# Strategies to calculate water binding free energies in protein-ligand complexes

Michael S. Bodnarchuk,[†] Russell Viner,[‡] Julien Michel,[¶] and Jonathan W. Essex[*,†]

*School of Chemistry, University of Southampton, Highfield, Southampton, SO17 1BJ, U.K.,*

*Syngenta, Jealott's Hill International Research Centre, Bracknell, RG42 6EY, U.K., and*

*EaStCHEM School of Chemistry, University of Edinburgh, The Kings Buildings, West Mains*

*Road, Edinburgh EH9 3JJ, U.K.*

E-mail: jwe1@soton.ac.uk

**Abstract**

Water molecules are commonplace in protein binding pockets, where they can typically form a complex between the protein and a ligand or become displaced upon ligand binding. As a result, it is often of great interest to establish both the binding free energy and location of such molecules. Several approaches to predicting the location and affinity of water molecules to proteins have been proposed and utilized in the literature, although it is often unclear which method should be used under what circumstances. We report here a comparison between three such methodologies; Just Add Water Molecules (JAWS), Grand Canonical Monte Carlo (GCMC) and double-decoupling, in the hope of understanding the advantages and limitations of each method when applied to enclosed binding sites. As a result, we have adapted the JAWS scoring procedure, allowing the binding free energies of strongly bound water molecules to be

---

[*]To whom correspondence should be addressed
[†]University of Southampton
[‡]Syngenta
[¶]University of Edinburgh

calculated to a high degree of accuracy, requiring significantly less computational effort than more rigorous approaches. The combination of JAWS and GCMC offers a route to a rapid scheme capable of both locating and scoring water molecules for rational drug design.

# Introduction

The role of water molecules in protein-ligand structures and drug design has become of considerable interest in recent years.[1–5] Fuelled by the early work of Poornima and Dean in 1995,[6–8] water molecules are increasingly being included in a number of stages in the drug discovery process, primarily in virtual screening and docking approaches.[9–11] Traditionally, it has been thought that there are two major roles that water molecules play in ligand binding. The first is to stabilize a protein-ligand complex through creating a hydrogen bonding network, as seen in the binding of zanamivir to N9-neuraminidase.[12,13] The second is the ability of a water molecule to be displaced upon ligand binding, demonstrated through the development of cyclic urea analogs to target HIV-1 protease.[14] This is typically advantageous since the release of an ordered water molecule into the bulk carries an entropic gain, coupled with an enthalpic gain of strong protein-ligand interactions.[15]

More recent studies have helped to shed even more light upon the multiple, and often complex, roles which water molecules play in protein binding sites. Seemingly subtle changes in water-based hydrogen bonding networks have been shown to affect ligand-protein interaction energies,[16] highlighting the need for an accurate representation of solvation within a protein binding site. Studies by Setny, Baron and McCammon[17,18] have shown that water molecules play an active role in determining ligand-protein binding or rejection, whether this be due to mediating direct interactions, or providing an electrostatic screening effect. A recent communication by Shan *et al.* has also shown that water molecules can form a solvation shell between the protein and ligand in the binding site prior to the binding event, with the formation of this solvation shell providing a kinetic barrier to binding.[19] Bren and Janežič have developed a molecular dynamics method which

decouples the individual degrees of freedom of water molecules, allowing for the simulation of waters at any given combination of rotational, translational and vibration temperatures.[20] Such an approach allows for the identification of the delicate interactions between water molecules and their environment, and has shown promise in the calculation of the hydration free energy of water and predicting the structural properties of water under different excited states and environments.

Whilst individual, and clusters of, water molecules can directly influence ligand binding events, the structure of the water networks can also play a critical role. Experimental studies by Homans[21] have challenged the commonly held belief that protein-binding sites are fully solvated. Applying isothermal titration calorimetry to the major uninary protein system alongside Molecular Dynamics (MD), it was found that the pocket was suboptimally hydrated, with ligand binding driven by favourable protein-solute dispersion interactions rather than the expected entropic gain in displacing water from the pocket. Based upon this, it is argued that shape complementarity between ligand and protein is a viable method for drug design to target the suboptimal hydration patterns. A similar argument has been proposed by Englert *et al.* looking at the binding of phosphonamidate to thermolysin.[22]

Owing to the complex roles which water can perform in protein binding sites,[23] the reliable incorporation of water molecules into computational drug design is of critical importance. Indeed, studies have shown that incorporating water molecules into docking and virtual screening approaches can dramatically improve the predictions formed when analysing the predicted poses.[10,11] Crystallographic approaches have long been used to identify waters in and around protein binding sites, although this approach is often limited. Carugo has suggested that protein resolution is typically a limiting factor in determining the number of water molecules in a protein structure.[24] Through analysis of 873 known crystal structures, the authors indicated that, on average, a protein structure with a resolution of 2 Å had one water molecule per residue, whilst at a resolution of 1.0 Å around 1.6-1.7 waters are resolved. The same behaviour has been noted by Abel *et al.*[25] Another issue with relying on crystallographic methods to predict waters lies in the role of the crystallographer. It has been demonstrated that two independently resolved structures of the transforming

3

growth factor-$\beta2$ were found to have a different number of crystallographic waters with varying temperature factors[26] when analyzed by different crystallographers, suggesting that the addition of water molecules into a crystal structure can be problematic.

It is therefore apparent that relying upon crystallographic evidence is not always sufficient if the location of water molecules in protein binding sites is of interest. As such, various simulation methods have been developed to locate water molecules in protein binding sites. The Grand Canonical Monte Carlo method (GCMC)[27,28] is capable of locating water molecules in protein binding sites using the $\mu$VT ensemble.[29,30] During a GCMC simulation, the number of water molecules is allowed to fluctuate according to the defined chemical potential $\mu$. The chemical potential can be related directly to the binding free energy of the molecule, allowing both the location and affinity of water molecules to be found during a simulation.[31] One major drawback to the method lies in the poor acceptance rates, with water molecule insertions typically accepted with a probability of $< 1\%$ even if cavity-biasing[32,33] and configurational bias[34] schemes are utilized.

Based upon inhomogeneous fluid solvation theory (IFST),[35–38] the WaterMap method has shown promise in both locating water molecules in protein binding sites, and also assessing their free energies.[39,40] A MD simulation is performed which tracks the location of water molecules, followed by a clustering procedure which places water molecules. These sites are then subjected to the IFST analysis to estimate the enthalpy and entropy of each water site. One problem with this approach is that low density solvent regions are not accounted for, and that the clustering does not discriminate between short and long lived high-density hydration sites. In addition, sampling occluded binding sites can be problematic,[41] due to the long timescales required to allow the passage of water between the bulk and the binding site. The WaterMap methodology has been used on a number of systems, including a range of kinases,[42] the PDZ domain[43] and factor Xa.[40]

To help sample solvent inaccessible protein binding sites, the Just Add Water Molecules (JAWS) method was developed.[41] The JAWS method is based upon $\lambda$-dynamics,[44] and is capable of both locating the position of waters within a protein binding site and also providing an estimate for the binding affinity of these waters compared to the bulk. The approach works by initially simulating

so called $\theta$ water molecules which can appear and disappear across a user-defined volume typically centered on the binding site, a process referred to as JAWS stage 1. $\theta$ is an energy scaling parameter which controls the interaction energy $U_i(r)$ between $\theta$-water $i$ and the rest of the system. If the value of $\theta_i$ is 0 then the molecule acts like a ghost particle and does not interact with the system. Equally, if $\theta_i = 1$ then the molecule interacts fully with its surroundings. A simulation is performed whereby the set of $\theta$ parameters are sampled using a Metropolis Monte Carlo scheme, alongside full motion of the protein sidechains and any adjacent ligands, which results in a population density map of favourable hydration sites on the grid.

The resulting population densities are then clustered into an integer number of hydration sites, with water molecules restrained at each site using a hardwall potential. A biasing potential based upon the hydration free energy of water, together with a correction term for the incorporation of the restraint, is applied to each molecule, and a new simulation performed to observe the probability of a water molecule experiencing states with high ($\theta > 0.95$) or low ($\theta < 0.05$) values, a process termed JAWS stage 2. The cubic volume of the restraint is typically close to that of the volume which a water molecule occupies in the bulk, with the hardwall potential necessary to prevent the water molecule from leaving its site as the intermolecular interactions are removed. The ratio of these probabilities are then used to derive the binding free energy. As with the WaterMap approach, one drawback lies in the clustering of population density, whilst another lies in the adequate sampling of the $\theta < 0.05$ state for strongly bound water molecules, preventing an estimate of the binding free energy.[45] The JAWS methodology has been used by Michel *et al.* to evaluate the energetics of water displacement from three protein binding sites, finding that ligand affinity changes strongly correlate with the binding free energy of the displaced water molecules.[46] Lucarrelli *et al.* have employed the JAWS methodology to optimize the placement of water molecules in the binding site of p38a MAP kinase. They reported significant improvements in the accuracy of subsequent MC/FEP relative binding free energy predictions for 17 ligands in comparison with standard solvent equilibration protocols, highlighting the importance of accurate water placement for structure-based molecular modelling.[47]

Although the method is not capable of predicting the location of water molecules, the double-decoupling method[48] is seen as the 'gold standard' in calculating the binding free energy of water molecules in protein-ligand complexes, owing to its rigorous methodology. The method involves running two simulations, whereby a water molecule is first decoupled from a box of bulk water, with a second simulation decoupling the water molecule from the protein-ligand complex. The free energy changes from these simulations can then be used to construct a free energy cycle, from which the binding free energy can be found.

Barillari *et al.* have utilized double-decoupling to help understand the nature of water molecules which can be displaced on protein-ligand binding and those which cannot.[49] Utilising the method with replica exchange thermodynamic integration (RETI),[50,51] the study focussed on understanding whether the binding affinity of a water molecule was related to its propensity to be displaced. The paper demonstrated that, on average, water molecules which are more tightly bound are less likely to be displaced than those which are more weakly bound. With this knowledge the medicinal chemist can decide whether to target a particular water molecule for displacement, or to try and design a ligand which is capable of utilising the hydrogen bonding opportunities afforded by that water. The double-decoupling method has also been used more recently by Fadda to investigate the stability of conserved water molecules in Concanavalin A.[4]

Given the different approaches reported in the literature, it can be difficult to determine which method should be used given a particular problem. In addition, no study has ever addressed whether the methods themselves give comparable results. For example, will the binding free energy of a water molecule calculated by GCMC be comparable to that calculated by double-decoupling? In this study three freely available methods are compared; GCMC, JAWS and double-decoupling, and applied to N9 neuraminidase. This allows the various methods to be fairly assessed, and highlights the advantages and disadvantages of each method for locating and scoring water molecules. In the case of JAWS, a solution to the problem of evaluating the binding free energies of tightly bound water molecules is proposed.

# Methods

## System Setup

The crystal structure chosen for the simulations was *1nnc* (resolution = 1.80 Å).[12] Polar hydrogens were added onto the structure using the HBONDS option in whatif v7.0 (the protonation states of the ionizable residues can be found in the supplementary information),[52] with non-polar hydrogens added using LEaP. The zanamivir ligand was parameterized using the antechamber module in AMBER, with the partial charges assigned using the AM1-BCC model.[53] The partial charges obtained (found in the supplementary information) are broadly similar to those reported in a molecular dynamics study by Udommaneethanakit *et al.*,[54] which used the RESP approach to assign charges. Any differences are due to the different charge-fitting procedure used in the two studies, making direct comparisons difficult.

The complex was then minimized using the Sander module of AMBER to remove bad contacts, using an igb keyword of 1. To reduce the computational cost, only protein residues that have a heavy atom within 15 Å of zanamivir were retained. Crystallographic waters within this region were retained, except for those which were predicted in this study. The complex was solvated by a sphere of TIP4P water molecules of 23 Å radius centered upon zanamivir.[55] To prevent evaporation, a half-harmonic potential with a 1.5 kcal mol$^{-1}$ force constant was applied to all water molecules whose oxygen atom distance to the crystallographic zanamivir center of geometry was greater than 23 Å. The resulting complex was then equilibrated for 10 million moves in the NVT ensemble to remove bad contacts. During the equilibration, solvent moves were attempted with a probability of 85.7 %, protein side-chain moves with a probability of 12.9 % and solute moves with a probability of 1.4 %. For the free energy calculations, the amber99 forcefield was used,[56] with a temperature of 25 $^o$C and a residue based non-bonded cutoff of 10 Å, feathered over the last 0.5 Å. Brunsteiner has demonstrated that similar hydration free energies of water are obtained for group-based cut-off schemes as Ewald-summation schemes,[57] and consequently long-range electrostatic interactions were not modeled in this work. Any ligands used in the studies were modeled

7

using the GAFF forcefield.[58]

All simulations were performed using a modified version of the Monte Carlo code, ProtoMS 2.2.[59] For the free energy calculations, the protein backbone was constrained whilst only the sidechains were sampled. Constraining the backbone allows for consistency when comparing the different free energy methods, and the strong interactions between zanamivir, neuraminidase and the bound waters in the pocket suggest that the inclusion of backbone motion in these simulations is unlikely to induce a statistically significant shift in the reported water binding free energies. For more flexible and open binding pockets the inclusion of backbone motion is necessary to calculate reliable binding free energies.

## JAWS protocol

The JAWS stage 1 simulation was performed upon the entire binding site, encompassing a grid of size 13x9x9 $\mathring{A}^3$. 48 TIP4P JAWS waters were added to the simulation region,[55] with these molecules allowed to move freely around the grid region for one million moves whilst turned off ($\theta = 0$). Unless stated otherwise, the $\theta$ threshold applied for water molecules being classed as 'on' was 0.95 and 'off' as 0.05. Statistics were then collected on the grid region for 40 million MC moves using a grid spacing of 1 $\mathring{A}$, in line with the original JAWS study.[41] The resulting data was analyzed using AstexViewer,[60] and each grid point normalized according to the number density of the most populated grid coordinate. During the simulation, the JAWS waters were allowed to move and sample their associated $\theta_i$ value, with full sampling of the ligand angles and dihedral and bulk solvent performed. The bond angles and torsions for the side chains of residues within 10 $\mathring{A}$ of any heavy atom of zanamivir were also sampled, with the protein backbone constrained throughout the simulation. For the JAWS stage 1 simulations, solvent moves were attempted with a probability of 23 %, protein side-chain moves with a probability of 3.6 % and solute moves with a probability of 0.4 %. Variations in $\theta_i$ were attempted with a probability of 50 %, in line with the original JAWS study,[41] with translations and rotations of the JAWS waters attempted with a probability of 23 %.

The JAWS stage 1 simulation identified 7 hydration sites which were then used as starting

points for the free energy methods.

JAWS stage 2 simulations were performed by placing a 3x3x3 $Å^3$ grid over the water molecule of interest. The biasing potential, as described in Eq. (1) and Eq. (2), was added for the $\theta$-water to the potential energy, and statistics regarding the value of $\theta$ collected for 40 million MC moves.

$$V(\theta_i) = (-\Delta G_{hyd} + \Delta G_{constr}(\text{ideal, site } i))\theta_i \tag{1}$$

In Eq. (1), $\Delta G_{constr}$ is the free energy for constraining an ideal particle in a volume of $V^{constr}$ instead of the bulk, $V^0$.[48]

$$\Delta G_{constr} = -k_B T \ln \frac{V^{constr}}{V^o} \tag{2}$$

The hydration free energy of water used in the biasing potential, $\Delta G_{hyd}$, was taken to be -6.4 kcal mol$^{-1}$ in line with previous studies.[41,49] A binding free energy for the water molecule was found from the ratio of probabilities of observing a $\theta$-water at high ($\theta > 0.95$) and low ($\theta < 0.05$) $\theta$ values, using Eq. (3).

$$\Delta G_{bind}(\text{water, site } i) = -k_B T \ln \left( \frac{P(\theta_i \to 1)}{P(\theta_i \to 0)} \right) \tag{3}$$

In Eq. (3), $k_B$ is the Boltzmann constant and T is the temperature of the simulation. The $\theta$ thresholds for the high and low $\theta$ states are arbitrary, and are consistent with the original JAWS study.[41] The dependence of the calculated free energies on the choice of threshold is investigated later in this paper.

For the JAWS stage 2 simulations, solvent moves were attempted with a probability of 23 %, protein side-chain moves with a probability of 3.6 % and solute moves with a probability of 0.4 %. Variations in $\theta_i$ were attempted with a probability of 50 %, consistent with the original JAWS study,[41] with translations and rotations of the isolated JAWS water attempted with a probability of 23 %.
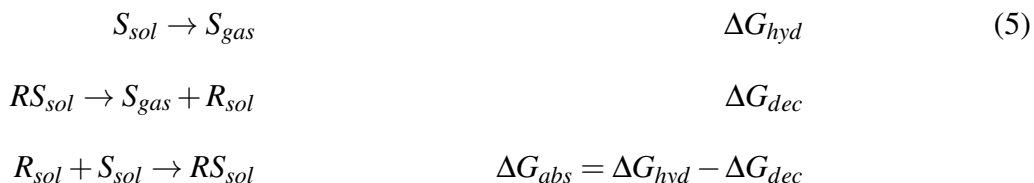
## Double-decoupling protocol

Double-decoupling[48] simulations were performed using Replica Exchange Thermodynamic Integration (RETI)[50,51] and the coordinates found from the JAWS stage 1 simulation. The binding free energy of a water molecule was found in two stages using a single topology approach: first the electrostatic terms between the water molecule and its environment were linearly perturbed to zero, followed by a gradual linear reduction in the Lennard-Jones parameters on the oxygen atom to reduce its size to zero. The water molecules were constrained by a hardwall constraint of radius 1.8 Å to allow direct comparison with the JAWS hardwall which has a similar size. The hardwall was applied to only the water in question and forbids it from leaving this spherical region. Furthermore, other water molecules, solute atoms and protein atoms were not permitted to diffuse into this excluded region. As shown in Eq. (4) the volume of this spherical hardwall, $V^{eff}$, can be calculated to be 24.43 $Å^3$, which is of similar size to the cubic 27 $Å^3$ hardwall used in JAWS stage 2 simulations.

$$\Delta G_{rest} = k_B T ln \frac{V^{eff}}{V^o} \tag{4}$$

In Eq. (4), $k_B$ is the Boltzmann constant, T is the temperature of the simulation, $V^{eff}$ the volume occupied by the hardwall and $V^0$ the standard state volume of water, 29.89 $Å^3$ at 55.56 M. From this the free energy correction term of the hardwall, $\Delta G_{rest}$ for double decoupling simulations can be found to be -0.12 kcal mol$^{-1}$.

For both the electrostatic and Lennard-Jones decoupling simulations, 16 equally spaced $\lambda$ windows were used with a value of $\Delta\lambda$ of 0.001 in conjunction with RETI. The decoupling of both the electrostatic and Lennard-Jones interactions was performed in 40 million MC steps divided into 400 blocks of 100K steps each. Data were collected and averaged over the last 30 million steps for both sets of simulations. At the end of the simulation, the computed free energies for the decoupling of the electrostatic terms of the molecule and decoupling the Lennard-Jones terms were summed, to give a value for $\Delta G_{comp}$.

Having calculated the values of $\Delta G_{comp}$ and $\Delta G_{rest}$, the binding free energy of a water molecule, $\Delta G_{abs}$, was found using Eq. (5) and Eq. (6).

$$S_{sol} \rightarrow S_{gas} \qquad\qquad \Delta G_{hyd} \qquad (5)$$

$$RS_{sol} \rightarrow S_{gas} + R_{sol} \qquad\qquad \Delta G_{dec}$$

$$R_{sol} + S_{sol} \rightarrow RS_{sol} \qquad \Delta G_{abs} = \Delta G_{hyd} - \Delta G_{dec}$$

$$\Delta G_{dec} = \Delta G_{comp} + \Delta G_{rest} - k_B T \ln \frac{\sigma_{RS}}{\sigma_R \sigma_S} + P^0 (V_R - V_{RS}) \qquad (6)$$

The third term in Eq. (6) is a symmetry related term. T is the temperature, $\sigma_{RS}$ is the symmetry number of the complex, $\sigma_R$ is the symmetry number of the protein and $\sigma_S$ is the symmetry number of water. Water has a symmetry number of 2 and, since the other two terms have a symmetry of 1, the term can be found to be -0.4 kcal mol$^{-1}$. The final term in Eq. (6) is taken to be negligible under standard pressures.[49]

For the double-decoupling simulations, solvent moves were attempted with a probability of 85.7 %, protein side-chain moves with a probability of 12.9 % and solute moves with a probability of 1.4 %. As with the JAWS simulations, only the bond angles and torsions for the side chains of residues within 10 Å of any heavy atom of zanamivir were sampled.

Error estimates from the double-decoupling simulations were obtained as the standard error across at least three independent simulations.

**GCMC protocol - Interacting Particle Method**

The GCMC simulations for the individual water molecules were performed using the interacting particle method.[61] Insertion and deletion attempts of water molecules were accepted using the following Metropolis tests.

$$P_{in} = min \left[ 1, \frac{exp(B)}{N+1} exp \left( \frac{-\Delta E}{k_B T} \right) \right] \quad (7)$$

$$P_{del} = min \left[ 1, N exp(-B) exp \left( \frac{-\Delta E}{k_B T} \right) \right] \quad (8)$$

In the above equations, N is the number of particles in the simulation and B is the Adams parameter ($B = \mu'/k_B T + \ln \bar{n}$). $\bar{n}$ is the expected number of particles in the system given the volume of the simulation region and is equal to $\bar{p}v$, where $\bar{p}$ is the number density of the particle and v the simulation volume.[30] $\mu'$ is the excess chemical potential and $\Delta E$ the change in energy between the new and old states.

Unlike the double-decoupling and JAWS simulations, no formal hardwall region is applied in a GCMC simulation. Each water molecule was looked at individually in this study to calculate its binding free energy; calculating the binding free energy of all of the waters within the same simulation is a demanding task which will result in sampling difficulties over short simulation lengths. Consequently, each water was allocated its own defined GCMC simulation region, where other water molecules were prohibited from entering the GCMC region using a hardwall but solute and protein atoms were allowed to occupy the same region. A smaller 2x2x2 $\text{Å}^3$ grid was defined around each water molecule to obtain sufficient sampling of the localized water occupancy, since it was observed that in some cases a larger volume occupied by the water molecule was filled with a solute atom, meaning that poor acceptance rates were observed.

Each B-value was simulated for 40 million MC moves, divided into 800 blocks of 50K steps each. At the end of each simulation the average water population across the entire simulation was recorded. The decoupling free energy of the water was found using Eq. (9).

$$\Delta G_{dec} = -k_B T \ln \left( \frac{[L_{sim}]}{[L_{ideal}]} \right) \quad (9)$$

In Eq. (9), $[L_{sim}]$ is found by recording the population at a particular B value. This population is converted into a localized concentration by dividing by the simulation volume, and then converting this into a molar concentration using Avogadro's number. $[L_{ideal}]$, the concentration of an ideal gas

in the simulation volume, is related to the B value of the simulation and is found using Eq. (10).

$$[L_{ideal}] = 55.56M \times exp(B - \ln \bar{n}) \tag{10}$$

In Eq. (10), $\bar{n}$ is the expected number of particles in the system given the volume of the simulation region and is equal to $\bar{p}v$, where $\bar{p}$ is the number density of the particle and v the simulation volume.[30]

Having calculated $\Delta G_{dec}$, the binding free energy of the water was found using Eq. (11).

$$\Delta G_{bind} = \Delta G_{dec} - \Delta G_{hyd} \tag{11}$$

For each water molecule, at least 5 B values were simulated to allow for a reliable estimate of the binding free energy, found as the average of the binding free energies across the range of B values.

For the GCMC simulations, solvent moves were attempted with a probability of 44 %, protein side-chain moves with a probability of 5.8 % and solute moves with a probability of 1.8 %. Insertion and deletions were each attempted with an equal probability of 2.2 %, with translations and rotations of the isolated GCMC water attempted with a probability of 44 %.

Using the interacting particle approach, an estimate of the binding free energy of a water can be found as standard error of the binding free energy across a number of independent simulations performed with different B values.

## GCMC protocol - Simulated Annealing Method

An alternative method for calculating the binding free energy of a water molecule through GCMC lies in the simulated annealing approach.[29] Rather than converting populations into localized concentrations, the populations obtained from simulating at a range of B values are instead used to make a free energy titration plot. The value of B can be related to the water binding free energy using Eq. (12).

$$\Delta G_{bind} = \Delta G_{hyd} + k_B T (B - \ln \bar{n}) \tag{12}$$

By plotting the average population of the water molecule as a function of the binding free energy, found from B, the value of $\Delta G_{bind}$ can be estimated at the equivalence point of the graph i.e. by interpolating at a population of 0.50.

Using the simulated annealing approach, an estimate of the binding free energy of a water can be found as the standard error of the population across a number of independent simulations performed at the same B value, and observing the subsequent range in the binding free energy.

The simulation protocol for the simulated annealing approach is the same as for the interacting particle method.

Unless stated otherwise, the above protocols were used for all of the subsequent studies. The JAWS and GCMC protocols were coded into the in-house Monte Carlo software, ProtoMS,[59] which was used to perform all the simulations.

## Results and discussion

### JAWS placement

A JAWS stage one simulation was performed upon N9-neuraminidase, incorporating a grid volume of 1053 $\text{Å}^3$. The native crystal structure contains a network of 5 crystallographic waters with two of these, Wat567 (subsequently referred to as Wat5) and Wat568 (subsequently referred to as Wat4), forming direct interactions with each other, the protein and the zanamivir ligand.[12] Wat4 is stabilized by a triad of nearby waters; Wat508 (Wat3), Wat543 (Wat2) and Wat575 (Wat1), which are all hydrogen-bonded to both each other and the protein. An additional structural water molecule, Wat507 (Wat6), is found in an isolated cavity and has no direct interactions with the zanamivir ligand. The JAWS simulation identified 7 possible hydration sites, shown in Figure 1, in good agreement with both the crystallographic data and the original simulations performed by
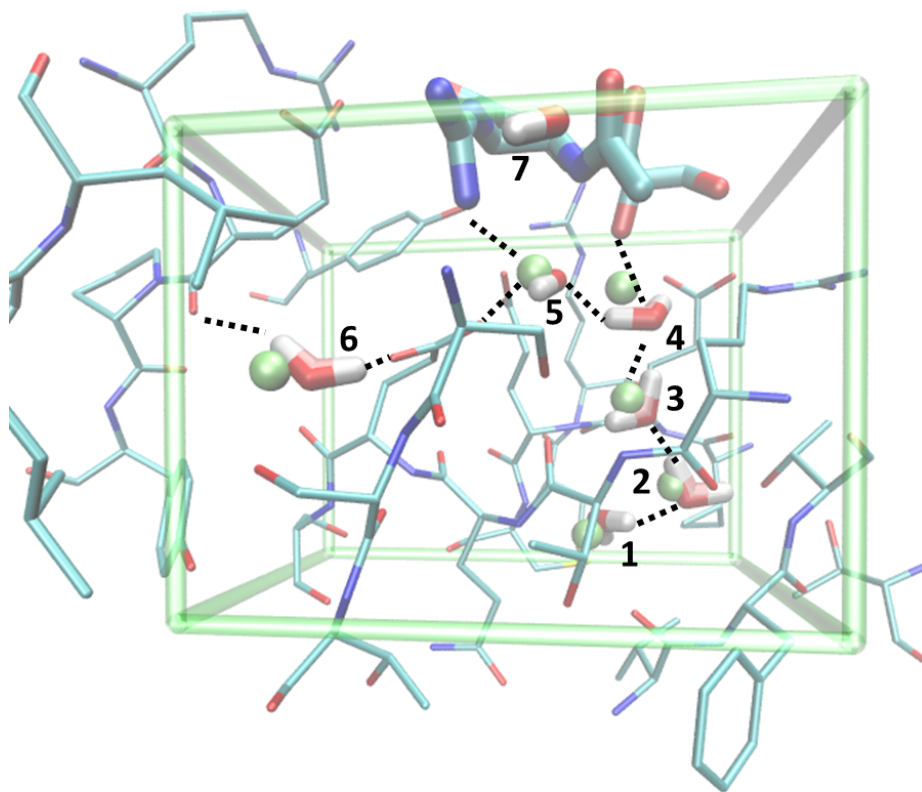
Michel *et al.*[41]



Figure 1: The 7 possible hydration sites identified by JAWS stage 1 in N9 neuraminidase. Crystallographic waters are shown in light green, with critical hydrogen bonds involving the protein, zanamivir and JAWS-waters shown as dashed lines. The zanamivir ligand is shown in liquorice.

Attempts were made to calculate the binding affinity of each of the waters using the JAWS stage 2 algorithm. It was found, however, that the majority of water molecules did not experience sufficient $\theta < 0.05$ transitions during the simulation timeframe. As a result, the binding free energies calculated by Eq. (3) were either poorly converged or unavailable. An example of this is shown in Figure 2, where the standard biasing term does not induce any $\theta < 0.05$ transitions for Wat 5 in N9 neuraminidase: the water is always on in the simulation, and hence the free energy cannot be estimated using Eq. (3).

It has been previously recognized that one the major drawbacks of the JAWS algorithm is that it cannot calculate the binding affinities of strongly bound waters.[45] To calculate the binding free
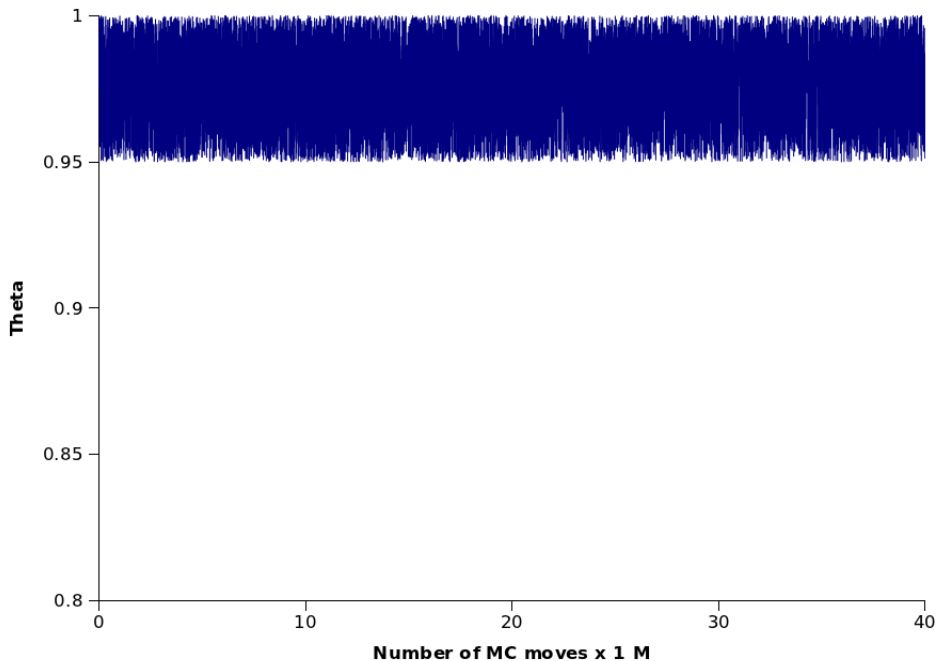
Figure 2: $\theta$ sampling for Wat5 in N9 neuraminidase using a biasing potential of 6.4 kcal mol$^{-1}$.

energies of strongly bound water molecules, we have modified the JAWS biasing term. These modifications are now discussed.

## Extension of the JAWS algorithm to calculate $\Delta G_{bind}$ for strongly bound waters

The calculation of the binding free energy of a water molecule is captured by Eq. (3). For weakly bound water molecules, the biasing potential applied in the second stage of the JAWS algorithm is sufficient to ensure that the $\theta$ water molecule can sample both the on and off states, ensuring that enough statistical sampling is performed to obtain a reliable free energy estimate. However, for strongly bound water molecules, the bias potential used in previous studies is not sufficient to induce transitions to the off state, resulting in either poor or no sampling and an unreliable estimate of the binding free energy.

One way of ensuring that sufficient sampling is performed at both end states is by changing the biasing potential applied in the second stage of the algorithm. Rather than basing this upon

the hydration free energy of water, the applied bias can be changed to one which induces sufficient transitions between the two states. The resulting free energies obtained are indicative of the new biasing potential, and hence must be corrected to take this into account, as shown in Eq. (13) where $\Delta G_{bias}$ is the value of the bias applied in the second stage of the algorithm.

$$\Delta G_{bind}(\theta_i) = -k_B T \ln\left(\frac{P(\theta_i \to 1)}{P(\theta_i \to 0)}\right) - \Delta G_{hyd} - \Delta G_{bias} \tag{13}$$

The use of a threshold to define the on and off states induces an error in the calculated values using Eq. (13), since the thresholds for the end states are taken to be 0.95 and 0.05 respectively rather than 1 and 0. Since the probability of observing a water molecule being completely turned on ($\theta$=1) or completely off ($\theta$=0) is unlikely during the simulation timeframe, an approximation is made to define the on and off states. Assuming that the water population distribution of $\theta$ across 0.95-1.00 and 0.05-0.00 is uniform, Eq. (13) can be modified to Eq. (14) and Eq. (15). In these equations, the biasing correction is approximated by the average of the applied threshold and the absolute end-points ($\theta$=1 and $\theta$=0):

$$\Delta G_{corr} = 0.975 * \Delta G_{bias} - 0.025 * \Delta G_{bias} \tag{14}$$

$$\Delta G_{bind}(\theta_i) = -k_B T \ln\left(\frac{P(\theta_i \to 1)}{P(\theta_i \to 0)}\right) - \Delta G_{hyd} - \Delta G_{corr} \tag{15}$$

The free energies obtained by this method are broadly independent of the applied bias, providing $\theta > 0.95$ and $\theta < 0.05$ are adequately sampled, as shown in Figure 3.

The free energies shown in Figure 3 were found using Eq. (15), with the associated error found as:

$$\Delta G_{error} = \Delta G_{bias} - \Delta G_{corr} \tag{16}$$

The need for changing the applied bias can be seen in Figure 4, whereby transitions between the on and off state are induced by increasing the applied bias. At a biasing potential, $\Delta G_{bias}$, of
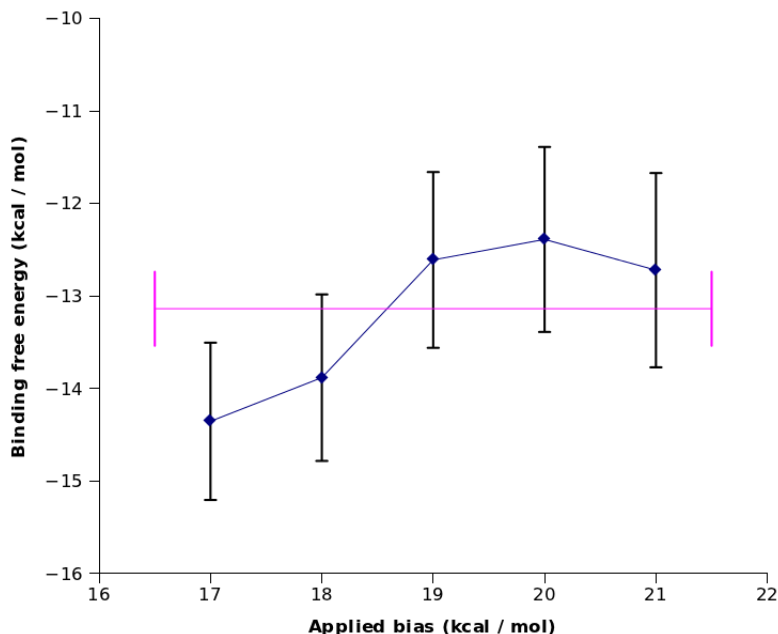
Figure 3: Effect of the applied biasing potential upon the threshold-corrected JAWS stage 2 binding free energy of Wat5 in N9 neuraminidase. The average binding free energy across the 5 simulations is shown in magenta.

10 kcal mol$^{-1}$ the water molecule experiences most of the simulation time in the on state, with no transitions to the off state being observed, meaning that a reliable free energy estimate cannot be obtained. However, upon using a $\Delta G_{bias}$ of 17 kcal mol$^{-1}$, the water molecule can sample both end states, allowing a reliable free energy estimate to be obtained.

There is also an error associated with the choice of $\theta$ threshold. For example, the threshold for an 'on' state could arbitrarily be either $\theta > 0.95$ or $\theta > 0.98$. To estimate the associated error, the binding free energy of Wat5 in N9-neuraminidase was calculated using different thresholds for both the on and off states. The calculated free energies can be found in Figure 5, with the error associated with the choice of $\theta$, found as the standard error across the 5 measurements, estimated to be $\pm$ 0.40 kcal mol$^{-1}$. It can be seen that there seems to be no systematic dependence on the theta threshold chosen.

The overall error for the JAWS simulations is thus the combination of the error in the $\theta$ threshold and Eq. (16).
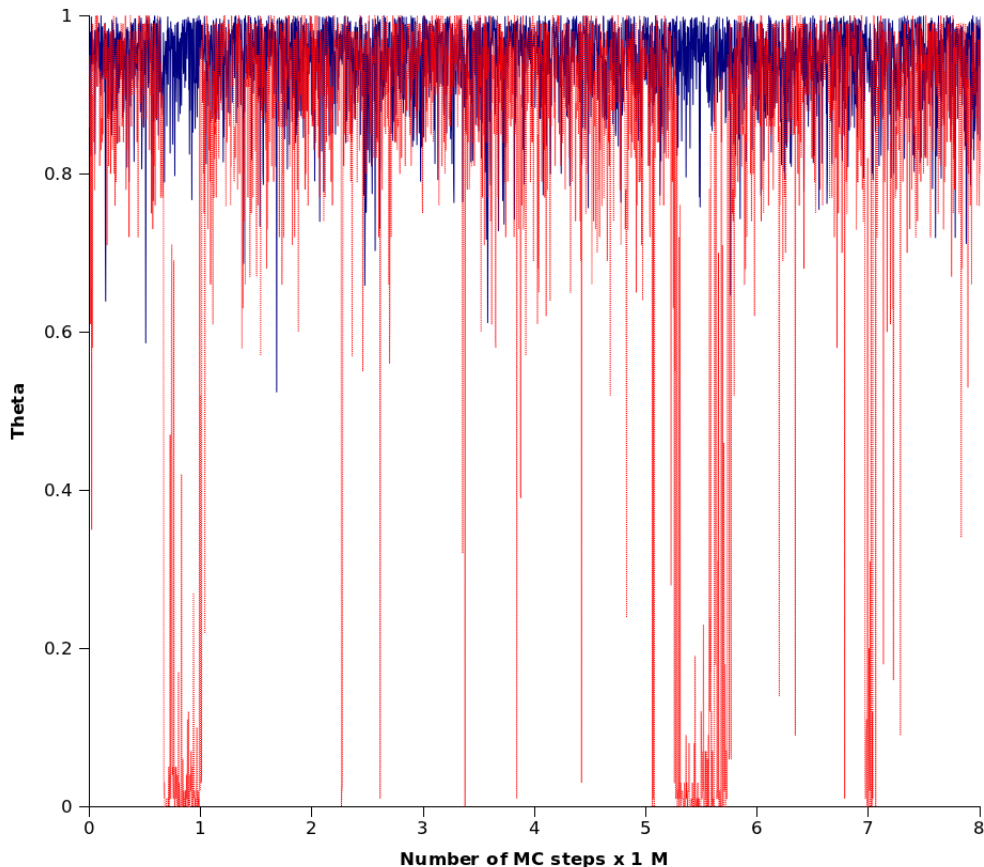
18

Figure 4: $\theta$ sampling as a function of the applied bias potential for Wat3 in N9 neuraminidase. Shown in blue is the sampling at a bias of 10 kcal mol$^{-1}$, whilst the sampling at 17 kcal mol$^{-1}$ is shown in red.

**Choice of biasing potential**

Since the binding free energy is broadly independent of the applied bias, provided the on and off states are adequately sampled, a JAWS stage 2 simulation needs to be run at only one value of the bias to extract the binding free energy. The ideal bias should induce an equal number of on and off states, meaning that the binding free energy becomes the difference between the standard hydration free energy of water and the applied bias, as in Eq. (15). An equal number of on and off states means that the water is, on average, present 50 % of the time and should give the most reliable estimate of the binding free energy. To achieve this, a simple optimization procedure of the biasing potential has been applied.
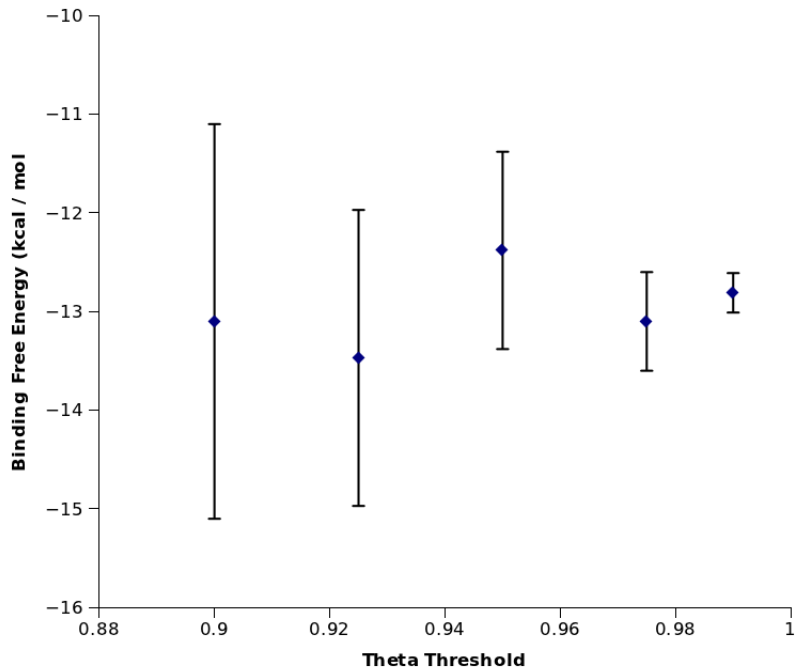
Figure 5: JAWS stage 2 binding free energy of Wat5 in N9 neuraminidase as a function of the $\theta$ threshold applied. Simulations were performed using a biasing potential of 20 kcal mol$^{-1}$, with errors found using Eq. (16).

A short JAWS stage 2 simulation, typically one million MC moves, is performed and an estimation of the binding free energy found. The biasing potential is then optimized to obtain a value of the bias which yields an equal number of on and off states. As an example, the process is illustrated for hydration site Wat3 in Table 1.

Table 1: Optimization procedure for Wat3 in N9 neuraminidase

| Iteration | $\Delta G_{bias}$ (kcal mol$^{-1}$) | Ln(On/Off) | Iteration | $\Delta G_{bias}$ (kcal mol$^{-1}$) | Ln(On/Off) |
|---|---|---|---|---|---|
| 1 | 15.0 | 12 | 6 | 17.5 | -1.0 |
| 2 | 15.5 | 10 | 7 | 17.0 | 11 |
| 3 | 16.0 | 11 | 8 | 17.5 | 3.0 |
| 4 | 16.5 | 4.0 | 9 | 18.0 | -2.0 |
| 5 | 17.0 | 3.0 | 10 | 17.5 | 0.0 |

In Table 1, the iterative optimization is trying to obtain a value of $\Delta G_{bias}$ which induces an equal number of on and states. As shown in Eq. (15), an equal number of on and off states should

result in the log term approaching zero. At the end of each simulation, the log ratio of on and off states is calculated. If the value is positive, suggesting more on states than off, then the value of $\Delta G_{bias}$ is increased by 0.5 kcal mol$^{-1}$ for the next iteration. If the value is negative, suggesting more off states than on, then the value of $\Delta G_{bias}$ is decreased by 0.5 kcal mol$^{-1}$. At the end of the process, typically incorporating 10 iterations, the value of $\Delta G_{bias}$ which gives a log term of zero is chosen for the main simulation. In this example, the value of $\Delta G_{bias}$ was taken to be 17.5 kcal mol$^{-1}$. The variance in the calculated log terms shows that this an approximate method and that longer runs are required to achieve convergence, but they provide guidance to the most suitable value of $\Delta G_{bias}$.

# Binding Free Energy Calculations

## Comparison of JAWS and RETI

Using the new modifications, a JAWS stage 2 simulation was performed upon each hydration site as identified in Figure 1, with each site also studied using double-decoupling. The binding free energy for JAWS stage two simulations was found using Eq. (15). The free energy comparison between the two methods can be seen in Figure 6.

Figure 6 clearly demonstrates that the two methods give excellent agreement with each other. Binding affinities of both strongly and weakly bound water molecules are calculated up by the two methods, showing that the modification to the JAWS algorithm has been successful in predicting the binding free energy of water molecules which previously were incalculable. The biasing potentials used to calculate the JAWS free energies in Figure 6 are given in Table 2.

The convergence for the JAWS stage 2 simulations as a function of the simulation length are given in Figure 7. It can be seen that all of the simulations are well converged after ca. 10 million MC moves, demonstrating further that the biasing potentials used yield precise results.
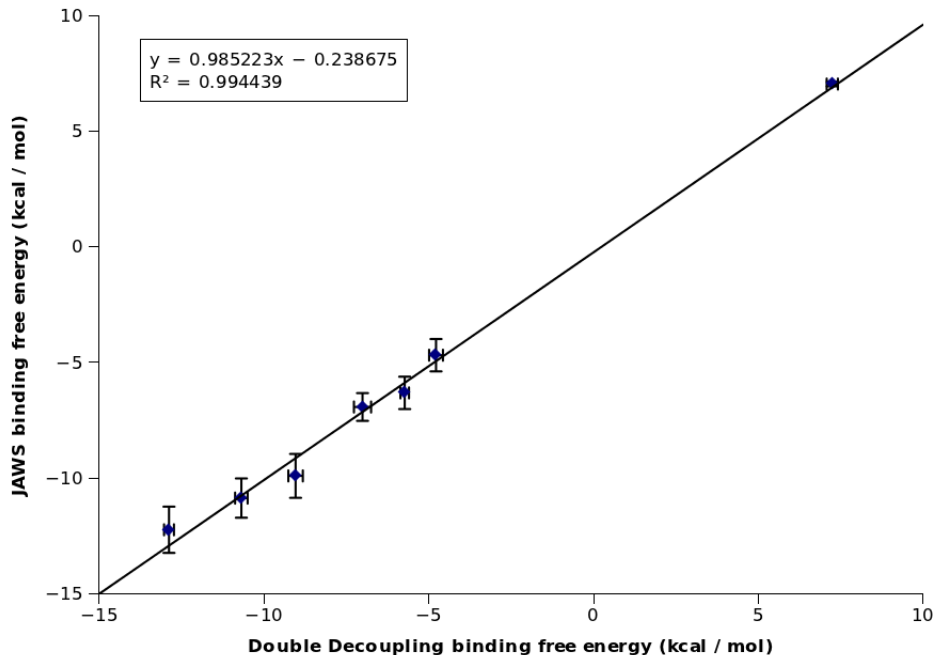
21

Figure 6: Binding free energies for the 7 hydration sites in N9 neuraminidase, found using JAWS stage 2 and RETI double-decoupling

Table 2: Binding free energies for the 7 water molecules in N9 neuraminidase, found using JAWS stage 2 and RETI double-decoupling. Statistical errors are shown in parentheses. The individual free energy components for these double-decoupling simulations can be found in the supplementary information.

| Water Molecule | JAWS Bias (kcal mol$^{-1}$) | JAWS $\Delta G_{bind}$ (kcal mol$^{-1}$) | RETI $\Delta G_{bind}$ (kcal mol$^{-1}$) |
|---|---|---|---|
| 1 | 14 | -4.7 (0.7) | -4.5 (0.3) |
| 2 | 12 | -6.9 (0.6) | -7.6 (0.3) |
| 3 | 17.5 | -10.9 (0.9) | -10.9 (0.1) |
| 4 | 14 | -6.3 (0.7) | -6.1 (0.1) |
| 5 | 20 | -12.4 (1.0) | -13.2 (0.1) |
| 6 | 19 | -9.9 (1.0) | -9.8 (0.1) |
| 7 | 2 | 7.0 (0.1) | 7.1 (0.3) |

## Comparison of GCMC and RETI

Using the hydration sites identified by the JAWS stage 1 simulations, the free energy of binding of each site was calculated using the interacting particle method of Clark *et al.*[61] Eq. (9) yields the decoupling energy of the water molecule from the protein, and can be corrected with the hydration
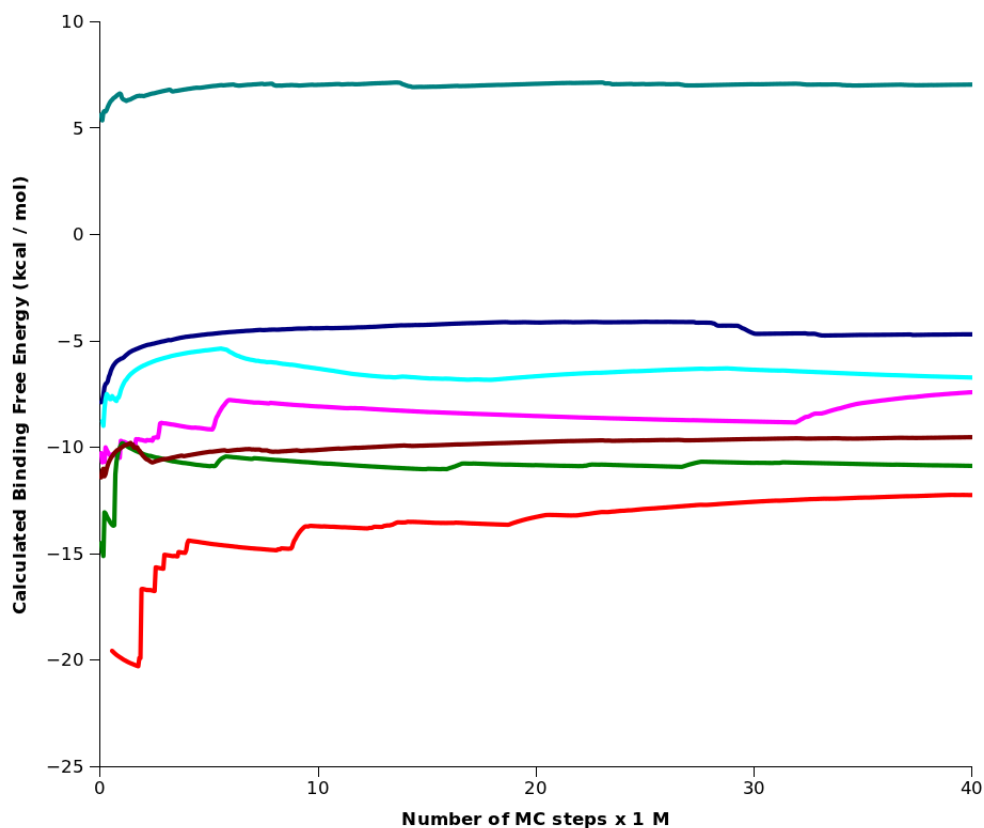
Figure 7: Convergence of the binding free energies in the JAWS stage 2 simulations. Key: Wat 1 (Dark blue), Wat 2 (Magenta), Wat 3 (Dark green), Wat 4 (Cyan), Wat 5 (Red), Wat 6 (Maroon), Wat 7 (Turquoise).

free energy of water to arrive at a binding free energy using Eq. (11). Figure 8 shows the predicted binding affinity of the 7 molecules using the two methods and shows an excellent correlation.

The reported GCMC binding free energies were calculated as the average of the binding free energy across a range of B values. An example of this for Wat7, the weakest binder in the series, can be seen in Table 3. The table shows that the calculated binding free energy is consistent and approximately 7 kcal mol$^{-1}$ once the average population recorded throughout the simulation, used to calculate [L$_{sim}$], drops below 0.60, corresponding to sufficient sampling of the water populations. This behaviour is consistent with that demonstrated by Clark *et al.* in the calculation of benzene-T4 lysozyme binding free energies.[61,62] Running a GCMC simulation at different B values can be viewed as a free energy titration, and thus the binding free energy of the water should correspond
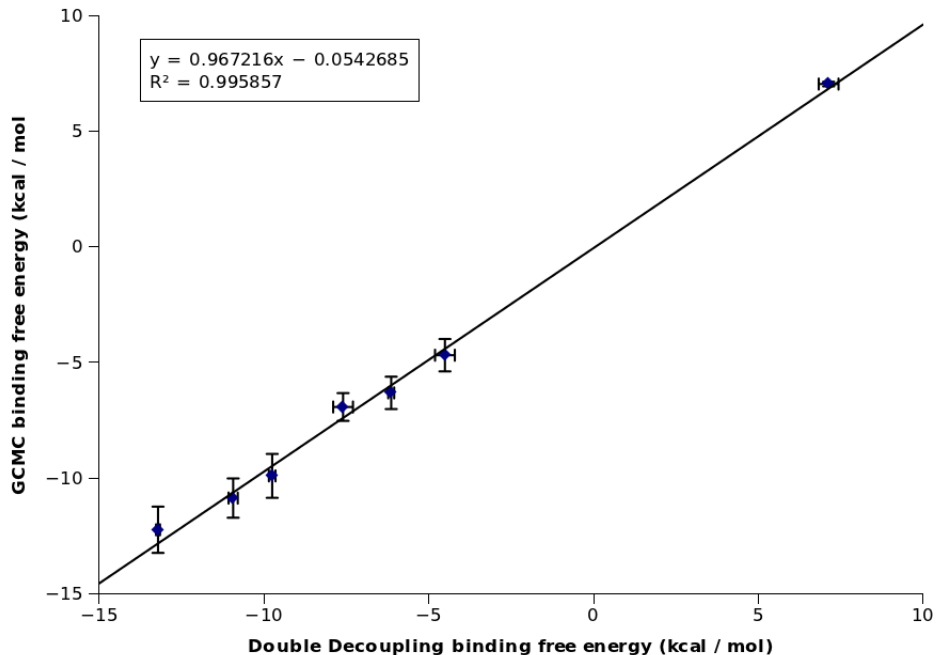
Figure 8: Binding free energies for the 7 hydration sites in N9 neuraminidase, found using GCMC and RETI double-decoupling

to the equivalence point of the titration, i.e. 0.50. The maximum occupancy for the water site is unity, meaning statistically significant occupancies of less than one need to be obtained to calculate a reliable estimate of the binding free energy using Eq. (9).

Table 3: Calculated free energies for Wat7 in N9 neuraminidase, found at different B levels

| B | Average number of molecules | $[L_{sim}]$ (M) | $[L_{ref}]$ (M) | $\Delta G_{dec}$ (kcal mol$^{-1}$) | $\Delta G_{bind}$ (kcal mol$^{-1}$) |
|---|---|---|---|---|---|
| 4 | 0.97 | 200 | 11000 | 2.39 | 8.79 |
| 3 | 0.93 | 190 | 4100 | 1.82 | 8.22 |
| 2 | 0.73 | 150 | 1500 | 1.37 | 7.77 |
| 1 | 0.57 | 120 | 570 | 0.93 | 7.33 |
| 0 | 0.37 | 80 | 210 | 0.59 | 6.99 |
| -1 | 0.15 | 31 | 80 | 0.54 | 6.94 |
| -2 | 0.04 | 8 | 28 | 0.76 | 7.16 |

Using the data for B values less than 2 in Table 3, the binding free energy can be estimated as $7.24 \pm 0.17$ kcal mol$^{-1}$.

One major potential drawback associated with the GCMC method lies in the acceptance rate of insertion and deletion moves. For an insertion to be accepted it is important that the orientation of the water molecule is correct, since otherwise it is likely that the intermolecular interactions between the water and its environment will be unfavorable.[34] As a result, insertion rates as low as 0.1 % are seen in GCMC simulations, which in turn leads to poor sampling. It is this poor sampling which could potentially lead to an increase in the uncertainty in the free energies derived from GCMC simulations compared to both double-decoupling and JAWS, but no evidence of this has been seen in this study. Although an insertion rate of 0.1 % at first sight is worryingly low, this corresponds to successfully sampling 800 insertion events over the simulation, which appear sufficient to converge the free energy estimate.

Since the interacting particle method can generate populations as a function of B, this information can also be used to derive free energies via the simulated annealing approach.[29] As previously described, the method can be considered to be analogous to a chemical titration whereby the decoupling free energy of the water molecule is the equivalence point at which the average water population is 0.50. The data from Table 3 has been used to generate such a titration profile, seen in Figure 9.

Figure 9 shows that the estimated binding free energy of Wat7 is $7.4 \pm 0.3$ kcal mol$^{-1}$ at the 0.50 equivalence point, in good agreement with the value calculated by the interacting particle method. Since either method can be used to derive the same result to within error, the question arises as to which of the GCMC methods is to be preferred when calculating binding free energies. Whilst the simulated annealing approach gives information regarding the behaviour of the system as a function of B, the interacting particle approach is significantly faster since it only requires the simulation to be performed at one value of B. However, the optimal B value to choose is not always known *a priori*, meaning that it can require several different simulations to identify a B value that yields precise binding free energies. It therefore appears to be more advantageous to use the simulated annealing approach, since this gives information on how the population changes as a function of the applied B value and can also be used to calculate the free energy. In addition the
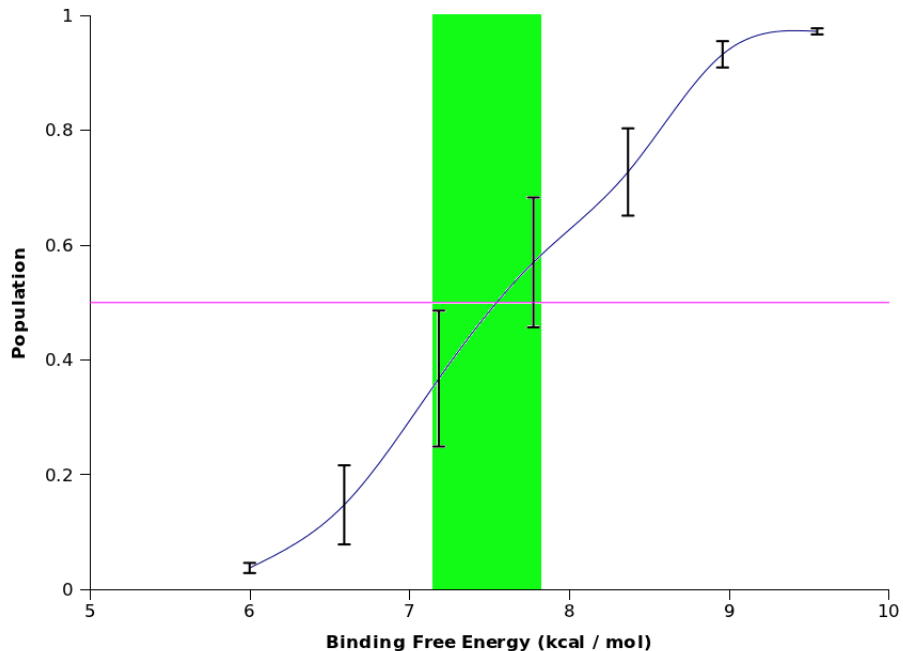
Figure 9: Free energy titration plot for Wat7 in N9 neuraminidase, found using the GCMC simulated annealing approach. The error bars were found as the standard error across two independent simulations. The shaded green area shows the error in the interpolated value, found as the region incorporating the minimum and maximum equivalence points

method is potentially suitable to be used on larger systems, whereby multiple waters can be looked at once, although longer simulation timescales will be required to derive converged binding free energies.

## Method Comparison

The N9-neuraminidase system was chosen as a test case for the different free energy methodologies since there are a large number of different studies which have used the system in free energy calculations.[41,49] The fact that all three free-energy methods give near identical results, to within statistical error, is clearly encouraging and lends itself to the question of which method is best suited to a particular problem. Whilst RETI double-decoupling is the most rigorous method of the three it is also the most computationally expensive and requires knowledge of the water binding sites. A typical simulation requires in excess of 300 CPU hours on 16 2.6 GHz processors, whilst

both GCMC and JAWS require an order of magnitude less computational time and can be run on a single processor. As a result it is suggested that double decoupling is used in cases where precise free energies are required.

One drawback of the double-decoupling approach is that it requires prior knowledge of the water binding positions; something which is found dynamically in both JAWS and GCMC. As such, if systems without clearly defined crystallographic waters are studied then either JAWS, GCMC or both methods should be employed to identify potential hydration sites. JAWS has already been employed in free energy studies, whereby changes in hydration as a function of ligand perturbations evaluated.[46,47] With the extension of the JAWS biasing potential, both JAWS and GCMC can calculate the binding free energy for both strongly and weakly bound waters during a single simulation and take similar simulation times. However, one possible advantage in the JAWS approach is that once the optimal biasing potential is found no further simulations need to be run whilst the GCMC approach requires several simulations at different B-values to arrive at the binding free energy. This in itself, however, highlights one of the advantages in the GCMC approach; that it can give information on the binding of molecules as a function of the chemical potential.

One potential problem with the JAWS approach to calculate binding free energies is how to deal with the intermediate $\theta$ states. Since only the end points are considered when calculating the binding free energy, it is unclear whether or not the intermediate data should be considered. The probability that a $\theta$-water adopts an ill-defined intermediate state can be significant, as indicated in Figure 10.

Based upon the excellent agreement between JAWS, GCMC and double-decoupling, it seems that ignoring the intermediate states is acceptable. This is, however, a waste of potentially useful data, although it is unclear how the data could be analyzed.
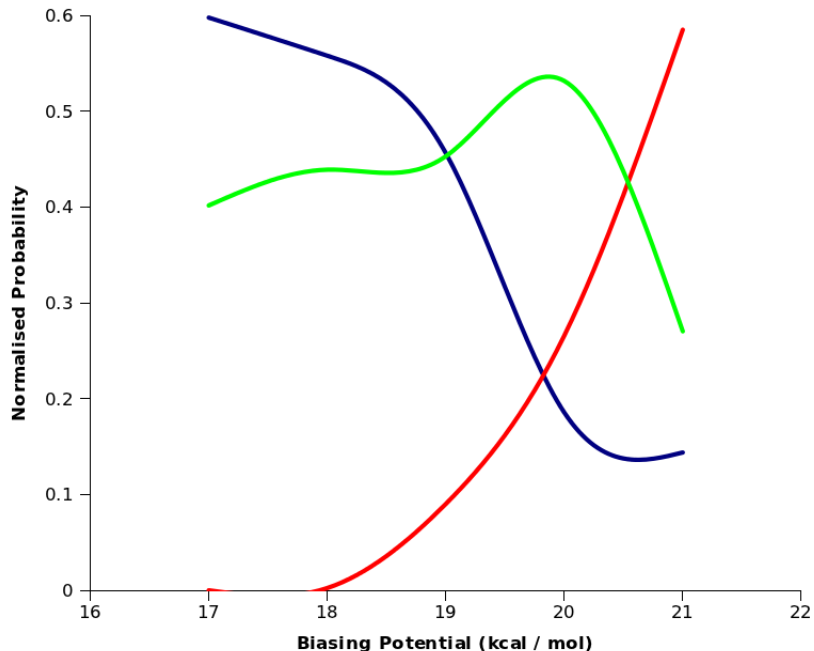
Figure 10: Normalised sampling probabilities for the on, off and intermediate states as a function of the applied biasing potential for Wat 5 in N9 neuraminidase. The on state (Dark blue) was defined as $\theta > 0.95$, the off state (Red) defined as $\theta < 0.05$ and the intermediate state (Green) accounted for the remaining $\theta$ values

## Conclusions

The application of JAWS, GCMC and double-decoupling Monte Carlo has demonstrated that all three methods yield water binding free energies which are consistent with each other when applied to a single, isolated water site in N9-neuraminidase. The double-decoupling approach is by far the slowest of the three methods, since it requires free energy changes for intermediate states between $\lambda$=0 and $\lambda$=1 to be simulated. The method is the most rigorous of the three studied, and we propose that it is used to score water molecules when precise free energies are required.

The extended JAWS stage 2 and GCMC simulations both require, on average, an order of magnitude less computational resources than the double-decoupling method, meaning they should be preferred for rapid estimation of the binding free energy of a single water molecule. Despite this, the JAWS methodology typically requires fewer simulations to be run than GCMC, and on this basis we propose that JAWS should be preferred in the calculation of water binding free energies.

Although both JAWS and GCMC are capable of predicting the location of water molecules in protein cavities, the resulting data from a JAWS simulation can be difficult to interpret. For well defined sites, such as the N9 neuraminidase system, presented in this paper, the population maps are clear and easy to interpret. However, if the population density contours overlap it can prove difficult to resolve the population densities into individual hydration sites. The population density does not reveal whether the binding of two or more water molecules are correlated, or whether a single water molecule hops between two adjacent sites. In comparison the GCMC method is capable of clearly identifying hydration sites, since a representative simulation snapshot with the desired number of molecules can be easily extracted. Crucially, none of the GCMC snapshots contain ill-defined intermediate $\theta$-water states and each snapshot is a valid structural representation for the particular binding free energy at which the simulation was run. As such cooperative binding between water molecules can also be identified, something which is harder to infer from a JAWS population density analysis.

In cases where the JAWS density map does not indicate clear, isolated water binding sites, we propose that GCMC should be used to identify potential hydration sites. Since the results from a GCMC simulation are dependent upon the chemical potential used, the chemical potential can be tuned to observe hydration patterns as a function of the binding free energy. The resulting hydration sites could then be scored in a JAWS stage 2 simulation if desired.

# Acknowledgments

## Supporting Information Available

Protonation state analysis for the 1NNC protein structure, and partial charges used to model the zanamivir ligand. Free energy decomposition for the decoupling of water molecules bound to 1NNC and zanamivir using RETI double-decoupling.

This material is available free of charge via the Internet at `http://pubs.acs.org/`.

## References

(1) Cappel, D.; Wahlstrom, R.; Brenk, R.; Sotriffer, C. A. Probing the dynamic nature of water molecules and their influences on ligand binding in a model binding site. *J. Chem. Inf. Model.* **2011**, *51*, 2581–2594.

(2) Wallnoefer, H. G.; Liedl, K. R.; Fox, T. A GRID-derived water network stabilizes molecular dynamics computer simulations of a protease. *J. Chem. Inf. Model.* **2011**, *51*, 2860–2867.

(3) Li, Y.; Sutch, B. T.; Bui, H.-H.; Gallaher, T. K.; Haworth, I. S. Modeling of the water network at protein-RNA interfaces. *J. Chem. Inf. Model.* **2011**, *51*, 1347–1352.

(4) Fadda, E.; Woods, R. J. On the role of water models in quantifying the binding free energy of highly conserved water molecules in proteins: The case of Concanavalin A. *J. Chem. Theory Comput.* **2011**, *7*, 3391–3398.

(5) Barillari, C.; Duncan, A. L.; Westwood, I. M.; Blagg, J.; Van Montfort, R. L. M. Analysis of water patterns in protein kinase binding sites. *Proteins: Struct. Funct. Bioinf.* **2011**, *79*, 2109–2121.

(6) Poornima, C. S.; Dean, P. M. Hydration in drug design. 1. Multiple hydrogen-bonding features of water molecules in mediating protein-ligand interactions. *J. Comput.-Aided Mol. Des.* **1995**, *9*, 500–512.

(7) Poornima, C. S.; Dean, P. M. Hydration in drug design. 2. Influence of local site surface shape on water binding. *J. Comput.-Aided Mol. Des.* **1995**, *9*, 513–520.

(8) Poornima, C. S.; Dean, P. M. Hydration in drug design. 3. Conserved water molecules at the ligand-binding sites of homologous proteins. *J. Comput.-Aided Mol. Des.* **1995**, *9*, 521–531.

(9) Cheng, T.; Li, Q.; Zhou, Z.; Wang, Y.; Bryant, S. H. Structure-based virtual screening for drug discovery: a problem-centric review. *Amer. Asso. Pharma. Sci.* **2012**, *14*, 131–141.

(10) Wong, S. E.; Lightstone, F. C. Accounting for water molecules in drug design. *Exp. Opin. Drug Discovery* **2011**, *6*, 65–74.

(11) Huggins, D. J.; Tidor, B. Systematic placement of structural water molecules for improved scoring of protein-ligand interactions. *Prot. Eng. Des. Sel.* **2011**, *24*, 777–789.

(12) Varghese, J. N.; Epa, V. C.; Colman, P. M. Three-dimensional structure of the complex of 4-guanidino-Neu5Ac2en and influenza virus neuraminidase. *Prot. Sci.* **1995**, *4*, 1081–1087.

(13) Smith, B. J.; Colman, P. M.; Von Itzstein, M.; Danylec, B.; Varghese, J. N. Analysis of inhibitor binding in influenza virus neuraminidase. *Prot. Sci.* **2001**, *10*, 689–696.

(14) De Lucca, G. V.; Erickson-Viitanen, S.; Lam, P. Y. S. Cyclic HIV protease inhibitors capable of displacing the active site structural water molecule. *Drug Discovery Today* **1997**, *2*, 6–18.

(15) De Beer, S. B. A.; Vermeulen, N. P. E.; Oostenbrink, C. The role of water molecules in computational drug design. *Curr. Top. Med. Chem.* **2010**, *10*, 55–66.

(16) Hummer, G. Molecular binding: Under water's influence. *Nat. Chem.* **2010**, *2*, 906–907.

(17) Setny, P.; Baron, R.; McCammon, J. A. How can hydrophobic association be enthalpy driven? *J. Chem. Theory Comput.* **2010**, *6*, 2866–2871.

(18) Baron, R.; Setny, P.; McCammon, J. A. Water in cavity-ligand recognition. *J. Am. Chem. Soc.* **2010**, *132*, 12091–12097.

(19) Shan, Y.; Kim, E. T.; Eastwood, M. P.; Dror, R. O.; Seeliger, M. A.; Shaw, D. E. How does a drug molecule find its target binding site? *J. Am. Chem. Soc.* **2011**, *133*, 9181–9183.

(20) Bren, U.; Janežič, D. Individual degrees of freedom and the solvation properties of water. *J. Chem. Phys.* **2012**, *137*, 024108–024119.

(21) Homans, S. W. Water, water everywhere–except where it matters? *Drug Discovery Today* **2007**, *12*, 534–539.

(22) Englert, L.; Biela, A.; Zayed, M.; Heine, A.; Hangauer, D.; Klebe, G. Displacement of disordered water molecules from hydrophobic pocket creates enthalpic signature: binding of phosphonamidate to the S4-pocket of thermolysin. *Biochim. Biophys. Acta* **2010**, *1800*, 1192–1202.

(23) Ball, P. Biophysics: More than a bystander. *Nature* **2011**, *478*, 467–468.

(24) Carugo, O.; Bordo, D. How many water molecules can be detected by protein crystallography? *Acta Cryst. Sec. D Biol. Cryst.* **1999**, *55*, 479–483.

(25) Abel, R.; Salam, N. K.; Shelley, J.; Farid, R.; Friesner, R. A.; Sherman, W. Contribution of explicit solvent effects to the binding affinity of small-molecule inhibitors in blood coagulation factor serine proteases. *Chem. Med. Chem.* **2011**, *6*, 1049–1066.

(26) Davis, A. M.; Teague, S. J.; Kleywegt, G. J. Application and limitations of X-ray crystallographic data in structure-based ligand and drug design. *Angew. Chem. Int. Ed.* **2003**, *42*, 2718–2736.

(27) Adams, D. J. Chemical potential of hard-sphere fluids by Monte-Carlo methods. *Mol. Phys.* **1974**, *28*, 1241–1252.

(28) Adams, D. J. Grand canonical ensemble Monte Carlo for a Lennard-Jones fluid. *Mol. Phys.* **1975**, *29*, 307–311.

(29) Guarnieri, F.; Mezei, M. Simulated annealing of chemical potential: a general procedure for locating bound waters. Application to the study of the differential hydration propensies of the major and minor grooves of DNA. *J. Am. Chem. Soc.* **1996**, *118*, 8493–8494.

(30) Woo, H.-J.; Dinner, A. R.; Roux, B. Grand canonical Monte Carlo simulations of water in protein environments. *J. Chem. Phys.* **2004**, *121*, 6392–6400.

(31) Bortolato, A.; Tehan, B.; Bodnarchuk, M. S.; Essex, J. W.; Mason, J. Water network perturbation in ligand binding: Adenosine A(2A) antagonists as a case study. *J. Chem. Inf. Model.* **2013**, *53*, 1700–1713.

(32) Mezei, M. A cavity-biased (T,V,mu) Monte-Carlo method for the computer-simulation of fluids. *Mol. Phys.* **1980**, *40*, 901–906.

(33) Mezei, M. Grand-canonical ensemble Monte Carlo study of dense liquid Lennard-Jones, soft spheres and water. *Mol. Phys.* **1987**, *61*, 565–582.

(34) Shelley, J. C.; Patey, G. N. A configuration bias Monte-Carlo method for water. *J. Chem. Phys.* **1995**, *102*, 7656–7664.

(35) Lazaridis, T. Inhomogeneous fluid approach to solvation thermodynamics. 1. Theory. *J. Phys. Chem. B* **1998**, *102*, 3531–3541.

(36) Lazaridis, T. Inhomogeneous fluid approach to solvation thermodynamics. 2. Applications to simple fluids. *J. Phys. Chem. B* **1998**, *102*, 3542–3550.

(37) Nguyen, C. N.; Young, T. K.; Gilson, M. K. Grid inhomogeneous solvation theory: Hydration structure and thermodynamics of the miniature receptor cucurbit[7]uril. *J. Chem. Phys.* **2012**, *137*, 044101.

(38) Prabhu, E. P.; MacKerell, A. D.; Rapid estimation of hydration thermodynamics of macromolecular regions. *J. Chem. Phys.* **2013**, *139*, 055105.

(39) Young, T.; Abel, R.; Kim, B.; Berne, B. J.; Friesner, R. A. Motifs for molecular recognition exploiting hydrophobic enclosure in protein-ligand binding. *Proc. Natl. Acad. Sci. USA* **2007**, *104*, 808–813.

(40) Abel, R.; Young, T.; Farid, R.; Berne, B. J.; Friesner, R. A. Role of the active-site solvent in the thermodynamics of factor Xa ligand binding. *J. Am. Chem. Soc.* **2008**, *130*, 2817–2831.

(41) Michel, J.; Tirado-Rives, J.; Jorgensen, W. L. Prediction of the water content in protein binding sites. *J. Phys. Chem. B* **2009**, *113*, 13337–13346.

(42) Robinson, D. D.; Sherman, W.; Farid, R. Understanding kinase selectivity through energetic analysis of binding site waters. *Chem. Med. Chem.* **2010**, *5*, 618–627.

(43) Beuming, T.; Farid, R.; Sherman, W. High-energy water sites determine peptide binding affinity and specificity of PDZ domains. *Prot. Sci.* **2009**, *18*, 1609–1619.

(44) Kong, X.; Brooks III, C. L. Lambda-Dynamics: A new approach to free energy calculations. *J. Chem. Phys.* **1996**, *105*, 2414–2423.

(45) Ross, G. A.; Morris, G. M.; Biggin, P. C. Rapid and accurate prediction and scoring of water molecules in protein binding sites. *PLoS ONE* **2012**, *7*, e32036.

(46) Michel, J.; Tirado-Rives, J.; Jorgensen, W. L. Energetics of displacing water molecules from protein binding sites: consequences for ligand optimization. *J. Am. Chem. Soc.* **2009**, *131*, 15403–15411.

(47) Luccarelli, J.; Michel, J.; Tirado-Rives, J.; Jorgensen, W. L. Effects of water placement on predictions of binding affinities for p38$\alpha$ MAP kinase inhibitors. *J. Chem. Theory Comput.* **2010**, *6*, 3850–3856.

(48) Gilson, M. K.; Given, J. A.; Bush, B. L.; McCammon, J. A. The statistical-thermodynamic basis for computation of binding affinities: a critical review. *Biophys. J.* **1997**, *72*, 1047–1069.

(49) Barillari, C.; Taylor, J.; Viner, R.; Essex, J. W. Classification of water molecules in protein binding sites. *J. Am. Chem. Soc.* **2007**, *129*, 2577–2587.

(50) Woods, C. J.; Essex, J. W.; King, M. A. Enhanced configurational sampling in binding free-energy calculations. *J. Phys. Chem. B* **2003**, *107*, 13711–13718.

(51) Woods, C. J.; Essex, J. W.; King, M. A. The development of replica-exchange-based free-energy methods. *J. Phys. Chem. B* **2003**, *107*, 13703–13710.

(52) Vriend, G. WHAT IF: a molecular modeling and drug design program. *J. Mol. Graph.* **1990**, *8*, 52–56, 29.

(53) Jakalian, A.; Bush, B. L.; Jack, D. B.; Bayly, C. I. Fast, efficient generation of high-quality atomic charges. AM1-BCC model: I. Method. *J. Comput. Chem.* **2000**, *21*, 132–146.

(54) Udommaneethanakit, T.; Rungrotmongkol, T.; Bren, U.; Frecer, V.; Stanislav, M. Dynamic behavior of avian influenza A virus neuraminidase subtype H5N1 in complex with os-eltamivir, zanamivir, peramivir, and their phosphonate analogues. *J. Chem. Inf. Model.* **2009**, *49*, 2323–2332.

(55) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.* **1983**, *79*, 926–935.

(56) Wang, J.; Cieplak, P.; Kollman, P. A. How well does a restrained electrostatic potential (RESP) model perform in calculating conformational energies of organic and biological molecules? *J. Comput. Chem.* **2000**, *21*, 1049–1074.

(57) Brunsteiner, M.; Boresch, S. Influence of the treatment of electrostatic interactions on the results of free energy calculations of dipolar systems. *J. Chem. Phys.* **2000**, *112*, 6953–6955.

(58) Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. Development and testing of a general amber force field. *J. Comput. Chem.* **2004**, *25*, 1157–1174.

(59) Woods, C.; Michel, J. *ProtoMS2.2*. 2007.

(60) Hartshorn, M. J. AstexViewer: a visualisation aid for structure-based drug design. *J. Comput.-Aided Mol. Des.* **2002**, *16*, 871–881.

(61) Clark, M.; Meshkat, S.; Wiseman, J. S. Grand canonical free-energy calculations of protein-ligand binding. *J. Chem. Inf. Model.* **2009**, *49*, 934–943.

(62) Clark, M.; Guarnieri, F.; Shkurko, I.; Wiseman, J. Grand canonical Monte Carlo simulation of ligand-protein binding. *J. Chem. Inf. Model.* **2006**, *46*, 231–242.

# For Table of Contents Only

**Strategies to calculate water binding free energies in protein-ligand complexes.**

*Michael S. Bodnarchuk, Russell Viner, Julien Michel, and Jonathan W. Essex*