McCammon, Protocols Utilizing Constant pH Molecular Dynamics to Compute pH-Dependent Binding Free Energies, *J. Phys. Chem. B*, 2015, **119**, 861–872.
- H. Li, A. D. Robertson and J. H. Jensen, The determinants of carboxyl pK<sub>a</sub> values in turkey ovomucoid third domain, *Proteins: Struct., Funct., Bioinf.*, 2004, **55**(3), 689–704.
- J. Ho and M. L. Coote, A universal approach for continuum solvent pK<sub>a</sub> calculations: are we there yet?, *Theor. Chem. Acc.*, 2010, **125**(1–2), 3–21.
- R. A. Alberty, *Thermodynamics of biochemical reactions*, John Wiley & Sons., 2005.
- R. A. Alberty, A. Cornish-Bowden, R. N. Goldberg, G. G. Hammes, K. Tipton and H. V. Westerhoff, Recommendations for terminology and databases for biochemical thermodynamics, *Biophys. Chem.*, 2011, **155**(2), 89–103.
- A. Jinich, D. Rappoport, I. Dunn, B. Sanchez-Lengeling, R. Olivares-Amaya, E. Noor and A. Aspuru-Guzik, Quantum Chemical Approach to Estimating the Thermodynamics of Metabolic Reactions, *Sci. Rep.*, 2014, **4**, 7022.
- C. P. Kelly, C. J. Cramer and D. G. Truhlar, Aqueous solvation free energies of ions and ion-water clusters based on an accurate value for the absolute aqueous solvation free energy of the proton, *J. Phys. Chem. B*, 2006, **110**(32), 16066–16081.
- M. Cossi, V. Barone, B. Mennucci and J. Tomasi, *Ab initio* study of ionic solutions by a polarizable continuum dielectric model, *Chem. Phys. Lett.*, 1998, **286**(3), 253–260.
- C. J. Cramer and D. G. Truhlar, General parameterized SCF model for free energies of solvation in aqueous solution, *J. Am. Chem. Soc.*, 1991, **113**(22), 8305–8311.
- A. Ben-Naim, Standard thermodynamics of transfer, Uses and misuses, *J. Phys. Chem.*, 1978, **82**(7), 792–803.
- A. Ben-Naim and Y. Marcus, Solvation thermodynamics of nonionic solutes, *J. Chem. Phys.*, 1984, **81**(4), 2016–2027.
- J. Ho, Are thermodynamic cycles necessary for continuum solvent calculation of pK<sub>a</sub>s and reduction potentials? *Phys. Chem. Chem. Phys.*, 2015, **17**(4), 2859–2868.
- C. Curutchet, A. Bidon-Chanal, I. Soteras, M. Orozco and F. J. Luque, MST continuum study of the hydration free energies of monovalent ionic species, *J. Phys. Chem. B*, 2005, **109**(8), 3565–3574.



- 29 V. Barone, M. Cossi and J. Tomasi, A new definition of cavities for the computation of solvation free energies by the polarizable continuum model, *J. Chem. Phys.*, 1997, **107**(8), 3210–3221.
- 30 Y. Takano and K. N. Houk, Benchmarking the conductor-like polarizable continuum model (CPCM) for aqueous solvation free energies of neutral and ionic organic molecules, *J. Chem. Theory Comput.*, 2005, **1**(1), 70–77.
- 31 A. V. Marenich, C. J. Cramer and D. G. Truhlar, Universal solvation model based on solute electron density and on a continuum model of the solvent defined by the bulk dielectric constant and atomic surface tensions, *J. Phys. Chem. B*, 2009, **113**(18), 6378–6396.
- 32 F. Eckert and A. Klamt, Fast solvent screening *via* quantum chemistry: COSMO-RS approach, *AIChE J.*, 2002, **48**(2), 369–385.
- 33 A. V. Marenich, R. M. Olson, C. P. Kelly, C. J. Cramer and D. G. Truhlar, Self-consistent reaction field model for aqueous and nonaqueous solutions based on accurate polarized partial charges, *J. Chem. Theory Comput.*, 2007, **3**(6), 2011–2033.
- 34 J. Ho, A. Klamt and M. L. Coote, Comment on the correct use of continuum solvent models, *J. Phys. Chem. A*, 2010, **114**(51), 13442–13444.
- 35 R. F. Ribeiro, A. V. Marenich, C. J. Cramer and D. G. Truhlar, Use of solution-phase vibrational frequencies in continuum models for the free energy of solvation, *J. Phys. Chem. B*, 2011, **115**(49), 14556–14562.
- 36 S. Genheden, J. Kongsted, P. Söderhjelm and U. Ryde, Non-polar Solvation Free Energies of Protein–Ligand Complexes, *J. Chem. Theory Comput.*, 2010, **6**(11), 3558–3568.
- 37 S. Genheden and U. Ryde, Improving the Efficiency of Protein–Ligand Binding Free-Energy Calculations by System Truncation, *J. Chem. Theory Comput.*, 2012, **8**(4), 1449–1458.
- 38 V. S. Bryantsev, M. S. Diallo and W. A. Goddard III, Calculation of solvation free energies of charged solutes using mixed cluster/continuum models, *J. Phys. Chem. B*, 2008, **112**(32), 9709–9719.
- 39 A. C. Chamberlin, C. J. Cramer and D. G. Truhlar, Extension of a temperature-dependent aqueous solvation model to compounds containing nitrogen, fluorine, chlorine, bromine, and sulfur, *J. Phys. Chem. B*, 2008, **112**(10), 3024–3039.
- 40 R. E. Plata and D. A. Singleton, A Case Study of the Mechanism of Alcohol-Mediated Morita Baylis-Hillman Reactions. The Importance of Experimental Observations, *J. Am. Chem. Soc.*, 2015, **137**, 3811–3826.
- 41 F. Dullweber, M. T. Stubbs, D. Musil, J. Stürzebecher and G. Klebe, Factorising ligand affinity: a combined thermodynamic and crystallographic study of trypsin and thrombin inhibition†, *J. Mol. Biol.*, 2001, **313**(3), 593–614.
- 42 M. D. Hanwell, D. E. Curtis, D. C. Lonie, T. Vandermeersch, E. Zurek and G. R. Hutchison, Avogadro: An advanced semantic chemical editor, visualization, and analysis platform, *J. Cheminf.*, 2012, **4**, 17.
- 43 D. G. Fedorov, T. Nagata and K. Kitaura, Exploring chemistry with the fragment molecular orbital method, *Phys. Chem. Chem. Phys.*, 2012, **14**(21), 7562–7577.
- 44 H. Li, C. S. Pomelli and J. H. Jensen, Continuum solvation of large molecules described by QM/MM: a semi-iterative implementation of the PCM/EFP interface, *Theor. Chem. Acc.*, 2003, **109**(2), 71–84.
- 45 D. G. Fedorov, K. Kitaura, H. Li, J. H. Jensen and M. S. Gordon, The polarizable continuum model (PCM) interfaced with the fragment molecular orbital method (FMO), *J. Comput. Chem.*, 2006, **27**(8), 976–985.
- 46 Y. Nishimoto, D. G. Fedorov and S. Irle, Density-Functional Tight-Binding Combined with the Fragment Molecular Orbital Method, *J. Chem. Theory Comput.*, 2014, **10**(11), 4801–4812.
- 47 A. Ghysels, D. Van Neck, V. Van Speybroeck, T. Verstraelen and M. Waroquier, Vibrational modes in partially optimized molecular systems, *J. Chem. Phys.*, 2007, **126**(22), 224102.
- 48 H. Li and J. H. Jensen, Partial Hessian vibrational analysis: the localization of the molecular vibrational energy and entropy, *Theor. Chem. Acc.*, 2002, **107**(4), 211–219.
- 49 R. Sure and S. Grimme, Corrected small basis set Hartree–Fock method for large systems, *J. Comput. Chem.*, 2013, **34**(19), 1672–1685.

