

PCCP

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

A Review of Methods for the Calculation of Solution Free Energies and the Modelling of Systems in Solution

Cite this: DOI: 10.1039/x0xx00000x

Received 00th January 2012,
Accepted 00th January 2012

DOI: 10.1039/x0xx00000x

www.rsc.org/

R.E. Skyner^{a†}, J.L. McDonagh^{a†}, C.R. Groom^b, T. van Mourik^a and J.B.O. Mitchell^{a*}

Over the past decade, pharmaceutical companies have seen a decline in the number of drug candidates successfully passing through clinical trials, though billions are still spent on drug development. Poor aqueous solubility leads to low bio-availability, reducing pharmaceutical effectiveness. The human cost of inefficient drug candidate testing is of great medical concern, with fewer drugs making it to the production line, slowing the development of new treatments. In biochemistry and biophysics, water mediated reactions and interactions within active sites and protein pockets are an active area of research, in which methods for modelling solvated systems are continually pushed to their limits. Here, we discuss a multitude of methods aimed towards solvent modelling and solubility prediction, aiming to inform the reader of the options available, and outlining the various advantages and disadvantages of each approach.

1. Introduction

Poor aqueous solubility is a major cause of attrition (failure) in the pharmaceutical development process and remains a vital property to quantify in the development of agrochemicals, and in the identification and quantification both of metabolites and of potential environmental contaminants. It is estimated that around 70% of pharmaceuticals in development are poorly soluble with 40% of those currently approved also being poorly soluble.^{1,2} Solubility is determined by structural and energetic components emanating from solid phase structure and packing interactions, in addition to relevant solute–solvent interactions and structural reorganisation in solution. In this review, we focus on the methods currently available to model the solution phase and to predict solubility for a wide range of applications including ligand binding, molecular property prediction and molecular design.³ Readers specifically interested in solubility prediction are also referred to the solubility challenge.⁴ Accurate and timely prediction of solubility could save time and money in drug development, agrochemical development and environmental monitoring. An early-stage analysis of drug and agrochemical candidates allows organisations to focus on those molecules most likely to meet their required solubility

criteria. Many models exist in this area, with differing levels of accuracy, physical interpretability, and calculation time.

Quantitative Structure Activity Relationships (QSAR) and Quantitative Structure Property Relationships (QSPR) are very successful in this field, providing good predictive results at a reasonably low computational cost. These models, however, tend to be limited to molecules similar to those used in their training set. Moreover, these models lack a full physical interpretation, although some do allow assessments of descriptor importance that can perhaps to some extent be physically interpreted.

Several fitted or derived general equations, which take only a few pieces of empirical data as arguments, have also been produced. One of the most successful is the General Solubility Equation (GSE),⁵ taking the melting point and the base ten logarithm of the partition coefficient (logP; partition coefficient for neutral molecules in octanol and water) as empirical input.

The field has also seen the revival of old ideas as new automated data driven design protocols, such as Matched Molecular Pair Analysis (MMPA).⁶ MMPA allows one to acquire previously ‘unknown’ data from existing data sets by exploring how a single molecular change can impact a particular property or activity of interest. We now see large

46 scale data mining following these kinds of protocols, consortia
47 such as SALT MINER, and programs developed by individual
48 companies such as GSK's BioDig^{7,8}. 100

49 In addition to these approaches, we see physics based
50 models ranging from classical simulations to quantum chemical
51 calculations being applied to solubility prediction. These
52 methods vary greatly in complexity. Classical simulations
53 encompass simple Molecular Dynamics (MD), studying
54 interactions between solute and solvent, to more complex
55 perturbations of solutes in the solution phase to a gas phase.
56 Recent advances have seen a new generation of polarisable
57 force fields emerging with a greater capacity to account
58 changes in the electronic charge distribution. Many of these
59 forcefields utilise multipole moments, as opposed to point
60 charges, to capture the anisotropy of the charge distribution.
61 Forcefields such as Atomic Multipole Optimised Energetics
62 Bimolecular Applications (AMOEBA) have been used to study
63 the solvation dynamics of ions⁹. Newer, polarisable forcefields
64 such as the quantum chemical topology forcefield (QCTHFF)
65 use multipolar electrostatics calculated based on quantum
66 chemical topology supplemented with machine learning
67 (Kriging) to model the system. This forcefield has been used to
68 model amino acids with small water clusters¹⁰. These models
69 can be mixed with a quantum chemical core region in mixed
70 Quantum Mechanics – Molecular Mechanics (QM/MM)
71 approaches. Other common models include those representing
72 the solvent as a continuous field with no explicit solvent
73 coordinates. In most cases, these models come at much higher
74 computational cost than their informatics counterparts, and
75 often at lower accuracy. However, if such a method were
76 feasible and accurate enough to predict solubility, it would
77 have a domain of applicability restricted by the molecules
78 within a training set and would also be physically interpretable.
79 Thus, there is a continuing search for such physical methods.
80 These methods have proven useful for modelling the solution
81 approximating the solution phase, hence their applications are
82 diverse and widespread outside of solubility prediction. 130

83 1.1 Thermodynamics and Solubility 131

84 A solution is considered as an equilibrium state between
85 solute and solvent, reaching equilibrium when the number of
86 molecules transferred from the solution to a non-solute state
87 equal to the transfer of molecules from a non-solute state to
88 solution, i.e. when the forward rate is equal to the backward
89 rate and both phases are in equilibrium. Solubility is a
90 quantitative term, most simply describing the amount of
91 substance that will dissolve in a given amount of solvent, and
92 a property of thermodynamic equilibrium. A second process
93 involved in solvation is dissolution; a kinetic term describing
94 the rate at which a substance is transferred from a non-solute
95 phase into solution. Solubility and dissolution are fundamental
96 terms describing the process of solvation, and are related by the
97 Noyes-Whitney equation¹¹; 146

$$\frac{dW}{dt} = \frac{kA(C_s - C)}{L} \quad (1.)$$

where dW/dt is the rate of dissolution, A is the solute surface
area in contact with the solvent, C is the instantaneous solute
concentration in the bulk solvent, C_s is the diffusion layer solute
concentration (given from the solubility of the molecule with
the assumption that the diffusion layer is saturated), k is the
diffusion coefficient, and L is the diffusion layer thickness.

As solubility is a thermodynamic term, it is inherently
affected by factors such as temperature and pressure, as well as
ionisation, solid state effects, and gaseous partial pressure for
solvated gases.

pH is considered to have a significant effect on solubility, as
many organic molecules can behave as weak acids or weak
bases, due to ionisable basic or acidic functional groups, with
polarisation of ionisable groups in solution increasing or
decreasing the overall solubility. The pH of the aqueous
solution in which such molecules are dissolved determines
whether the molecule exists in its neutral or ionised form. The
charged form of a molecule is more soluble, and thus the
aqueous solubility of a substance is pH-dependent¹². This
dependence is described by the Henderson-Hasselbalch (HH)
equations as follows;

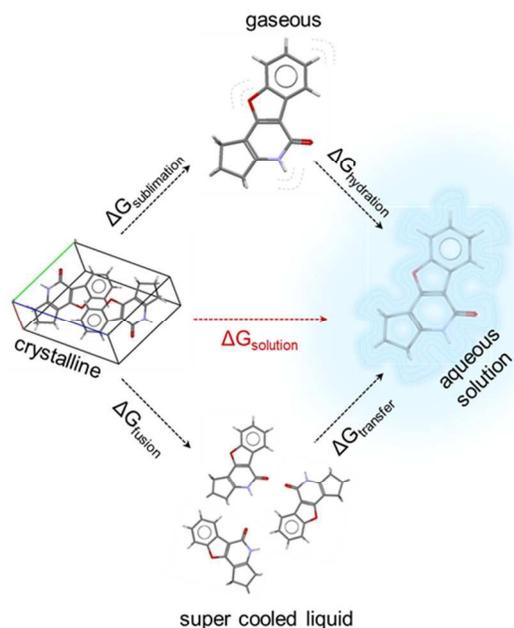
$$\log S_{total}^{acidic} = \log S_0 + \log(1 + 10^{pH-pKa}) \quad (2.)$$

$$\log S_{total}^{basic} = \log S_0 + \log(1 + 10^{pKa-pH})$$

where S_{total} is the equilibrium (thermodynamic) solubility, $\log S_0$
is the intrinsic solubility, defined as the solubility of an
unionised species in a saturated solution, pKa is the negative
logarithm of the ionisation constant of the molecule, and the
final term on the right hand side is the solubility of the ionised
form¹². The HH relationship can be utilised in the prediction of
pH-dependent aqueous solubility of drugs when the pKa and
 $\log S_0$ values of a compound are known¹³. The intrinsic
solubility is a particularly important quantity as it can be used
to find the pH dependent profile and estimate the pKa , it is a
quantity required by industry and hence the focus of several
prediction methods¹⁴. The pH dependant profile of a drug is
particularly important in pharmaceuticals, as it has a direct effect
on the absorption profile of a drug once it has entered the body.
A basic drug-like molecule at a high pH (>2 pH units above the
 pKa) will be fully unionised with solubility at a minimum
(intrinsic solubility). Protonation of the base increases as pH
becomes more acidic, and solubility increases. When pH and
 pKa are equal, half of the solute molecules are protonated and
the solubility of the drug becomes double the intrinsic
solubility. According to the HH equation, this rise in solubility
increases indefinitely with decreased pH, however in practice a
limit is reached at the salt solubility. Two intersecting
concentration curves for the base solubility and the salt
solubility can be combined to give a composite curve for base
solubility as a function of pH. If any one point on this curve is
known (solubility and pH at which it was measured), the whole
curve can be predicted providing pKa and the acid solubility
factor C_{OA}/C_{OB} (the ratio of S_0 of acid to S_0 of base) is known¹⁵.

148 Intermolecular interaction strengths play an important
 149 in the solvation of substances from the solid state. Solutes
 150 which exhibit weak intermolecular forces (i.e. are weakly
 151 bound) tend to have a higher solubility, as the energy cost of
 152 breaking up the lattice is lower. Polymorphic effects can
 153 lead to complications in solubility prediction. A classic
 154 cited example of this is the case of the anti-HIV drug
 155 Ritonavir^{16,17}, in which a polymorphic shift led to a significant
 156 change in solubility, leaving the drug with a greatly reduced
 157 bio-availability. This exemplifies the consideration of solubility
 158 as a property which is dependent upon solid, solute, solvent
 159 and solution state properties and interactions.

160 Two common approaches to the calculation of the Gibbs free
 161 energy of solution utilise a thermodynamic cycle approach.
 162 A first approach calculates the free energy of solution by the
 163 addition of the free energy of sublimation (taking the molecule
 164 in the crystalline phase and subliming it into the gaseous phase)
 165 and free energy of solvation (taking the molecule in its gaseous
 166 phase and solvating it into aqueous solution). An example of
 167 this approach is shown in section 5 of this review, and other
 168 examples are also cited within the literature^{14,18,19}. A second
 169 approach involves calculation of the free energy of solution by the
 170 addition of the free energy of fusion (taking a molecule from
 171 the crystalline state to a hypothetical supercooled liquid)
 172 the free energy of transfer (transfer from a supercooled liquid
 173 into aqueous solution). This method is widely cited within the
 174 literature, and common GSE methods are also derived from this
 175 approach⁵. Both thermodynamic cycle approaches are depicted
 176 in Figure 1.



177
 178 Fig. 1- Calculating the Gibbs free energy of solution is often achieved through the
 179 utilisation of thermodynamic cycles. Two routes are depicted here. The first
 180 route is shown at the top of the diagram, whereby a molecule is taken in
 181 crystalline form and sublimed, and then hydrated. The addition of the Gibbs
 182 energy terms of these processes gives the free energy of solution. The second
 183 thermodynamic cycle is represented at the bottom of the diagram, whereby the

molecule is taken in its crystalline form and undergoes fusion into a hypothetical supercooled liquid, and then is transferred into aqueous solution. The addition of the free energy terms for these two processes also gives the Gibbs free energy of solution.

The solid state is an important consideration for the initial crystalline phase calculated within thermodynamic cycle approaches. Lattice minimisation calculations and periodic DFT provide excellent tools for modelling these systems. Recent advances in these methods show promise for improving predictions, these include updated codes and improved dispersion corrections in periodic DFT^{20,21}.

Complete polymorphic screening and prediction still eludes our capabilities and hence hampers our ability to predict solubility from purely first principles.

A further consideration is that of the standard states used in the different physical states. Typically sublimation data is reported in a 1 atmosphere standard state. Solvation is typically quoted in the Ben-Naim standard state of 1 mol/L with a fixed centre of mass. The difference between the two standard states is a constant 1.89 kcal/mol (7.91 kJ/mol), calculated as $\Delta G_{\text{atm} \rightarrow \text{mol/L}} = RT \ln(24.46)$, where 24.46 is the molar volume at ambient conditions).

The free energy of solution can be calculated directly by the following formula:

$$\Delta G_{\text{solution}} = -RT \ln(S_0 V_m) \quad (3)$$

$$\log(S_0 V_m) = \frac{-\Delta G_{\text{solution}}}{2.303RT}$$

where S_0 is the intrinsic solubility V_m is the crystalline molar volume, R is the gas constant and T is the temperature in Kelvin (K).

A convenient formula¹⁹ allows the solution free energy to be calculated using the native standard states, and removes the dependence on the crystalline molar volume.

$$S_0 = \frac{-p_0}{RT} \exp\left(\frac{\Delta G_{\text{sub}}^{1 \text{ atm}} + \Delta G_{\text{solv}}^{1 \text{ mol L}^{-1}}}{RT}\right) \quad (4)$$

214 2. Informatics – ‘Smart’ Machines in Solubility 215 Prediction

216 Informatics is the science of information processing,
 217 storage, and data mining. There are many applications and
 218 methodologies available for this type of task. Commonly used
 219 methods in chemistry are QSAR/QSPR in which are models
 220 built from data. These models correlate structural features of
 221 molecules with physical properties of interest. A major
 222 supposition of QSPR is that molecules similar in structure will
 223 have similar physical properties, and for QSAR models,
 224 perhaps chemical or biological similarities. Therefore it is
 225 possible to train a model defining a specific relationship
 226 between structure and property/activity on a training dataset,
 227 and apply it to similar molecules to predict their properties and
 228 activities. For this reason, QSAR/QSPR models are not broadly
 229 applicable (i.e., they cannot be applied to molecules differing

considerably from the training set). While QSPR was dominated by multiple linear regression, nowadays machine learning represents the state of the art. Both regression machine learning protocols can identify these structure-property relationships by correlating structural features with experimentally determined physical data. A brief introduction to some of these methods is provided below, and for a more detailed account, see “An Introduction to Cheminformatics” and references therein. Initially, one must represent a molecule in a machine readable format to enable the calculation of molecular descriptors. Two of the most common methods doing this are the Simplified Molecular Input Line Entry System (SMILES)²⁴ and the IUPAC International Chemical Identifier (InChI)²⁵.

2.2 Molecular Descriptors

Descriptors represent physical, chemical, topological energetic features of chemical structures, and can vary greatly in form and derivation. In general, a descriptor is a vector of single numerical values (features), each encoding specific information about an individual molecule.²² This information can be a simple number, such as the molecular weight or count of a specific atom type, or they can be a prediction corresponding experimental quantities, such as the octanol-water partition coefficient (usually expressed as $\log P$). Alternatively, they can also be derived from semi-empirical quantum chemistry. Clearly the cost of calculating different descriptors can vary dramatically. It is often the case that descriptors offering higher levels of refinement, and therefore more useful molecular discrimination, incur a higher computational cost.²² There are many different molecular descriptors and numerous pieces of software to calculate them.²²

2.3 Methods

2.3.1 REGRESSION

Regression analysis is a fundamental tool in informatics. Simple linear regression expresses a relationship between a scalar dependent variable Y and a single explanatory

independent variable X. Multiple Linear Regression (MLR) extends this to allow for multiple dependent y_i variables or explanatory independent variables x_i , expressed as;

$$y = \sum_i^j \alpha_i x_i \quad (5)$$

These methods have seen widespread use in many fields.²⁶ A disadvantage of MLR is the apparent ease of over-fitting. It is suggested that a useful rule of thumb is that the number of data points should be in excess of five times the number of explanatory variables^{22,23}.

2.3.2 RANDOM FOREST

Random Forest (RF), is a learning method based on decision trees. These are stacked sets of binary separators following a tree like graph structure. RF uses a ‘forest’ of these decision trees, making use of “*the wisdom of crowds*”; hence, is considered an ensemble learning method. RF can be used for classification or regression. For application to classification problems, the binary splitting is based upon the Gini index, which is a calculation of the maximal discrimination of the data points. For regression, splitting is generally based on a minimisation of the root mean squared error (RMSE). The initial node is known as the root node, with subsequent nodes being called branch nodes. The final nodes are referred to as leaf nodes and contain molecules with similar predictions of the property or activity.^{14,23}

2.3.3 SUPPORT VECTOR MACHINES

Another commonly used machine learning method is that of Support Vector Machines (SVM). SVM supports both regression and classification tasks, and is capable of handling multiple continuous and categorical variables. Methods for handling classification tasks are based on typically non-linear kernel functions. These kernel functions allow the transformation of datapoints into a higher dimensional feature space.

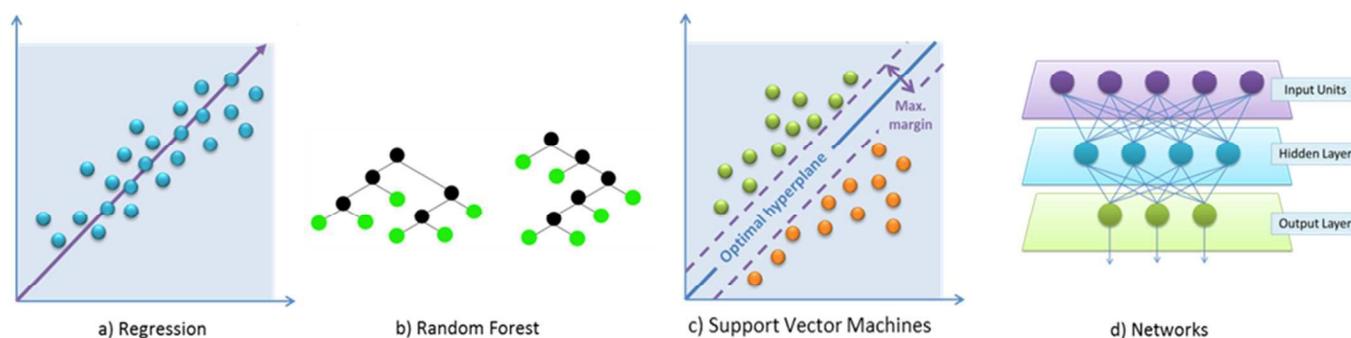


Fig. 2 – Machine learning methods; a) Regression analysis aims to describe how the typical value of the dependent variable changes as the independent variables are changed. The regression function (purple arrow) characterises variation; b) Decision trees consisting of a binary separation at the nodes, leading to predictions or classifications at the leaf nodes (green circles); c) An example of SVM separates data into distinct categories by an optimal hyperplane, which should have optimal margins either side for a clear distinction in data categorisation; d) A typical network consists of layers of nodes. All nodes have connections with all other nodes in adjacent layers. The input units (top) do not count as a layer of nodes, as they do not carry out any typical arithmetic operations. A typical arithmetic operation is the

300
301
302
303
304

305
306

generation of a net signal and transformation by a transfer function into an output signal. The input units distribute input values to all of the neurons in the layer below. The connections between nodes each have a different weight, representing different descriptors used in machine learning.

307
308

309 SVM training algorithms are built up of binary categorisation
310 data, whereby a particular data point belongs to one of
311 categories. Thus, the test set data is also categorised, producing
312 a clear separation, which should be as wide as possible, in
313 feature space. Alternatively, in the case of regression, the
314 surface behaves analogously to a regression line, providing the
315 maximal explanation of the data within the bounds of an
316 acceptable error margin whilst attempting to remain relatively
317 flat to avoid overfitting.^{22,23}

318 2.3.4 NETWORKS

319 Artificial Neural Networks (ANNs) and deep learning
320 architectures are another common form of machine learning
321 method in chemistry. These are models conceptually based on
322 the brain's neuron network (although a great simplification).
323 ANNs contain an input layer which receives the molecular
324 information, an output layer which provides the prediction
325 the user, and between these at least one hidden layer which is
326 trained using data to link the neurons of the input layer to the
327 output layer in a suitable fashion for the problem at hand. The
328 training generally involves weighting specific paths between
329 the neurons.^{7,8,13} Deep learning architectures aim to enhance
330 the learning capabilities of machine learning methods such as
331 ANNs. Deep learning algorithms attempt to abstract data at a
332 high-level through model architectures comprising multiple
333 non-linear transformations. In the case of ANNs, enhanced
334 abstraction can be achieved through the addition of hidden
335 layers, capturing the interaction of many factors which
336 contribute to the observed data.

337 2.4 THE GENERAL SOLUBILITY EQUATION (GSE)

338 GSE (as briefly mentioned in the introduction) is a QSPR
339 model based on the melting point and the octanol-water
340 partition coefficient $\log P$ of a chemical substance, used to
341 predict the aqueous solubility of non-ionisable compounds
342 and acts as a useful guide for ionisable compounds using
343 lipophilicity $\log D$ at the pH of the aqueous buffer employed.
344 The equation states that;

$$345 \log S = 0.5 - 0.01(m.p. \text{ } ^\circ\text{C} - 25) - \log P$$

346 Or in terms of $\log D$;

$$347 \log S_{pH(x)} = 0.5 - 0.001(m.p. \text{ } ^\circ\text{C} - 25) - \log D_{pH(x)}$$

348 GSE is a simple QSPR model, with powerful predictive
349 ability (coefficient of determination (r^2) = 0.96 and root mean
350 squared error (RMSE) = 0.53 (units) for a data set of 1000
351 organic molecules²⁹), and the simplicity of the model means
352 it has found wide application in the pharmaceutical industry.
353 However, the reliance of the GSE on experimentally
determined descriptors limits its applicability, and data

sparingly populated at their limits can lead to overestimation of the model's predictive power³⁰.

Ali *et al.*³⁰ have revisited the GSE and have attempted to relieve the reliance of the GSE on the experimentally determined melting point by replacing it with a descriptor that describes the topological polar surface area (TPSA). They demonstrate the effects of inflated predictive power of the GSE by using a subset of an initial dataset, which reduced the overall predictive power of the GSE by approximately 6.4%. TPSA was included in a revised model to account for the fact that 88.5% of poorly performing compounds contained polarisable groups. The pure GSE model employed provided $r^2 = 0.818$, and the TPSA replacement of melting point model provided $r^2 = 0.813$, showing a comparable effectiveness. The number of compounds containing polarisable groups with $\log S$ predicted within ± 1 log unit of experimentally determined values was also higher for the revised TPSA model (83.2% TPSA; 79.6% GSE). A final model combining melting point, $\log P$ and TPSA was also tested, and was found to have a better predictive power than both of the previously employed models ($r^2 = 0.869$) with 90.8% of compounds containing polarisable groups predicted within ± 1 log unit of experimentally determined values.

The work of Ali *et al.*³⁰ highlights the importance of reliable descriptors in improving the overall performance of QSPR models, particularly when polar or polarisable functionality is included in test sets, and when experimentally determined values are required. As such, experimentally determined values may be best suited only for comparative analysis of predictive models to experimental data as a measure of performance in many cases.

349 2.5 Other Cheminformatics Applications

A recent approach to predict solubility proposed by McDonagh *et al.*¹⁴ applied three models, exploiting both cheminformatics descriptors and theoretically derived thermodynamic properties. The initial models use theoretical chemistry and QSPR models alone, with further development combining the two approaches into a unified QSPR model. The developed models aim to calculate solubilities in agreement with experiment and in a reasonable time period. It was found that quantitatively accurate solvation free energies are unobtainable from the specific simple theoretical chemistry approach applied. The authors suggest that QSPR models are the most effective method, when both time and accuracy are considered. The machine learning methods employed, which use a modest number of cheminformatics descriptors, predict solubility values comparable to those obtained with currently available commercial software. Notably, only a small improvement in accuracy was found on combining the two approaches. This suggests that the cheminformatics descriptors

404 and the theoretically derived quantities are not very
405 complementary, but duplicate much of the same information.
406 Another recent approach, by Lusci *et al.*²⁷, applies deep
407 learning to the solubility prediction problem. The deep learning
408 method is based on recursive neural networks adapted to
409 undirected graph representations of molecules. The method
410 produces good predictions of solubility on a number of standard
411 datasets in the field²⁷.

412 A further example of a cheminformatics approach is
413 demonstrated by Shayanfar *et al.*³¹ who apply a simple QSPR
414 model to the prediction of aqueous solubility of drugs, validated
415 by cross-validation. A training set of 220 drug-like molecules
416 was used to build a model with MLR. Seven descriptors
417 (aqueous solubility from the literature, solute, melting point,
418 experimental logP, calculated Abraham solvation parameters,
419 calculated ClogP values and calculated melting points) were
420 used to develop a two-variable model. The two variables used
421 gave an R^2 value of 0.934 and a standard error estimate of
422 0.893. The proposed model was compared to a GSE model and
423 a linear-solvation-energy-relationship (LSER) model.
424 Correlations between each model's computationally determined
425 values of aqueous solubility with corresponding experimental
426 values gave an $R^2=0.62$ for GSE, $R^2=0.57$ for LSER and
427 $R^2=0.66$ for the proposed MLR method.

428 Recent work has also suggested that, contrary to popular
429 arguments, the quality of the experimental data available is not
430 the limiting factor for the predictive accuracy of solubility
431 predictions obtained from cheminformatics models.³² This
432 work may suggest that inherent limitations within the models
433 are responsible for the largest part predictive errors.

434 3. Implicit solvation – An isotropic field as a solvent 435 representation

436 Continuum solvation models consider solvent as a
437 continuous isotropic medium. An underlying assumption of
438 implicit solvation models is that explicit solvent molecules may
439 be removed from the model; provided that the continuous
440 medium replacing them sufficiently represents equivalent
441 properties.

442 A simplification of continuum models can be thought of in
443 terms of a Hamiltonian as;

$$\hat{H}^{tot}(r_M) = \hat{H}^M(r_M) + \hat{H}^{MS}(r_M) \quad (8.)$$

444 where M refers to a single solute molecule, S refers to the
445 solvent, and r refers to position. Solvent coordinates do not
446 appear within the Hamiltonian term, exemplifying the
447 representation of solute in a continuum, rather than as definite
448 atoms, as with explicit models. \hat{H}^{MS} is a sum of different
449 interaction operators, which can be expressed in terms of
450 solvent response functions, indicated by $Q_x(\vec{r}, \vec{r}')$ where \vec{r}
451 indicates a position vector, and x represents a contribution
452 interaction. More in-depth discussions are available in
453 textbooks specific to computational chemistry, such as that by
454 Cramer³, and reviews by Tomasi *et al.*¹⁵

In a standard continuum model, generally represented by
Polarisable Continuum Models (PCM), solute-solvent
interaction energies can be represented by a number of Q_x
operators. The free energy of M is therefore described by an
expression of five terms;

$$G(M) = G_{cav} + G_{el} + G_{dis} + G_{rep} + G_{tm} \quad (9.)$$

with the order of terms corresponding to the best performing
order of the 'charging processes', integration processes
coupling a distribution function with a potential function. The
terms are the free energy of cavitation, electrostatic energy,
dispersion energy, repulsion energy and thermal fluctuation,
respectively.

3.1 Continuum Models for Electrostatic Interactions

PCM models are advantageous in that they can represent a
statistically averaged (continuum) solvent so that meaningful
results can be acquired within a single calculation. PCM models
have been particularly useful in modelling reactivity and
spectroscopy of various solvents with different polarities.³³

In a solvent-solute system where atom Q (solute) has a
positive charge, solvent water molecules will preferentially
orientate their negative dipoles towards the solute's positive
charge (Fig. 3, left). For a single water molecule, there is only a
slight preference in orientation, which is smaller than that of its
average thermal fluctuations. Therefore, this effect is averaged
over the long range of electrostatic interactions of water in the
bulk (Fig. 3, right). For an isotropic solvent with random
thermal motion, the average electric field is zero at any given
point. However, introduction of a solute gives a net change in
orientation, introducing an overall change in electric field,
known as the 'reaction field'.

Accounting for the reaction field increases the solute's
polarity proportionally to the solute polarisability, and the
strength of the external electric field. This causes an increase
in the dipole moment of Q, consequently polarising and increasing
the change in orientation of the solvent to oppose the dipole
moment of Q.³

There are energy costs associated with both the orientation
and polarisation of the solvent, and the dipole moment of Q. As
solvent molecules oppose the dipole moment of Q, they interact

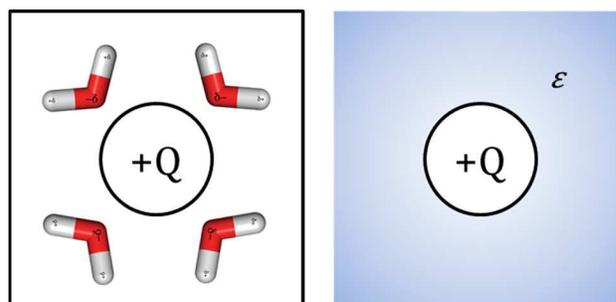


Fig. 3 - Left - water molecules reorient themselves to preferentially point the negative end of their dipole towards the positive solute charge (+Q). Right - The system is modelled with a continuous polarisable field. Polarisability is represented by the bulk dielectric constant, ϵ .

498 unfavourably with the reaction field. They also 535
 499 configurational freedom, with an associated free-energy cost 536
 500 a continuum model, the charge distribution of a solvent 537
 501 represented as a continuous electric field, statistically averaged 538
 502 over all degrees of freedom at thermodynamic equilibrium. 539
 503 electric field at any given point is the gradient of 540
 504 electrostatic potential. The work required to create the charge 541
 505 distribution is determined from the interaction of solute charge 542
 506 density ρ with the electrostatic potential ϕ from;

$$G = \frac{1}{2} \int \rho(r)\phi(r)dr \quad (15.46)$$

507 The polarisation component of G , which we call G_p , is 549
 508 difference between charging the system in gas and solution 550
 509 phases; thus only the electrostatic potentials in both gas 551
 510 solution phases are needed to calculate G_p . 552

511 PCM methods are generally applied through two models; 553
 512 the Poisson-Boltzmann (PB) model, and the Generalised Born 554
 513 (GB) models. Both models are advantageous for different 555
 514 systems, and the accuracy of either model is mostly dependent 556
 515 upon the suitability of the cavity type used to surround 557
 516 solute molecule within an ideal solvent system. 558

517 3.1.1 THE POISSON-BOLTZMANN (PB) MODEL 559

518 The Poisson equation (eqn. 11) combines the terms 560
 519 electrostatic potential and the differential form of Gauss's 561
 520 to define the electrostatic potential ϕ as a function of 562
 521 dielectric constant ϵ and charge density ρ . When a surrounding 563
 522 dielectric medium responds linearly to an embedded charge 564
 523 Poisson's equation states that; 565

$$\nabla^2 \phi(r) = -\frac{4\pi\rho(r)}{\epsilon} \quad (11.)$$

526 Continuum solvation models represent the charge 566
 527 distribution on the basis of two separate areas: inside (solute) 568
 528 and outside (solvent) of a cavity. For this case, the Poisson 569
 529 equation states; 570

$$\nabla\epsilon(r) \cdot \nabla\phi(r) = -4\pi\rho(r) \quad (12.)$$

530 The Poisson equation as expressed above is valid only for 571
 531 systems under non-ionic conditions. In a real solution, 572
 532 dissolving a solute produces mobile electrolytes. This effect is 573
 533 accounted for by an expansion of the Poisson equation, known 574
 534 as the Poisson-Boltzmann (PB) equation;

$$\nabla\epsilon(r) \cdot \nabla\phi(r) - \epsilon(r)\lambda(r) \frac{8\pi q^2 I}{\epsilon\kappa_B T} \frac{\kappa_B T}{q} \sinh\left[\frac{q\phi(r)}{\kappa_B T}\right] = -4\pi\rho(r) \quad (13.)$$

where q gives the magnitude of electrolyte ionic charge, λ is a function equal to 0 in areas inaccessible to electrolyte ions and 1 for accessible areas, and I indicates the ionic strength of the electrolyte system.

PB equations are best used to calculate the electrostatic potential of systems where the cavitation of solute is near-spherical or ellipsoidal (*ideal cavitation*), as the convergence of the predicted electrostatic component of the solvation free energy ΔG_E is computationally expensive and often inaccurate. Thus, derivations applying approximations of the Poisson equation are often used in continuum models³³, the most common of which are Self-Consistent Reaction Field (SCRf) models, such as the Onsager model.³⁴

A further limitation of PB based models is the definition of cavitation. A number of variational SCRf models have been proposed in order to optimise cavitation parameters, most commonly using tessellation (tiling) of the cavity surface to simplify and reduce iterations of the PB equation.³³

517 3.1.2 THE GENERALISED BORN (GB) MODEL 559

For systems in which ideal cavitation is not accurate, arbitrary cavitation can be applied. Arbitrary cavitation refers to the construction of a cavity around the solute similar to the shape represented by space-filling models generated from the overlap of atomic spheres at volumes representing van der Waals (vdW) radii. An alternative method to SCRf models involves an approximation of the Poisson equation that can be analytically solved, known as the Generalised Born (GB) approach.

A conducting sphere with charge q can be considered representative of a monatomic ion. If the surface of the sphere is assumed to be entirely smooth, the charge distribution around it will be uniform, and the charge density at any point is given by;

$$\rho(s) = \frac{q}{4\pi a^2} \quad (14.)$$

where s is a point on the sphere's surface, and a is the spherical radius. Integrating over the entire outside surface and adding a term for the electrostatic potential, the energy term G , with $|r| = a$, becomes;

$$G = -\frac{1}{2} \int \left(\frac{q}{4\pi a^2}\right) \left(-\frac{q}{\epsilon a}\right) ds = \frac{q^2}{2\epsilon a} \quad (15.)$$

The Born equation for the polarisation of a monatomic ion is calculated from the difference in the required work in the gas and solution phases applied to equation 8;

$$G_p = -\frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \frac{q^2}{a} \quad (16.)$$

The GB method extends the Born equation to polyatomic molecules to express polarisation energy as;

$$G_P = -\frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{k,k'}^{atoms} q_k q_{k'} \gamma_{kk'}$$

578 where k and k' run over all atoms, each with a partial charge
 579 The determination of suitable parameters for γ for polyatomic
 580 systems involves a radial integration of the charge q_k to
 581 determine the interaction of atom k with the surrounding
 582 medium. γ has units of reciprocal length, thus representing an
 583 inverse Coulomb integral. γ is given a suitable functional form
 584 in order to approximate the PB equation, and has a limiting
 585 behaviour, becoming closer to the exact reciprocal length r^{-1} at
 586 large interatomic distances.

587 3.2 Continuum Models for Non-electrostatic Interactions

588 Similarly to the electrostatic components of solvation free
 589 energy, non-electrostatic contributions to the solvation free
 590 energy are not experimentally measurable. The solubility of
 591 experimental systems may be more susceptible to some effects
 592 than others. Various neutral model systems have been
 593 developed in accordance with this.

594 3.2.1 SPECIFIC COMPONENT MODELS

595 Pierotti³⁵ developed a model formula, based on scaled
 596 particle theory, for the calculation of cavitation free energy,
 597 through the observation of the solvation energy for noble gases.
 598 Scaled particle theory is a statistical-mechanical theory of fluids
 599 derived from exact radial distribution functions, to give
 600 expression for the work required to place a spherical particle
 601 into a fluid of spherical particles. Noble gas atoms do not
 602 exhibit permanent electrical moments, thus their transfer into
 603 solution is considered to be the most analogous example of
 604 perfect cavitation.

605 The experimental data from Pierotti's work has been
 606 complemented by simulation data,³⁶ including free energy of
 607 formation data of molecular-sized cavities in 12 common
 608 solvents obtained from free energy perturbation simulations.
 609 Pierotti's formula has since been expanded for molecular
 610 cavities by Colominas *et al.*³⁷

611 A further, specific contributing factor to solvation free
 612 energy is dispersion. A somewhat simplistic explanation of
 613 dispersion is as follows. The average electron cloud of an atom
 614 is spherically symmetrical, but at any instantaneous time point
 615 there may be a polarisation of charge causing an instantaneous
 616 dipole moment. This dipole moment interacts with
 617 neighbouring atoms, inducing a second instantaneous dipole
 618 and so on, and an interaction occurs between these. The in-
 619 phase correlation of instantaneous and induced dipoles means
 620 the overall interaction energy does not average to zero over
 621 time.³ The average interaction energy falls off (largely)
 622 proportionally to r^{-6} (where r is the distance between
 623 interacting particles). The multipole expansion of the dispersion
 624 interaction is written;

$$V(r) = -\frac{C_6}{r^6} - \frac{C_8}{r^8} - \frac{C_{10}}{r^{10}} \dots$$

625 where C_6 , C_8 and C_{10} are dispersion coefficients dependent on
 626 the atomic species. This is normally evaluated as a sum over all
 627 pairs of atoms in different interacting molecules.

3.2.2 ATOMIC SURFACE TENSIONS

Another approach for the evaluation of the non-electrostatic
 components of solvation free energy assumes the non-
 electrostatic component to be atom or group specific, and
 proportional to atomic surface area. A recent review by Wang
*et al.*³⁸ (2009) considers four QSPR aqueous solubility models
 developed on the principle of weighted atom type counts and
 Solvent Accessible Surface Areas (SASA). They note that
 models considering SASA are often developed with small test-
 sets, and are therefore, in common with QSAR/QSPR models,
 poor performers for test molecules dissimilar to the original
 training set. The authors found that SASA descriptors did not
 enhance model performance any further than weighted atom
 type counts. This suggests the influences upon the non-
 electrostatic components of solvation free energy may be more
 complex than simple surface area considerations.

A further notable feature of continuum models based on
 surface tension is the neglect of any other contribution; that is,
 the development of these models assumes surface area as the
 sole determinant of solvation free energy, and that electrostatic
 components are implicit within the calculation parameters
 used.³³

3.3 The Current State of Continuum Models

There are a large number of available continuum solvent
 models, all with relative merits and shortcomings. The
 following is a brief description of those most commonly
 applied.

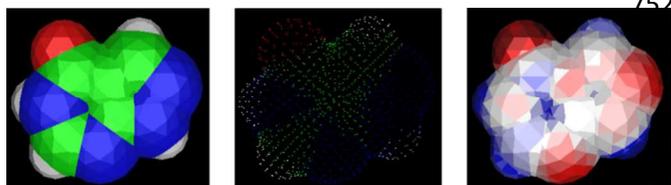
Integral Equation Formalism PCM (IEFPCM) is the current
 version of PCM applied in common quantum chemistry
 packages. IEFPCM is a reformulation of dielectric PCM
 (DPCM) in terms of the integral equation formalism. One of the
 biggest challenges to PCM methods is that they are all derived
 assuming the solute charge density is entirely encapsulated in
 the cavity. This is often not the case, as the electron
 distributions often extend beyond the cavity. IEFPCM has been
 shown to cope well with this effect when compared to other
 PCM based methods.³³

A further variation of PCM is the conductor-like polarisable
 continuum model (CPCM), which is often considered one of
 the most successful solvation models³⁹. The Conductor-like
 screening model/Conductor-like screening model for real
 solvents (COSMO/COSMO-RS)⁴⁰ is a variation on Poisson-
 Boltzmann PCM and CPCM. In COSMO the dielectric
 permittivity (ϵ) is set to infinity ($\epsilon = \infty$). This defines the
 solvent as a conductor, which is suggested as a more realistic
 approximation for strong dielectric media such as water, with
 the first version of COSMO⁴⁰ having values of the dielectric
 constant with a relative error of less than $\frac{1}{2}\epsilon^{-1}$. COSMO has
 been shown to be a reliable and readily available method for
 calculations on the liquid and solution phases. The use of a
 boundary condition for the calculation of total potential in place

679 of a traditional dielectric boundary condition for the electric
680 field found values within 10% of the exact results obtained
681 from dielectric boundary condition methods⁴¹. COSMO-RS
682 extends the COSMO code to also define the ability of the
683 solvent to screen the surface charge on the cavity of the solute.
684 Parametrisation of COSMO and COSMO-RS performed by the
685 software developers tested 217 small to medium neutral
686 molecules, spanning a vast functionality of H, C, N, O and Cl.
687 An overall accuracy of 0.4(rms) kcal/mol for chemical potential
688 differences was achieved⁴¹.

689 A recent addition is the solvation model based on density
690 (SMD). This model applies the IEFPCM protocol, solving the
691 non-homogeneous Poisson equation using a set of optimised
692 atomic Coulomb radii. The non-electrostatic contributions are
693 calculated on the basis of a parameterised function which
694 includes terms for atomic and molecular surface tensions as
695 well as the solvent accessible surface area.⁴²

696 A recent investigation of gas to solution phase standard
697 state Gibbs free energies of solution compares energies
698 obtained for six combustion gas flue compounds at the
699 Gaussian-4 level of theory using IEFPCM, CPCM and SMD
700 implicit solvent models for 178 organic solvents. It is found
701 that IEFPCM and CPCM produce similar ΔG_S values for all six
702 flue compounds, with maximum absolute intra-solvent
703 deviations of <1.6 kJ mol⁻¹. Intra-solvent deviations between
704 the IEFPCM and SMD models up to 45.5 kJ mol⁻¹ were
705 observed. IEFPCM and CPCM also showed strong correlation
706 between calculated solvent ϵ and ΔG_S for all solvents, whereas
707 SMD showed a much more varied relationship⁴³.



708 **Fig. 4 - The PCM cavity of allopurinol.**
709 *Left: The solvent accessible surface of allopurinol from a PCM calculation.*
710 *Middle: The reaction field evaluation points.*
711 *Right: Surface polarisation as a result of reaction field.*

713 4. Explicit Solvation Models

714 Explicit Solvation models are the primary choice of
715 solubility models where solvent-specific effects are considered.
716 The explicit treatment of water should, in principle, provide the
717 most descriptive and realistic model for the investigation of
718 solvation⁴⁴, however it intrinsically requires a large number of
719 degrees of freedom and thus is associated with a phase space of
720 high dimensionality. This requires statistical averaging over the
721 entire phase space, particularly when extracting specific
722 underlying physical behaviour, such as thermodynamic
723 properties.

724 Statistical thermodynamics relates all observable
725 thermodynamic properties to the partition function, Q . The
726 partition function is summarised as;

$$Q = \iint e^{-\frac{E(q,p)}{k_B T}} dq dp \quad (19.)$$

where Q is the classical formulation integrated over all phase space of all spatial q and momentum p coordinates.

Explicit models consider solvation in terms of free energy calculations, with different models for water available, as discussed below.

4.1 Free Energy Calculations – Monte Carlo (MC) and Molecular Dynamics (MD) Simulations

Free energy considerations are distinctly different for intramolecular and intermolecular degrees of freedom. For intramolecular components, free energy contributions rely on vibrational and librational motions on an intramolecular energy surface⁴⁵. For well-defined energy-minima, the free energy is easily accessible from the partition function (eqn. 19) from vibrational frequencies treated with the harmonic approximation. The harmonic approximation estimates the nuclear potential of a molecular system in its equilibrium geometry at a potential energy surface minimum in terms of normal vibrational modes, each governed by a 1D harmonic potential. Anharmonic effects are accounted for with MC or MD simulations for the calculation of entropy on the intramolecular energy surface⁴⁵. Due to diffusion, the particles of a solution system do not exhibit motion definable by harmonic approximations. Thus, conventional MC and MD methods do not involve the direct determination of Q , and exhibit an extremely slow convergence for densities of typical chemical systems, due to the exponential dependence of the Boltzmann factor on the occupation of available energy levels at a given temperature.

4.1.1 FREE ENERGY PERTURBATION (FEP) METHODS

Free Energy Perturbation (FEP) methods were first introduced by Zwanzig⁴⁶ in 1954, who related the thermodynamics of two different systems, in order to evaluate differences in intermolecular potentials. Zwanzig notes that at high temperatures, the forces of repulsion between molecules determine the equation of state of a gas, and that at lower temperatures the equation of state should be determinable by considering forces of attraction as perturbations on the forces of repulsion. The energy change from state A to state B is calculated by;

$$\Delta G(A \rightarrow B) = G_B - G_A = -k_B T \ln \left\langle \exp \left(-\frac{E_B - E_A}{k_B T} \right) \right\rangle_A \quad (20.)$$

where T is temperature, and the triangular brackets indicate an average over the simulation runs for A . A normal simulation run for A coincides with a new energy state of B on each optimisation run. The energy difference between A and B is either between the atoms in each state, or in an isomeric difference, for example A may be the cis- isomer of a structure, and B the trans- isomer, with A and B in different energy states due to different intra- and/or intermolecular interaction. For

774 isomeric differences, the free energy map is calculated along
 775 reaction coordinates. The convergence of FEP calculations
 776 only reliable for a small difference between A and B ,
 777 traditional perturbation theory only holds true for systems
 778 which remain similar upon dissolution.
 779 More recent derivations of Zwanzig's model allow
 780 division of perturbations into smaller calculations, allowing
 781 parallelisation. These models involve breaking the reaction
 782 pathway down into a series of intermediate TS steps, allowing
 783 better convergence between the initial and final structures
 784 investigated.⁴⁷ However, FEP calculations remain one of the
 785 most computationally expensive methods for calculating free
 786 energy differences.
 787 An example of this is shown by Lüder *et al.*⁴⁸ who have
 788 investigated the effectiveness of FEP methods for the
 789 calculation of free energy of solvation in pure melts for 46
 790 molecules. Simulations were performed in two stages, starting
 791 down the Coulomb and Lennard-Jones (LJ) interactions
 792 independently. Results were interpreted under the assumption
 793 that the free energy of the liquid to vapour process ΔG_{vl} can
 794 be calculated from the sum of the free energy term for cavity
 795 formation ΔG_{cav} and the energy associated with LJ interactions
 796 and Coulomb interactions (over 2). ΔG_{cav} is obtained from hard
 797 sphere theories. Interaction energies and molar volumes for each
 798 of the 64 drug molecules were compared for systems
 799 comprising 260 molecules. Deviations between systems were
 800 found to be an average of 2.9% for intermolecular interaction
 801 energy, and 1.4% for molar volume, suggesting the dataset
 802 selected would provide reliable results. Predicted and simulated
 803 ΔG_{cav} values are found to be systematically underestimated
 804 by approximately 15%. An overall average deviation of calculated
 805 ΔG_{vl} values in comparison to experiment is -1.8 kJ/mol, within
 806 reasonable errors expected in the range -1 to 1 kJ/mol. This
 807 investigation suggests that overall, FEP methods require more
 808 work at the theory level, particularly due to systematic errors
 809 that occur in phase space relationships between reference and
 810 perturbed systems.
 811 An alternative approach to calculating the free energy
 812 difference from one state to another is to treat the change from
 813 A to B as a transformation; rather than to calculate free energies
 814 of independent structures, and calculate an energetic difference
 815 as in traditional FEP methods.³
 816 A recent application of this method, derived from FEP,
 817 has been demonstrated by Liu *et al.*⁴⁹ for the calculation of the
 818 solubility of gases in ionic liquids. The Bennett acceptance ratio
 819 (BAR) method utilises the method of transferring between
 820 states instead of treating each state as an individual structure.
 821 The Coulomb and LJ terms are calculated separately. It is found
 822 that simulated solubilities are found in good agreement with
 823 Henry's law constants. However, comparison to experimental
 824 data finds poorly soluble gases to have larger errors, with
 825 underestimated and overestimated gas solubilities found with
 826 similar calculation methods in complementary studies.

4.1.2 ENTHALPY-ENTROPY DECOMPOSITION

827

A further offshoot of free energy calculations is the
 decomposition of the free energy term into enthalpic and
 entropic components. Entropy and enthalpy complement free
 energy as they provide interpretive information to link
 molecular perturbations and thermodynamic changes. Two
 solutes may have similar hydration free energies (HFE), but
 may have solubilities dependent on distinct chemical function.⁴⁴
 As both enthalpy and entropy are experimentally measurable,
 the difference between theory and experiment is ascertainable,
 and may be applied as benchmarks for force field
 optimisations,⁴⁴ and give insight into the mechanism of
 solvation. Levy and Gallicchio have reviewed a variety of
 different approaches to the thermodynamic decomposition of
 free energies.⁴⁴

Wyczalkowski *et al.*⁵⁰ recently proposed two new methods
 for the estimation of entropy and enthalpy decomposition of
 free energy calculations, evaluated for the solvation of N -
 methylacetamide (NMA). The methods investigated found
 thermodynamic contributions to be in disagreement with
 experimental data, highlighting the difficulty in obtaining
 decompositions comparable in quality to free energy estimates,
 with thermodynamic decomposition of computational
 Helmholtz free energies of solvation (ΔF at fixed volume)
 values yielding errors approximately two orders of magnitude
 larger than the initial ΔF values found. It is noted that ΔF
 values are statistically reliable and can be used for quantitative
 comparison to experimental data. The calculation of entropic
 and enthalpic contributions is also extremely computationally
 demanding, as every temperature point of a simulation requires
 recalculation of the overall free energy.³ The authors highlight
 that where calculation of free energies of solvation has
 advanced so that computational errors are on par with
 experimental ones, thermodynamic decomposition calculations
 suffer from statistical errors 10-100 times larger than free
 energy of solvation calculations.

A recent study by Ahmed and Sandler⁵¹ uses the
 decomposition of free energies of hydration and self-solvation
 of low polarity nitrotoluenes to consider an array of
 thermodynamic terms and physicochemical properties. These
 include: solid-phase vapour pressures, solubilities, Henry's law
 constants, hydration and self-solvation entropies, enthalpies,
 heat capacities and enthalpies of vaporisation or sublimation.
 Their study focuses on the temperature-dependence of various
 terms. Decomposition of hydration free energies into enthalpic
 and entropic contributions is performed by a method utilising
 polynomial fitting of temperature-dependent self-solvation free
 energies (with respect to temperature). The use of fitting
 increases the sensitivity of derived values of hydration free
 energies. Self-Solvation enthalpy (ΔH_{self}) values and entropy
 ($T\Delta S_{\text{self}}$) values are calculated within approximately 2 kcal/mol
 of experimentally determined values.

4.2 Combined Quantum Mechanical / Molecular Mechanical Methodologies (QM/MM)

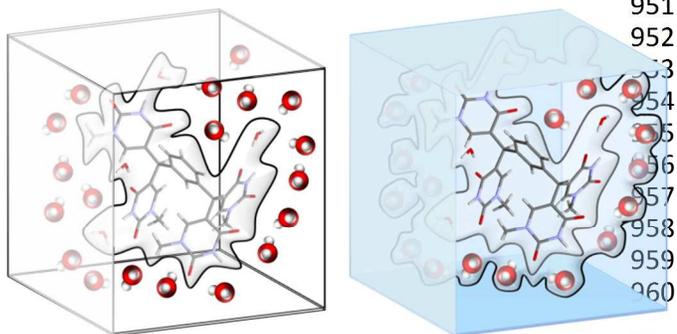
Explicit solvation models are often developed with respect
 to biological systems, due to the role of water in catalytic

883 mechanisms, protein folding and protein-DNA recognition⁹²⁵
 884 name but a few, which all require the specific detail of exp⁹²⁶
 885 water-substrate interactions to hold descriptive meaning.⁹²⁷
 886 particular interest are combined QM/MM models, with ⁹²⁸
 887 describing electronic system changes (where precise syst⁹²⁹
 888 description is needed) and the rest of the system (where ⁹³⁰
 889 precision is required) being described by a MM force field⁹³¹
 890 Applications of QM/MM combined models are discussed ⁹³²
 891 recent review.⁵² ⁹³³

892 The foundational concepts involve the partitioning of ⁹³⁴
 893 desired system into two subsystems: the QM subsystem ⁹³⁵
 894 containing a small number of atoms and described by quantum ⁹³⁶
 895 mechanics, with the remainder of the system described by a ⁹³⁷
 896 suitable MM force field. The Hamiltonian of the whole system ⁹³⁸
 897 is simply written; ⁹³⁹

$$H = H_{QM} + H_{MM} + H_{QM/MM}$$

898 where H_{QM} is a QM Hamiltonian, H_{MM} is an empirical force ⁹⁴²
 899 field and $H_{QM/MM}$ describes interactions at the QM/MM ⁹⁴³
 900 interface. The energy of the system is also described as the sum ⁹⁴⁴
 901 of QM, MM and QM/MM contributions. This model is often ⁹⁴⁵
 902 referred to as a two-layered approach (Fig. 5, left). A derivative ⁹⁴⁶
 903 of this model involves adding a third “layer” as a continuum ⁹⁴⁷
 904 solvent representation around the MM region, and is known as ⁹⁴⁸
 905 a three-layered approach (Fig. 5, right). ⁹⁴⁹



906
 907 Fig 5. – Left – two-layered approach to the QM/MM method. The solute ⁹⁵¹
 908 molecule and a few water molecules are treated with QM (centre) and the rest ⁹⁵²
 909 of the solvent system is represented by MM up to a user-defined distance. ⁹⁵³

910 Right – three-layered approach – an additional layer surrounds the MM region ⁹⁵⁴
 911 and uses a continuum approach to describe the long range solvent in the bulk. ⁹⁵⁵

912 Theoretically, any desired level of accuracy can be used ⁹⁶⁶
 913 within the QM region of the simulated system, within the scope ⁹⁶⁷
 914 of available methods. However, more accurate methods are ⁹⁶⁸
 915 susceptible to high computational cost. Thus, careful ⁹⁶⁹
 916 consideration is required by the user as to what level of ⁹⁷⁰
 917 accuracy is required, and at what cost. A succinct overview of ⁹⁷¹
 918 different available QM methods is provided by Friesner and ⁹⁷²
 919 Guallar⁵² for QM/MM methods applied to enzymatic catalysis. ⁹⁷³
 920 with descriptions, advantages and disadvantages of respective ⁹⁷⁴
 921 QM methods available in textbooks such as the one by ⁹⁷⁵
 922 Cramer.³ ⁹⁷⁶

923 A primary consideration when selecting a QM/MM method ⁹⁷⁷
 924 is the interactions at the QM/MM interface. Two aspects must ⁹⁷⁸
 925 ⁹⁷⁹

be considered; i) the presence of covalent bonds across the ⁹²⁵
 interface – a particular concern for large (*e.g.*, biomolecular) ⁹²⁶
 molecules, ii) the influence of the MM solvent region on the ⁹²⁷
 QM region – electrostatic and van der Waals interaction terms ⁹²⁸
 must be included. ⁹²⁹

In order to treat covalent bonds at the interface, it is possible ⁹³⁰
 to introduce “link atoms”. Link atoms are QM hydrogen atoms ⁹³¹
 that fill free valencies of QM atoms connected to MM atoms. A ⁹³²
 disadvantage of this method is the debate about inclusion of ⁹³³
 Coulombic interaction terms for the link atoms. Other methods ⁹³⁴
 developed in order to avoid the use of link atoms include the ⁹³⁵
 Local Self-Consistent Field (LSCF) method, which applies a ⁹³⁶
 mixture of hybrid and atomic orbitals to represent the QM ⁹³⁷
 system, and the “connection atom” method, where MM and ⁹³⁸
 QM interface atoms are described as QM methyl groups with a ⁹³⁹
 free sp^3 valence. ⁹⁴⁰

A recent three-layered approach aiming to tackle the issues ⁹⁴¹
 associated with the QM/MM interface and the interaction terms ⁹⁴²
 for MM solvent effects has been proposed by Steindal *et al.*⁵³. ⁹⁴³
 This approach is described as the fully polarisable ⁹⁴⁴
 QM/MM/PCM method (see section 3 for a description of ⁹⁴⁵
 PCM), and is designed for the effective inclusion of a medium ⁹⁴⁶
 in a QM calculation. Short range solvent electrostatic potentials ⁹⁴⁷
 are described by an atomistic model (QM/MM) whilst the long ⁹⁴⁸
 range potentials are described by a continuum. The method is ⁹⁴⁹
 implemented in combination with linear response techniques ⁹⁵⁰
 with a non-equilibrium formulation of environmental response. ⁹⁵¹
 The authors find a faster convergence with respect to system ⁹⁵²
 size for QM/MM/PCM than for QM/MM methods. This ⁹⁵³
 approach allows for reduction of the MM part of the calculation ⁹⁵⁴
 with PCM, allowing less demanding calculations, and reduced ⁹⁵⁵
 sampling. However, three-layered approaches such as this often ⁹⁵⁶
 require much more user input and method manipulation, for ⁹⁵⁷
 example, considerations for MM/PCM interactions have to be ⁹⁵⁸
 considered in addition to QM/MM interactions, and so such ⁹⁵⁹
 methods are suited only to advanced users. ⁹⁶⁰

4.3 Explicit Representations of Water Atoms

When solvent is represented explicitly, solvent molecules ⁹⁶¹
 usually greatly outnumber solute molecules. Thus, in order for ⁹⁶²
 a model to be efficient, it is advantageous to use the simplest ⁹⁶³
 possible solvent representation.⁴⁴ Water is often considered the ⁹⁶⁴
 most useful solvent system, and thus is the solvent most widely ⁹⁶⁵
 used in explicit solvent models. The macroscopic properties are ⁹⁶⁶
 well established, yet the microscopic forces that determine ⁹⁶⁷
 water structure are not fully understood. ⁹⁶⁸

The treatment of water can be rigid or flexible. Rigid ⁹⁶⁹
 models often include a fictitious H-H bond to constrain bond ⁹⁷⁰
 angles in the water monomer.³ Three of the most common rigid ⁹⁷¹
 models for water are the TIP3P (transferable intermolecular ⁹⁷²
 potential 3P), SPC (simple point charge) and SPC/E (simple ⁹⁷³
 point charge extended) models, and their modified counterparts. ⁹⁷⁴
 These three models are effectively rigid pair potentials ⁹⁷⁵
 comprising LJ and Coulombic terms. However, the terms used ⁹⁷⁶
 differ in each model, and give rise to different calculated bulk ⁹⁷⁷
 properties for water.⁵⁴ Values for various properties of water ⁹⁷⁸
 979

980 obtained with different rigid models of water are shown below
 981 in table 1.

982 **Table 1** – Model vs. experimental (Exp.) values for bulk properties of water under
 983 standard conditions (298K; 1 bar), including dipole μ , density ρ , static dielectric
 984 constant ϵ_0 and heat capacity C_p .

Property	TIP3P ^{55,56}	TIP4PEw ⁵⁷	SPCE ^{58,56}	Exp.
μ (D)	2.348	2.32	2.352	2.51
ρ (g/cm ³)	0.980	0.995	0.994	0.997
ϵ_0	94	63.90	68	78.4
C_p (cal/(K.mol))	18.74	19.2	20.7	18

985

986 MD calculations require the integration of Newton's
 987 equations of motion for all atoms, which is achieved through
 988 the evaluation of all atomic forces at each time step. Non-
 989 bonded interactions, especially long-range electrostatic
 990 interactions, dominate computationally, requiring extensive
 991 CPU time. In order to minimise this to an acceptable level,
 992 approximations are necessary. Boundaries are introduced in
 993 water models to restrain the system to a finite size, which
 994 almost always leads to artefacts in the obtainable data.⁵⁴ The
 995 most commonly utilised method for cost-effective simulation
 996 computations is the application of a spherical cut-off, limiting
 997 the number of pairwise interactions to those within a specified
 998 radius.⁵⁴ The use of cut-offs for non-bonded interactions
 999 have undesirable effects. LJ interactions are susceptible to
 1000 small energetic effects, and large pressure effects induced by
 1001 cut-offs. Pressure scaling can be used to correct for pressure
 1002 related cut-off effects, usually to the order of several hundred
 1003 bar. Cut-off effects for systems with dipolar electrostatic
 1004 interactions are more prominent, with cut-offs selected with
 1005 the parameters of experimental radial distribution functions
 1006 to ~ 1.0 nm. However, computer simulations have shown
 1007 ordering within water up to ~ 1.4 nm, so the full structure
 1008 of water is not typically accounted for, resulting in a poor
 1009 description of dielectric properties. A further, and the most
 1010 prominent, effect of cut-offs occurs in systems with net
 1011 charges, where accumulation of the charge occurs at the cut-off
 1012 boundary.⁵⁹

1013 Spoel *et al.*⁵⁹ (1998) investigated the effectiveness of
 1014 TIP3P, TIP4P, SPC, and SPC/E models in describing the
 1015 density and energy, dynamic, dielectric and structural
 1016 properties of water. All simulations and analyses were identical
 1017 for each model investigated, allowing the evaluation of
 1018 simulation methodology independent of the model. It was
 1019 found that system size, cut-off length and reaction fields
 1020 have comparable effects on the overall calculated structural
 1021 properties of water.

1022 System size effects are considered through the comparison
 1023 of systems comprising a small (216) and a large (820) number
 1024 of molecules. The average thermodynamic properties (ρ , E_{pot}

P) are the same regardless of system size. Fluctuations in
 thermodynamic properties are known to be proportional to the
 square root of the system size, which is confirmed within the
 study. However, differences between large and small systems
 are observed, particularly for the dielectric constant, which is
 higher for all systems with a large number of molecules. The
 diffusion constant for large systems is also higher, attributed to
 periodic boundary conditions (PBC).

Cutoff effects are considered by the use of two different
 cutoff lengths (0.9 nm and 1.2 nm) for the large systems. It is
 found that density increases with an increased cutoff length,
 and energy decreases. There is no effect on dielectric
 behaviour.

In all simulations density decreased by approximately 1 kJ
 mol⁻¹ on application of a reaction field. The self-diffusion
 constant D , and rotational correlation times were found to
 increase, indicating the reaction field affects on both the
 translational and rotational mobility of molecules.

Quantum chemical MD simulations of water are often
 developed with Density Functional Theory (DFT) methods,
 applied with a plane wave basis set to determine the electronic
 structure and forces. These methods offer reasonable estimates
 of the structural and dynamic properties of water when
 compared to experimental measurements. However, problems
 exist in the description of electronic gradient corrections, and
 equilibrium pressure. The interatomic forces of early quantum
 simulations, including DFT based methods, were originally
 parameterised with classical mechanics, leading to an
 unsatisfactory agreement between quantum and experimental
 results. DFT models also tend to calculate liquid structure with
 too much order, and underestimate equilibrium density. This is
 often attributed to the inability of local functionals to describe
 dispersion effects.

A recent approach to water simulation has claimed to
 provide a model, called the electronically coarse-grained
 model, capable of accounting for the shortcomings of both
 existing classical and quantum models.⁶⁰ Jones *et al.*⁶⁰ (2013)
 base their method on the replacement of valence electrons of an
 atom with an embedded Quantum Drude oscillator (QDO).
 QDO treatment of water is based upon the TIP4P classical rigid
 model of water, with the three water atoms supplemented by a
 dummy atom with a negative charge, added along the \angle HOH
 bisector to create an additional interaction point. The QDO
 parameters aim to reproduce the isotropic parts of the dipole,
 polarisability, and the dispersion coefficient. The dispersion
 interaction is then adjusted by scaling, whilst preserving
 polarisability. The baseline unadjusted model produces a
 realistic, but over-structured liquid with a density that is too
 low by up to 20%, attributed to its underestimation of
 dispersion. Note also that the value of the enthalpy of
 vaporisation (at ambient pressure) Δh_{vap} was found at 40 ± 2
 kJ/mol, close to the experimental value of 43.91 kJ/mol.
 Scaling the dispersion term results in an increased equilibrium
 density for increased dispersion. This induces a weakening
 effect on the H-bonding network of water, bringing the overall
 structure closer to agreement with benchmark data. However,

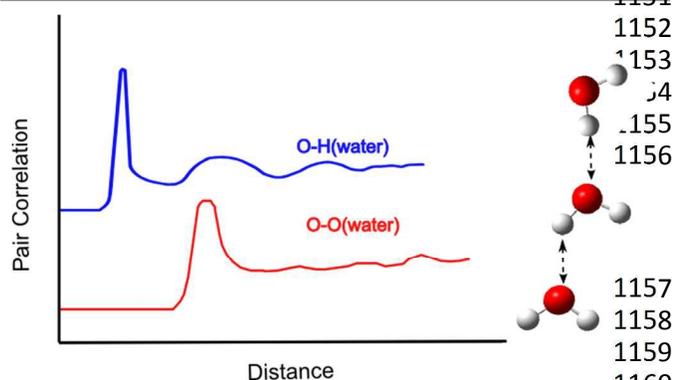
1081 the calculated Δh_{vap} increases to 46 ± 2 kJ/mol, which is
 1082 higher than the experimental value. It is also found that the
 1083 bond network is sensitive to changing polarisation at
 1084 dispersion, affirming the independent importance of
 1085 polarisation and dispersion effects on an overall explicit model.

1086 5. Efficient Hybrid Models – Statistical Mechanics

1087 Within an aqueous solution phase, single snapshot images
 1088 of structure are of limited use. Water is one of the few single
 1089 component liquids for which there are highly competitive
 1090 interactions at short range (hydrogen bonding), capable of
 1091 damping the effects of repulsion. For this reason, ensemble
 1092 averaging is required to identify the most probable geometric
 1093 configurations which most heavily contribute to the system's
 1094 interactions. This idea has already been introduced with
 1095 explicit models of solvation using ensembles taking snapshots
 1096 at specific time periods. However, the cost of calculating the
 1097 many configurations accessible in a solution is enormous,
 1098 hence, in this section we focus on statistical mechanics methods
 1099 which enable a more efficient calculation process.

1100 5.1 Correlation Functions

1101 From a chemical point of view, a solution is a highly mobile
 1102 system in which the dynamics are a vital contribution to the
 1103 system's properties and behaviour. Therefore, mathematically,
 1104 we wish to capture this. Attempting to quantify dynamics with
 1105 static properties is not sufficient; we must therefore provide
 1106 averages or probabilities of interactions occurring at given
 1107 distances. For this reason a natural choice is to represent the
 1108 solvent using Pair Correlation Functions (PCF), or equivalently
 1109 Radial Distribution Functions (RDF). These functions allow us
 1110 to determine a probabilistic structure of the solvent.



1111 Fig. 6 - A schematic representation of PCF for liquid water; water oxygen – water
 1112 hydrogen (blue) and water oxygen – water oxygen (red).
 1113
 1114

1115 PCF can be interpreted as showing the probability again
 1116 distance of there being an atom of interest at that distance
 1117 the atom under study. For example the first large blue peak
 1118 Figure 6 would correspond to either a water H at a distance
 1119 from an O atom under study or *vice versa*. These functions are
 1120 experimentally determinable from scattering experiments. We
 1121 would expect that the PCF/RDF would go to a constant value of

1 at large values of r (i.e. it would become isotropic, like a continuum model, as there are no solute interactions to perturb the system). However, at small values of r we would not expect this. At very small values (less than the van der Waals radii of the solute atoms) we expect zero as only one particle can occupy the space at a time. Just outside this distance we see sharp non-uniform behaviour as solvent in the space interacts favourably with the solute holding a more rigid form. This leads to troughs in the PCF/RDF just behind the peaks, thus deviating from the value of 1 for a uniform solvent.

5.1.1 COMPUTATIONAL USE AND DETERMINATION OF CORRELATION FUNCTIONS

The starting point for the use and determination of these functions for solvation modelling in statistical mechanics is integral equation theory (IET). In this theory a molecule is fully described by a six-dimensional vector (three degrees of freedom relate to position x,y,z and three degrees of freedom determine the orientation ψ,θ,ϕ). To refer to these two sets of variables collectively, we will use the following symbols $r=\{x,y,z\}$ and $\Theta=\{\psi,\theta,\phi\}$. These variables are conveniently incorporated into the fundamental 6D integral equation, the *Molecular Ornstein-Zernike* equation (MOZ). This equation utilises PCF/RDF between the various constituents of the liquid, $g(r_1,r_2,\Theta_1,\Theta_2)$. This simplifies for homogeneous solution to relative positions and orientation of the constituents, $g(r_1 - r_2, \Theta_1 - \Theta_2)$. This can most conveniently be written with reference to the total correlation function $h(r, \Theta)$.⁶¹

$$h_{ij}(r_1 - r_2, \theta_1 - \theta_2) = g_{ij}(r_1 - r_2, \theta_1 - \theta_2) - 1 \quad (22.)$$

We can simplify this equation by assuming spherical symmetry of molecules, hence removing consideration of orientational degrees of freedom by treating each water molecule as a hard sphere. We can now further separate the contributions to the total correlation function into direct and indirect components. To do this we must introduce the direct correlation function $c(r)$. We can now re-write the MOZ equation assuming spherical symmetry as follows:

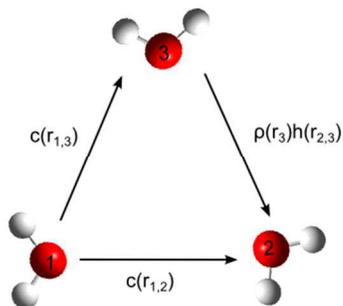
$$h(r_{1,2}) = c(r_{1,2}) + \int dr_3 c(r_{1,3})\rho(r_3)h(r_{2,3}) \quad (23.)$$

Two effects contribute to the total correlation function (eqn. 22); i) the direct correlation between r_1 and r_2 , and ii) an indirect correlation via a third body, r_3 . The indirect correlation via r_3 is weighted by the density at r_3 , and thus allows the consideration of all possible positions of the third body.⁶¹

To solve this equation, $h(r)$ and $c(r)$ need to be found. As we have only a single equation and two unknown functions, $h(r)$ and $c(r)$, another equation is required; a closure relation must be introduced. There are several such equations available from statistical mechanics. The exact closure relation is as follows:

$$g(r) = e^{-\beta U(r)+h(r)-c(r)+B(r)} \Rightarrow e^{-\beta U(r)+T(r)+B(r)} \quad (24.)$$

1168
 1169 where β is equal to $1/k_B T$ and $U(r)$ is the interaction potential
 1170 which is often of the following form:



1171
 1172 Fig. 7 – Illustration of the contributions, both direct and indirect, to the total
 1173 correlation function.

$$U(r) = 4\epsilon \left[\left(\frac{\sigma_{ab}}{r} \right)^{12} - \left(\frac{\sigma_{ab}}{r} \right)^6 \right] + \frac{q_a q_b}{r} \quad (25.)$$

1174 $T(r)$ is known as the indirect correlation function as it is the
 1175 difference between the total and direct correlation functions,
 1176 and quantifies the indirect contribution. $B(r)$ is the bridge
 1177 function, which comes from graph theory - its exact form is not
 1178 known. Several approximate closure relations exist; some of them
 1179 be discussed here, although others are available. Originally the
 1180 *HyperNetted-Chain* (HNC) approximate closure was used.

$$h(r) = e^{(-\beta U(r) + T(r))} - 1$$

1181 This closure works in principle for charged systems but
 1182 neglects the bridge function term completely, assuming it to be
 1183 zero. This can lead to poor convergence due to uncontrolled
 1184 growth in the argument of the exponent. An alternative is the
 1185 Partially Linearised HyperNetted Chain (PLHNC). This closure
 1186 linearises the HNC once a cut off value (C) is exceeded:

$$\Lambda = -\beta U(r) + T(r) \quad (27.)$$

$$h(r) = \begin{cases} e^{(-\beta U(r) + T(r))} - 1 & \text{When } \Lambda \leq C \\ -\beta U(r) + T(r) + e^C - C - 1 & \text{When } \Lambda > C \end{cases}$$

1187 This improves the convergence of the equations and is
 1188 regularly used in many applications for a variety of systems.

1189 Due to the spherical symmetry approximation, the MOA can
 1190 only be applied to simple solutions. Additionally, due to the
 1191 high dimensionality of the full equation, before the spherical
 1192 symmetry approximation was invoked, it is practically
 1193 incomputable. For this reason a number of approximations
 1194 have been developed which are collectively referred to as *Reference
 1195 Interaction Site Models* (RISM).

1196 5.2 3D-RISM: A Hybrid Solvation Model

1197 As we have seen, the explicit treatment of solvent
 1198 considered to be a necessary step in the understanding of
 1199 solvent structure. However, this naturally carries high
 1200 computational costs³. The alternative continuum treatment of

solvents lacks the ability to account for the underlying physical
 theory; energy contributions from solvation shell features are
 computable, but not transferable. Solvent structure features
 from the first and second solvation shells are lost in continuum
 models, and non-electrostatic energy terms are not described
 from first principles, thus are not transferable to more complex
 models.⁶³

The 3D derivation of RISM (3D-RISM)^{64,65} is a 3D
 molecular theory of solvation, applied through solvent
 distributions, rather than explicit solvent molecules, and
 conceives solvation structure and dynamics from the first
 principles of statistical mechanics.

3D-RISM is derived from a partial integration over the
 orientational degrees of freedom; this leaves a set of 3D integral
 equations (one equation per solvent site; N_{solvent}). This method
 utilises solvent site – solute total correlation functions and
 direct correlation functions in the solution of the RISM
 equations. The 3D-RISM equations take the following form:⁶²

$$h(\alpha) = \sum_{\xi} \int_{R^3} c_{\xi}(r_1 - r_2) \chi_{\xi, \alpha}(|r_2|) dr_2 \quad (28.)$$

Here $\chi_{\xi, \alpha}$ labels the solvent susceptibility function. This
 function models the bulk solvent mutual correlations. For the
 example of water, this function models the intermolecular
 correlation between water oxygen and water hydrogen. This
 function can be calculated from the intramolecular solvent
 correlation function ($\omega_{\zeta\gamma}^{\text{solvent}}(r)$), the radial site to site total
 correlation functions ($h_{\zeta\alpha}^{\text{solvent}}(r)$) and the number density at
 each solvent site (ρ_{α}):

$$\chi_{\xi, \alpha}(r) = \omega_{\zeta\gamma}^{\text{solvent}}(r) + \rho_{\alpha} (h_{\zeta\alpha}^{\text{solvent}}(r)) \quad (29.)$$

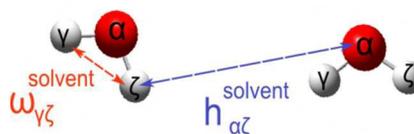


Fig. 8 – Illustration of the contributions to the solvent susceptibility function.

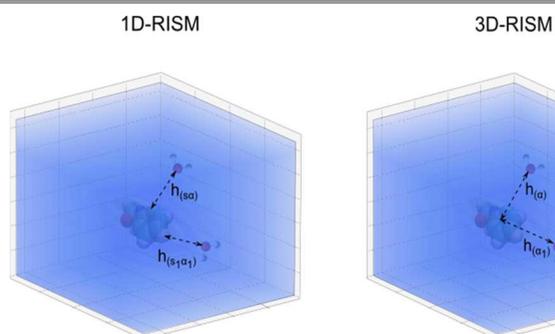
3D RISM can reliably account for the spatial correlation of the
 solvent density around the solute. As displayed above, the
 solvent molecules are modelled as a set of atomic sites, with 3D
 structure described by intramolecular correlation functions.^{62,66}

5.3 1D-RISM: A High Throughput Solvation Model

Another RISM method is 1D RISM, which separates the
 solute into a set of sites (generally the atoms) and utilises
 solvent site – solute site total correlation functions and direct
 correlation functions. This leads to a set of ($N_{\text{solute site}} \times N_{\text{solvent
 site}}$) closure relations. 1D RISM is extremely quick to calculate
 but does not account properly for spatial correlations of the
 solvent density around the solute:

$$h_{s'\alpha}(r) = \sum_{s'=1}^{N_{\text{solute}}} \sum_{\zeta=1}^{N_{\text{solvent}}} \int_{R^3} \int_{R^3} \omega_{ss'}(|r_1 - r'|) c_{s'\zeta}(|r' - r''|) \chi_{\xi,\alpha}(|r'' - r_2|) dr' dr''$$

1241 N_{solute} is the number of sites in the solute and N_{solvent} is number
 1242 of sites in the solvent molecule. $\omega_{ss'}$ are the intramolecular
 1243 correlation functions representing the solute molecule.⁶⁶
 1244 Implementations of both 1D- and 3D-RISM are available in
 1245 well-known computational packages such as AMBER. There
 1246 are also implementations in some quantum chemistry codes
 1247 such as ADF.



1248
 1249 **Fig. 9** - A schematic representation of 1D-RISM and 3D-RISM. The conceptual
 1250 difference in the models is that the total correlation functions are calculated
 1251 considering the solute as a set of sites (1D-RISM) or as a single site (3D-RISM). α
 1252 labels the solvent site in both models, s labels the solute site in the 1D-RISM
 1253 case.

1254 5.4 RISM Corrections and Derivations

1255 5.4.1 CORRECTION SCHEMES

1256 A well-known error in both 1D and 3D-RISM occurs due
 1257 to accounting for the cavitation term in the solution phase
 1258 incorrectly. Other limitations also exist, associated with the use
 1259 of approximations. Several schemes to correct these errors have
 1260 been developed for 3D-RISM, however these are beyond the
 1261 scope of this review, and thus are discussed in minimal detail.

1262 Many studies have been conducted over the last two
 1263 decades with a view to improving the accuracy of 3D-RISM for
 1264 a variety of applications. Modifications to the original
 1265 equations have included cavity corrections,⁶⁷ parallelisation
 1266 with fast Fourier transforms⁶⁸ and MD modifications,⁶⁷
 1267 amongst others.

1268 The universal correction (UC)⁶⁹ given in equation 25
 1269 is a two parameter correction derived by regression. $\Delta G_{\text{hydration}}^{GF}$
 1270 refers to the Gaussian fluctuation hydration free energy (HFE)
 1271 functional discussed below, a and b are regression coefficients
 1272 ($a = -3.2217$ and $b = 0.5783$), and ρV is the dimensionless partial
 1273 molar volume as calculated by 3D-RISM.

$$\Delta G_{\text{hydration}}^{3D-RISMUC} = \Delta G_{\text{hydration}}^{GF} + a(\rho V) + b \quad (31)$$

$$UC = a(\rho V) + b$$

1275 A second scheme known as cavity corrected 3D-RISM
 1276 single parameter calculated on the basis of a solution compo-

of spheres which interact exclusively by Lennard-Jones type
 interactions.⁷⁰ A very recent addition offers a theoretical
 justification for such schemes; applying a Thermodynamic-
 Ensemble Partial Molar Volume Correction.⁷¹

Correction schemes for 1D-RISM also exist. These
 correction schemes must correct for additional approximations
 from the 1D RISM theory. A recent addition is the Structural
 Descriptor Correction (SDC). This applies QSPR methods and
 group contributions to correct 1D-RISM.⁶⁶

A primary concern in the improvement of 3D-RISM
 remains its ability to describe the thermodynamic properties of
 solvation. One view adopted by Palmer *et al.*⁶⁹ is that solubility
 calculations should be considered in terms of a simple
 thermodynamic cycle, calculating the solvation free energy
 from summation of the free energy of sublimation, and the free
 energy of hydration, as illustrated in Fig. 10.

A recent investigation by Palmer *et al.*⁶² implements the
 thermodynamic cycle approach to the calculation of solubility,
 with sublimation free energies calculated from crystal lattice
 minimisation and HFEs calculated with 3D-RISM. Crystal
 lattice calculations are performed on known crystal structures.

The authors highlight a plethora of existing approximate
 functionals which can provide HFE values from the solvent
 site-solute total correlation functions and direct correlations of
 3D-RISM. However, the functionals investigated previously to
 Palmer *et al.*'s work often provide HFEs with RMSE errors
 higher than the standard deviation of experimental data, and
 worse than those reported in QSPR models.

The investigation⁶² implementing the thermodynamic cycle
 approach to the calculation of solubility applied the previous
 work of Palmer *et al.*⁶¹ and found that the thermodynamic cycle
 approach predicted HFEs in good agreement with experiment
 ($R = 0.94$, $\sigma = 0.99$ kcal mol⁻¹). However, the predictions did
 not perform as well as purely empirical approaches, and this
 was mostly attributed to a lack of parameterisation against
 experimental data.

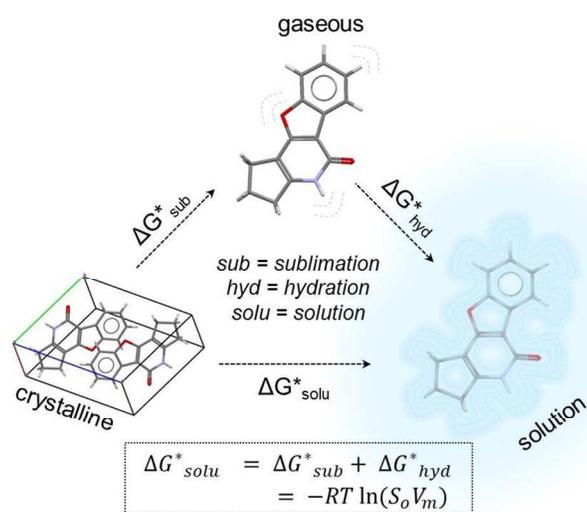
5.4.2 HYDRATION FREE ENERGY FUNCTIONALS

In order to calculate HFEs a HFE functional must be
 applied to the RISM output. There are a number of such
 functionals which vary in accuracy. Some of the correction
 schemes above recommend a specific HFE functional for use
 (UC recommends the Gaussian fluctuation HFE functional⁷²). It
 is suggested to the user that where possible several functionals
 are tested for accuracy. Where this is not possible, the guidance
 given for the selection of a HFE functional for specific schemes
 should be followed, as these are generally well documented by
 the developer groups.

5.4.3. RISM AND QUANTUM CHEMICAL APPLICATIONS

RISM has also been applied to quantum chemical
 applications. RISM was extended for applications to quantum
 chemistry - this extension is called RISM-SCF. This theory
 provides the following definition of the Helmholtz free energy
 of the system:

1330
 1331 $A = E_{\text{solute}} + \Delta\mu$
 1332
 1333 where A is the total Helmholtz free energy, E_{solute} is the solute
 1334 energy and $\Delta\mu$ is the solvation free energy from the RISM
 1335 equations. A is functionally connected to both the site-to-site
 1336 density correlation functions and the wavefunction of the
 1337 solute, hence mutual solution of E_{solute} and $\Delta\mu$ provide the
 1338 system's equilibrium energies.³³ 3D-RISM has been combined
 1339 with Kohn-Sham DFT, offering an alternative to continuum
 1340 solvents and *ab initio* MD.⁷³ These calculations have been
 1341 extended to higher levels of quantum mechanical theory (including
 1342 reference methods) which are currently unaffordable at the
 1343 QM/MM level.³³



1344
 1345 **Fig. 10** – Solubility prediction via a thermodynamic cycle. The free energy change
 1346 from crystalline to aqueous phase is calculated from the summation of the free
 1347 energy change of sublimation and the free energy change of hydration.

1348 5.4 Other Hybrid Models

1349 Combined implicit/explicit hybrid models work on a
 1350 common framework; the central part of the system contains
 1351 explicit solute and a few explicit solvent molecules, and the rest
 1352 of the system is treated as a dielectric continuum.

1353 The improvement associated with the insertion of explicit
 1354 water molecules within a dielectric continuum has been
 1355 demonstrated by Kelly, Cramer and Truhlar,⁷⁴ who use the
 1356 calculation of aqueous acid dissociation constants to
 1357 demonstrate the effects of inserting a single explicit solvent
 1358 molecule into a continuum solvent representation. Along with
 1359 previous work,⁷⁵ the authors show that in many cases an
 1360 implicit solvation method is sufficient for the calculation of pK_a
 1361 values. However, when strong and specific solute-solvent
 1362 hydrogen bonding interactions are expected to contribute
 1363 significantly to the aqueous phase, a single explicit molecule
 1364 inserted to the continuum significantly improves the
 1365 calculation. Using their own implicit continuum model (SM6),
 1366 it is found that addition of further explicit waters, up to three

significantly increases the accuracy of the calculation.
 However, the use of alternative continuum models, namely
 SM5.43R and PCM, finds a worsening of results when an
 increasing number of explicit atoms are added. This
 exemplifies the importance of choosing a suitable continuum
 representation in implicit/explicit hybrid models.

Zhu and Krilov⁷⁶ discussed two flexible boundary hybrid
 solvation models for biomolecular systems, based upon the
 traditional hybrid model with both explicit and implicit solvent
 regions. The proposed models aim to account for short-range
 solvent effects *via* elimination of PBC by limiting the number
 of explicit solvent molecules to two or three solvation shells.
 The first model, the dynamic boundary model, imposes a
 confining potential on the solvent, which responds dynamically
 to fluctuations in solvent distribution and solute conformation.
 The second model, the exchange boundary solvation model,
 allows pairwise exchanges between the explicit and implicit
 regions of the system, maintaining a uniform hydration of the
 solute. Comparison of the two methods with traditional PBC
 methods shows good agreement between calculated energies,
 and the two models are found to improve computational
 efficiency by up to two orders of magnitude, attributed to the
 reduced number of explicit solvent molecules in comparison to
 other models.

Chaudhury *et al.*⁷⁷ recently discussed the discrepancies
 between explicit and implicit methods for solvation models of
 biological systems such as proteins, and consequently
 investigate a Hybrid Replica Exchange Molecular Dynamics
 (REMD) method for protein solvation. Temperature-based
 REMD involves running multiple simultaneous simulations at a
 wide-range of temperatures, while allowing temperature
 exchange between simulation steps. This relates the relative
 probability of finding each conformation at a given temperature
 to conformational energy. Traditional REMD successfully
 models small peptides and proteins, but becomes more cost-
 constrained for larger systems. In order to account for
 discrepancies between implicit and explicit methods, the
 authors propose a hybrid implicit/explicit method with each
 simulation step run exclusively in explicit solvent. During
 exchange between time steps, the entire solvent system is
 replaced with an implicit solvent model. Finally, the explicit
 solvent is re-inserted for the next simulation step. The use of an
 implicit solvent model during exchange significantly reduces
 computational cost. Where implicit and explicit models give
 different behaviours, the hybrid method gives mixed results in
 terms of thermodynamic and structural descriptions. However,
 the explicit model of solvent molecules describes solvent-
 specific features of energy landscapes well.

A further emerging method that similarly attempts to reduce
 the cost-constraints of explicit methods is Grid Cell Theory
 (GCT).⁷⁸ GCT spatially resolves the enthalpic and entropic
 components of hydration on a 3D grid, covering a volume of
 space around a solute. The grid can be non-uniform and
 unevenly spaced. The solute is constrained to adopt a single
 conformation, speeding up convergence by only allowing rigid
 body translations and rotations of water molecules. A second

1423 benefit of GCT is that graphical analysis of a calculated grid is
 1424 possible. A drawback of GCT method development email is
 1425 from the fact that there does not exist a unique method for
 1426 partitioning a free energy into a sum of contributions. These
 1427 contributions are susceptible to coupling. Gerogiokas *et al.*
 1428 have recently proposed a GCT method, and evaluated the
 1429 enthalpic and entropic contributions to hydration, making
 1430 visualisation of hydration thermodynamics possible. GCT is a
 1431 slower method than other thermodynamic integration methods,
 1432 but such alternative methods are not as descriptive in terms of
 1433 thermodynamic contributions.

1434 6. Outlook and Conclusions

1435 The aim of this review is to introduce the multitude of
 1436 available methods and concepts for the calculation of solution
 1437 free energies, and the modelling of systems in solution.
 1438 Through the highlighting of many traditional and emerging
 1439 methods within explicit, implicit, informatics and hybrid
 1440 methods, it has become clear that each modelling category has
 1441 its own advantages and disadvantages. The trade-off between
 1442 the inaccuracies of implicit solvent models and the
 1443 computational cost-constraints of explicit models are a
 1444 prominent issue, and have conceived a number of hybrid
 1445 solvation methods, each of which aims to provide a model of
 1446 reasonable accuracy at an appropriate cost. The plethora of such
 1447 available methods exemplifies the importance of accurate
 1448 solvation models.

1449 We have placed particular emphasis on 3D-RISM and its
 1450 derived counterparts, as we believe that RISM based methods
 1451 are a strong contender in the challenge of finding a
 1452 computationally viable solubility prediction method which is
 1453 also descriptive enough for the theoretical study of a system's
 1454 thermodynamics. However, it is also noted that such methods
 1455 are a long way from perfection, and require further refinements
 1456 of solute-solvent correlation functions.

1457 With the increase of computing power, as described by
 1458 Moore's law, it is hard to predict how much of an issue
 1459 computational costs associated with solvation modelling will be
 1460 over the coming years. However, increases in computing power
 1461 will inevitably allow more accurate methods to be employed
 1462 within a faster timeframe. We predict the emergence of hybrid
 1463 models which describe the theoretical and physical components
 1464 of solvation at an ever increasing rate, with the need to trade off
 1465 accuracy over time becoming less as computing power
 1466 increases.

1467 Although future prospects for solvation modelling are
 1468 bright, we are also aware that there is a very present need for
 1469 good models. We would like to note that the best choice of
 1470 model for solvation is entirely dependent on the requirements
 1471 of the user. For high-throughput screening of molecules of
 1472 similar structural features, we suggest QSPR/QSAR as a
 1473 suitable and reliable approach for thermodynamic property
 1474 calculation (*e.g.*, solvation free energy). However, where
 1475 specific physical and mechanistic meaning is desired, it is best
 1476 to employ either explicit solvent representations, suitable for

relatively small solute sizes, or where larger solutes are used,
 hybrid models. The choice of hybrid models for such
 investigations is not intuitively obvious, as highlighted within
 this review, as some systems are described sufficiently with
 addition of a single solute molecule, whereas for other systems
 it is necessary to add enough explicit solvent molecules to
 describe full solvation shells. Thus, it is often necessary to
 consider whether solvent behaviour is a significant contributor
 to the property of interest. If so, explicit/hybrid methods are
 advisable, dependent upon available computing resources.
 Otherwise, continuum models could offer sufficient physical
 description of the solvent environment. Of course, where
 sufficient and trustworthy experimental data are available,
 several models should be tested and evaluated for correlation
 with available experimental data.

Acknowledgements

We are grateful for useful discussions with colleagues including
 the groups of Professor Maxim Fedorov and Dr David Palmer.
 JLMcD and JBOM are grateful to SULSA for funding; RES
 and JBOM thank the University of St Andrews, EPSRC (grant
 EP/L505079/1) and CCDC for funding.

Notes and references

^a School of Chemistry, University of St Andrews, Purdie Building, North
 Haugh, St Andrews, Fife, KY16 9ST, UK.

^b Cambridge Crystallographic Data Centre, 12 Union Road, Cambridge
 CB2 1EZ, UK.

[†]Signifies that these authors contributed equally to this work.

A recent machine learning method and dataset proposed by some of the
 authors is available from the Mitchell group web server:

http://chemistry.st-andrews.ac.uk/staff/jbom/group/Informatics_Solubility.html

1. S. Basavaraj and G. V. Betageri, *Acta Pharm. Sin. B*, 2014, **4**, 3–17.
2. H. D. Williams, N. L. Trevaskis, S. A. Charman, R. M. Shanker, W. N. Charman, C. W. Pouton, and C. J. H. Porter, *Pharmacol. Rev.*, 2013, **65**, 315–499.
3. C. J. Cramer, *Essentials of Computational Chemistry: Theories and Models*, John Wiley & Sons, Chichester, 2013.
4. A. Llinàs, R. C. Glen, and J. M. Goodman, *J. Chem. Inf. Model.*, 2008, **48**, 1289–1303.
5. S. H. Yalkowsky and S. C. Valvani, *J. Pharm. Sci.*, 1980, **69**, 912–922.
6. A. G. Leach, H. D. Jones, D. A. Cosgrove, P. W. Kenny, L. Ruston, P. MacFaul, J. M. Wood, N. Colclough, and B. Law, *J. Med. Chem.*, 2006, **49**, 6672–6682.
7. MedChemica, 2014.
8. J. Hussain, *GlaxoSmithKline*, 2011.
9. J. W. Ponder, C. Wu, P. Ren, V. S. Pande, J. D. Chodera, M. J. Schnieders, I. Haque, D. L. Mobley, D. S. Lambrecht, R. A. DiStasio, M. Head-Gordon, G. N. I. Clark, M. E. Johnson, and T. Head-Gordon, *J. Phys. Chem. B*, 2010, **114**, 2549–64.
10. S. Y. Liem and P. L. a Popelier, *Phys. Chem. Chem. Phys.*, 2014, **16**, 4122–34.

- 1532 11. A. A. Noyes and W. R. Whitney, *J. Am. Chem. Soc.*, 1897, **19**, 1587
1533 934. 1588
- 1534 12. A. Jouyban and M. A. A. Fakhree, *Toxicity and Drug Testing*,
1535 InTech, 2012. 1590
- 1536 13. G. Völgyi, E. Baka, K. J. Box, J. E. a Comer, and K. Takács-Nórágyi,
1537 *Anal. Chim. Acta*, 2010, **673**, 40–6. 1591
- 1538 14. J. L. McDonagh, N. Nath, L. De Ferrari, T. van Mourik, and J. Mitchell,
1539 *J. Chem. Inf. Model.*, 2014, **54**, 844–856. 1593
- 1540 15. B. A. Hendriksen, M. V. S. Felix, and M. B. Bolger, *AAPS PharmSciTech*,
1541 2003, **5**, 35–49. 1595
- 1542 16. J. Bauer, S. Spanton, R. Henry, J. Quick, W. Dziki, W. Porter, and
1543 Morris, *Pharm. Res.*, 2001, **18**, 859–866. 1597
- 1544 17. S. R. Chemburkar, J. Bauer, K. Deming, H. Spiwek, K. Patel, R. Henry,
1545 Morris, R. Henry, S. Spanton, W. Dziki, W. Porter, J. Quick, J. Bauer,
1546 J. Donaubauber, B. A. Narayanan, M. Soldani, D. Riley, and J. Mcfarland,
1547 *Org. Process Res. Dev.*, 2000, **4**, 413–417. 1600
- 1548 18. D. S. Palmer, A. Llinàs, I. Morao, G. M. Day, J. M. Goodman, J. B. O. Mitchell,
1549 Glen, and J. B. O. Mitchell, *Mol. Pharm.*, 2008, **5**, 266–279. 1603
- 1550 19. D. S. Palmer, J. L. McDonagh, J. B. O. Mitchell, T. van Mourik, and M. V. Fedorov,
1551 *J. Chem. Theory Comput.*, 2012, 3322–3337. 1605
- 1552 20. S. L. Price, M. Leslie, G. W. A. Welch, M. Habgood, L. S. Price, G. Karamertzanis,
1553 and G. M. Day, *Phys. Chem. Chem. Phys.*, 2010, **12**, 8478–90. 1607
- 1554 21. A. M. Reilly and A. Tkatchenko, *J. Phys. Chem. Lett.*, 2013, **4**, 1033. 1608
- 1555 22. A. R. Leach and V. J. Gillet, *An Introduction to Cheminformatics*,
1556 Springer, Dordrecht, 2007. 1610
- 1557 23. J. B. O. Mitchell, *Wiley Interdiscip. Rev. Comput. Mol. Sci.*, 2011, **1**, 468–481. 1612
- 1558 24. D. Weininger, *J. Chem. Inf. Comput. Sci.*, 1988, **28**, 31–36. 1613
- 1559 25. S. Heller, A. McNaught, S. Stein, D. Tchekhovskoi, and I. Pletnev, *Cheminform.*,
1560 2013, **5**, 7. 1614
- 1561 26. L. D. Hughes, D. S. Palmer, F. Nigsch, and J. B. O. Mitchell, *Chem. Inf. Model.*,
1562 2008, **48**, 220–232. 1615
- 1563 27. A. Luscì, G. Pollastri, and P. Baldi, *J. Chem. Inf. Model.*, 2011, **51**, 1563–1575. 1617
- 1564 28. Y. Ran and S. H. Yalkowsky, *J. Chem. Inf. Comput. Sci.*, 2001, **41**, 354–357. 1619
- 1565 29. Y. Ran, Y. He, G. Yang, J. L. H. Johnson, and S. H. Yalkowsky, *Chemosphere*,
1566 2002, **48**, 487–509. 1620
- 1567 30. J. Ali, P. Camilleri, M. B. Brown, A. J. Hutt, and S. B. Kirrane, *Chem. Inf. Model.*,
1568 2012, **52**, 420–8. 1621
- 1569 31. A. Shayanfar, M. A. A. Fakhree, and A. Jouyban, *J. Drug Deliv. Tech.*, 2010,
1570 **20**, 467–476. 1622
- 1571 32. D. S. Palmer and J. B. O. Mitchell, *Mol. Pharm.*, 2014, **11**, 2972. 1623
- 1572 33. J. Tomasi, B. Mennucci, and R. Cammi, *Chem. Rev.*, 2005, **105**, 2999–3093. 1624
- 1573 34. L. Onsager, *J. Am. Chem. Soc.*, 1936, **58**, 1486–1493. 1625
- 1574 35. R. A. Pierotti, *Chem. Rev.*, 1975, **76**, 717–726. 1626
- 1575 36. S. Höfner and F. Zerbetto, *J. Phys. Chem. A*, 2003, **107**, 11257. 1627
- 1576 37. C. Colominas, F. J. Luque, and M. Orozco, *Chem. Phys.*, 1999, **231**, 253–246. 1628
- 1577 38. J. Wang, T. Hou, and X. Xu, *J. Chem. Inf. Model.*, 2009, **49**, 571–581. 1629
- 1578 39. Y. Takano and K. N. Houk, *J. Chem. Theory Comput.*, 2005, **1**, 70–77. 1630
- 1579 40. A. Klamt and G. Schuurmann, *J. Chem. Soc. Perkin Trans. 2*, 1993, 799. 1631
- 1580 41. A. Klamt, V. Jonas, T. Bürger, and J. C. W. Lohrenz, *J. Phys. Chem. A*,
1581 1998, **102**, 5074–5085. 1632
- 1582 42. A. V. Marenich, C. J. Cramer, and D. G. Truhlar, *J. Phys. Chem. B*,
1583 2009, **113**, 6378–6396. 1633
- 1584 43. S. Rayne and K. Forest, *Available from Nat. Proc.*, 2010. 1634
- 1585 44. R. M. Levy and E. Gallicchio, *Annu. Rev. Phys. Chem.*, 1998, **49**, 531–67. 1635
- 1586 45. D. L. Beveridge and F. M. DiCapua, *Annu. Rev. Biophys. Biophys. Chem.*,
1989, **18**, 431–92. 1636
46. R. W. Zwanzig, *J. Chem. Phys.*, 1954, **22**, 1420. 1637
47. C. Chipot and A. Pohorille, in *Free Energy Calculations: Theory and Applications in Chemistry and Biology*, Springer, Berlin, Heidelberg, 2007, pp. 33–75. 1638
48. K. Lüder, L. Lindfors, J. Westergren, S. Nordholm, and R. Kjellander, *J. Phys. Chem. B*, 2007, **111**, 1883–92. 1639
49. H. Liu, S. Dai, and D. Jiang, *J. Phys. Chem. B*, 2014, **118**, 2719–25. 1640
50. M. A. Wyczałkowski, A. Vitalis, and R. V. Pappu, *J. Phys. Chem. B*, 2010, **114**, 8166–80. 1641
51. A. Ahmed and S. I. Sandler, *J. Chem. Eng. Data*, 2015, **60**, 16–27. 1642
52. R. A. Friesner and V. Guallar, *Annu. Rev. Phys. Chem.*, 2005, **56**, 389–427. 1643
53. A. H. Steindal, K. Ruud, L. Frediani, K. Aidas, and J. Kongsted, *J. Phys. Chem. B*, 2011, **115**, 3027–37. 1644
54. P. Mark and L. Nilsson, *J. Phys. Chem. A*, 2001, **105**, 9954–9960. 1645
55. M. W. Mahoney and W. L. Jorgensen, *J. Chem. Phys.*, 2000, **112**, 8910. 1646
56. C. Vega and J. L. F. Abascal, *Phys. Chem. Chem. Phys.*, 2011, **13**, 19663–88. 1647
57. H. W. Horn, W. C. Swope, J. W. Pitera, J. D. Madura, T. J. Dick, G. L. Hura, and T. Head-Gordon, *J. Chem. Phys.*, 2004, **120**, 9665–78. 1648
58. L. Wang, T. J. Martinez, and V. S. Pande, *J. Phys. Chem. Lett.*, 2014, **5**, 1885–1891. 1649
59. D. van der Spoel, P. J. van Maaren, and H. J. C. Berendsen, *J. Chem. Phys.*, 1998, **108**, 10220. 1650
60. A. Jones, F. Cipcigan, V. P. Sokhan, J. Crain, and G. J. Martyna, *Phys. Rev. Lett.*, 2013, **110**, 227801. 1651
61. D. S. Palmer, A. I. Frolov, E. L. Ratkova, and M. V. Fedorov, *J. Phys. Condens. Matter*, 2010, **22**, 492101. 1652
62. D. S. Palmer, J. L. McDonagh, J. B. O. Mitchell, T. van Mourik, and M. V. Fedorov, *JCTC*, 2012, 3322–3337. 1653
63. T. Luchko, S. Gusarov, D. R. Roe, C. Simmerling, D. A. Case, J. Tuszynski, and A. Kovalenko, *J. Chem. Theory Comput.*, 2010, **6**, 607–624. 1654
64. D. Chandler, J. D. McCoy, and S. J. Singer, *J. Chem. Phys.*, 1986, **85**, 5971. 1655
65. A. Kovalenko and F. Hirata, *J. Chem. Phys.*, 1999, **110**, 10095. 1656
66. E. L. Ratkova and M. V. Fedorov, *J. Chem. Theory Comput.*, 2011, **7**, 1450–1457. 1657
67. J.-F. Truchon, B. M. Pettitt, and P. Labute, *J. Chem. Theory Comput.*, 2014, **10**, 934–941. 1658
68. Y. Maruyama, N. Yoshida, H. Tadano, D. Takahashi, M. Sato, and F. Hirata, *J. Comput. Chem.*, 2014, **35**, 1347–55. 1659

Journal Name

- 1643 69. D. S. Palmer, A. I. Frolov, E. L. Ratkova, and M. V. Fedorov, *Mol.*
1644 *Pharm.*, 2011, **8**, 1423–9.
- 1645 70. J.-F. Truchon, B. M. Pettitt, and P. Labute, *J. Chem. Theory Comput.*,
1646 2014.
- 1647 71. V. P. Sergiievskiy, G. Jeanmairet, M. Levesque, and D. Borgis, *J.*
1648 *Phys. Chem. Lett.*, 2014, **5**, 1935–1942.
- 1649 72. F. Hirata, *Molecular theory of solvation*, Springer, 2003.
- 1650 73. A. Kovalenko and F. Hirata, *J. Mol. Liq.*, 2001, **90**, 215–224.
- 1651 74. C. P. Kelly, C. J. Cramer, and D. G. Truhlar, *J. Phys. Chem. A*, 2006,
1652 **110**, 2493–9.
- 1653 75. C. P. Kelly, C. J. Cramer, and D. G. Truhlar, *J. Chem. Theory.*
1654 *Comput.*, 2005, **1**, 1133–1152.
- 1655 76. W. Zhu and G. Krilov, *J. Mol. Struct. THEOCHEM*, 2008, **864**, 31–
1656 41.
- 1657 77. S. Chaudhury, M. A. Olson, G. Tawa, A. Wallqvist, and M. S. Lee, *J.*
1658 *Chem. Theory Comput.*, 2012, **8**, 677–687.
- 1659 78. G. Gerogiokas, G. Calabro, R. H. Henchman, M. W. Y. Southey, R.
1660 J. Law, and J. Michel, *J. Chem. Theory Comput.*, 2014, **10**, 35–48.

1661