# Protein-Solvent Interactions and Classical Density Functional Theory

by

Eric A Mills

BSc Physics, McMaster University, 2006

MSc Physics, McMaster University, 2008

A THESIS SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENTS FOR THE DEGREE OF

**Doctor of Philosophy**

in

THE FACULTY OF GRADUATE AND POSTDOCTORAL
STUDIES

(Physics)

The University of British Columbia

(Vancouver)

December 2015

# Abstract

We use classical density functional theory to investigate the interactions between solvents and proteins. We examine a diverse experimental literature to establish thermodynamic properties of protein-cosolute interaction, particularly the compensation between transfer entropy and transfer enthalpy. We develop a method of analysing the uncertainties in such measurements and use the method to resolve a long-standing debate over entropy-enthalpy compensation. We develop a classical density functional theory for interactions between proteins and cosolutes. The theory developed here ignores the solvent-solvent interaction but is nonetheless quite accurate. We use this approach to reproduce transfer free energies reported elsewhere, and show that the cDFT model captures the desolvation barrier and the temperature dependence of the transfer free energy. We use experimental values that we have analyzed to define the parameter space of a model density functional theory approach. We then extend the classical density functional theory to capture protein-water interactions, thus developing a new implicit solvent model. Along the way we give a proof that the free energy of a bath of particles in a finite external potential is independent of the external potential in the isothermal-isobaric ensemble. We finally discuss the challenges remaining in implementing our implicit solvent model.

# Preface

Some of the content of this thesis has been published previously or is currently in submission. As my supervisor and I are the sole authors of these papers they represent my original work, prepared with my supervisor's editing, advice, and direction.

- The bulk of chapter 3 was published in "Mills, E A and Plotkin, S S. 'Density functional theory for protein transfer free energy.' *The Journal of Physical Chemistry B*, 117(42):13278-13290, 2013."

- Chapter 2 was published as "Mills, E A and Plotkin, S S. 'Protein transfer free energy obeys entropy-enthalpy compensation.' *The Journal of Physical Chemistry B*, 119(44):14130-14144, 2015."

# Table of Contents

# List of Tables

# List of Figures

# List of Abbreviations

| Abbreviation | Description |
|---|---|
| ABSINTH | Assembly of Biomolecules Studied by an Implicit, Novel, and Tuneable Hamiltonian |
| AFM | Atomic Force Microscope |
| ALS | Amytrophic Lateral Sclerosis |
| AMBER | Assisted Model Build with Energy Refinement |
| ASA | Accessible Surface Area |
| BAR | Bennet Acceptance Ratio |
| cDFT | Classical Density Functional Theory |
| CHARMM | Chemistry at Harvard Macromolecular Mechanics |
| DFT | Density Functional Theory |
| DSC | Differential Scanning Calorimetry |
| ER | Endoplasmic Reticulum |
| GB/SA | Generalized Born/Surface Area |
| GROMACS | GROningen MAchine for Chemical Simulations |
| LJ | Lennard-Jones |
| MD | Molecular Dynamics |
| NMR | Nuclear Magnetic Resonance |
| OPLS | Optimized Potential for Liquid Simulations |
| O-Z | Ornstein-Zernike |
| RISM | Reference Interaction Site Model |
| TI | Thermodynamic Integration |
| UA | United Atom |
| vdW | van der Waals |

WHAM                    Weighted Histogram Analysis Method

# List of Chemical Nomenclature

| Abbreviation | Description |
| --- | --- |
| $\alpha$-CTgen | $\alpha$-chymotrypsinogen |
| $\alpha$-Lac | $\alpha$-Lactalbumin |
| $\beta$-hydrox | $\beta$-hydroxectoine |
| Ala | Alanine |
| Arg | Arginine |
| Asn | Asparagine |
| Asp | Aspartic acid |
| CSP | Cold-Shock protein |
| Cys | Cysteine |
| GdmHCl | Guanidine hydrochloride |
| Glu | Glutamine |
| Gln | Glutamic acid |
| Gly | Glycine |
| His | Histidine |
| HPr | Histidine-containing phosphocarrier protein |
| Ile | Isoleucine |
| KCl | Potassium Chloride |
| Leu | Leucine |
| Lys | Lysine |
| Met | Methionine |
| Phe | Phenylalanine |

| | |
|---|---|
| Pro | Proline |
| Ser | Serine |
| SOD1 | Super-oxide dismutase 1 |
| Thr | Threonine |
| TMAO | Trimethylamine N-oxide |
| Trp | Tryptophan |
| Tryps inh. | Trypsin inhibitor |
| Tyr | Tyrosine |
| Val | Valine |

# List of Symbols

| Symbol | Description |
|---|---|
| $\beta$ | The inverse of $k_B T$ |
| $\epsilon$ | van der Waals well depth; dielectric constant |
| $\lambda$ | Coupling parameter; thermal wavelength |
| $\mu$ | Chemical potential |
| $\rho$ | Charge density; equilibrium density |
| $\sigma$ | Standard deviation; Weighting function in WDA approach to DFT |
| $\Phi$ | Solvent-solvent interaction functional |
| $\phi$ | Density of Solvent Molecules |
| $\chi^2$ | Chi-squared test statistic |
| $\tau$ | Weighting function in WDA approach to DFT |
| $A$ | Surface Area |
| $C_p$ | Heat Capacity (at constant pressure) |
| $F$ | Force |
| $N$ | Number of particles in the system |
| $G$ | Free energy |
| $g(\mathbf{r})$ | Correlation function |
| $H$ | Enthalpy |
| $K_u$ | Equilibrium constant (of unfolding) |
| $k_B$ | Boltzmann's Constant |
| $P$ | Pressure |
| $p$ | probability |
| $S$ | Entropy |

| | |
|---|---|
| **R** | Atom Coordinate |
| $T$ | Temperature |
| $t$ | time; thickness |
| $U$ | Potential Energy |
| $V$ | Volume; potential energy terms |
| $\mathcal{V}$ | External potential |
| $W$ | Work |

# Acknowledgements

I would like to thank my supervisor, Dr. Steven S Plotkin, for his guidance and advice throughout my PhD.

I would like to acknowledge the following members of the Plotkin and Rottler groups for their feedback on this project at various stages: Miguel Garcia, Shazhad Ghanbarian, Mona Habibi, Liam Huber, Ali Mohazzab, Amanda Parker, and Anton Smessaert.

Most importantly, thank you Sara; your support made this possible.

# Chapter 1

# Introduction

Proteins are a class of biological molecules which facilitate essentially every biological process. The primary activity of DNA is coding for proteins; proteins thus act as the intermediary between genetic information and the physical outcome in an organism. Proteins act as hormones, enzymes, cellular structure elements, transportation, and play many other roles in the organism. A full understanding of their structure and function is thus sought both for foundational understanding of biology and for purposes of practical medical treatment of currently intractable diseases. For example, most new drug design efforts target proteins[1]. Understanding proteins is crucial for efforts to combat a variety of diseases, including cancer, Alzheimer's, and ALS[2].

Proteins are long molecules made up of a chain of amino acids. There are twenty different amino acids in eukaryotic proteins, forming an alphabet out of which the sequence of each protein is made. Once assembled into a protein, an individual amino acid is known as a residue, and consists of two parts: a backbone component that is the same for each amino acid and links them all together, and a side chain that differentiates the amino acids and provides much of their functionality. Many proteins "fold", adopting a well-defined three-dimensional structure. This structure then determines the function of the protein. The central dogma of molecular biology is that the folded structure of the protein is uniquely determined by its sequence[3].

Figure 1.1: **Protein structure elements.** Primary (A), secondary (B), tertiary (C) and quaternary (D) structure of a protein.

Understanding the way in which the folded structure is determined from the sequence and the protein's environment is a key challenge facing researchers.

In protein terminology the primary structure refers to the sequence of amino acid residues; see Figure 1.1. The secondary structure is a small set of structures which are seen in many proteins. These include $\alpha$-helices and $\beta$-sheets. These secondary structures and so-called unstructured regions (regions of the primary structure which do not have an identified secondary structure) fold together to form the tertiary structure. Individual protein molecules (monomers) often come together to form complexes, and this is referred to as quaternary structure.

## 1.1   Experimental Methods

Given the limits of *ab initio* structure prediction and the long time scales on which proteins fold, computational studies typically start with experimental structures. Further, since the structure of the protein is the most funda-

mental thing about it from a biochemical perspective, much work has gone into determining it. The most common ways of determining this structure are x-ray crystallography and nuclear magnetic resonance (NMR). X-ray crystallography methods rely on scattering x-rays from proteins which have been condensed into a regular crystal[4]. NMR methods use the relaxation of nuclear spins to determine the chemical environment of each nucleus (typically hydrogen) and thus the structure of the protein. Various structure motifs in a protein give rise to distinct NMR signals, which allows for the structure to be reconstructed with specialized software[5]. Both x-ray crystallography and NMR require a thermodynamically large number of copies of the protein in the same conformation. These methods are thus limited to well-structured proteins. Obtaining structural ensembles of proteins (or regions of proteins) that are intrinsically disordered is an ongoing challenge.

In addition to the folded structure of the protein, experiments can seek out the thermodynamics of folding and unfolding. Typically the protein of interest is placed in a buffer and either the temperature or the solution composition is varied to drive the transition from folded to unfolded. The techniques to measure the resultant behaviour most relevant to this thesis are differential scanning calorimetry (DSC) and circular dichroism.

Circular dichroism refers to the differential absorption of right and left circularly polarized light. At certain wavelengths in the UV range secondary structures of a protein have a characteristic difference in the amount of right and left circularly polarized light they absorb. Thus the circular dichroism at that wavelength of a sample indicates the percentage of proteins that have those secondary structures intact–$i.e$ the fraction of folded proteins[6].

This analysis typically assumes a two-state folding process, so that if a single secondary structure element is missing, the entire protein is unfolded. As we will discuss in chapter 2, the two-state assumption is quite good for a broad class of proteins.

Calorimetry refers to any of the variety of ways of determining the heat capacity of a sample. Differential scanning calorimetry involves heating two samples, one of which contains the protein of interest and the other of which is a reference. The difference in heat required to keep both samples

to the same temperature *vs* time is then recorded and the heat capacity extracted from this. Of particular interest is the heat capacity as a function of temperature during the process of unfolding the protein. This information can be used to determine the stability of a protein (that is, the minimum work required to unfold it at a certain temperature) as well as the extent to which that stability is determined by enthalpic and entropic effects[7].

There are obviously many other ways proteins are studied experimentally. Experiments that pull on the protein with AFM or optical tweezers are used to gain information about the stability of the protein structure and the unfolding and folding pathways. Binding experiments can also be performed, in which an agent which binds to the protein in one state (*eg* the unfolded state) but not another. This agent can then be measured and serve as a proxy for the amount of unfolded protein.

## 1.2   Entropy Enthalpy Compensation

Enthalpy and entropy play an intimately connected role in the free energy change during numerous biochemical processes. It has long been known, for example, that the transfer of hydrocarbons such as alkanes or alcohols from pure solvent to water is generally exothermic or enthalpically favorable ($\Delta H \approx -(3\text{-}6)\text{kcal/mol}$ for methane-butane) but entropically unfavorable, with $T\Delta S$ typically 2-3$\times$ the magnitude of the enthalpic contribution [8, 9]. These opposing thermodynamic forces result in a net free energy change that is smaller than either of the enthalpic or entropic contributions.

The first part of this thesis concerns an effect known as entropy-enthalpy compensation. For a variety of physical processes including solute transfer [10], unfolding of various proteins [11], and ligand binding, ionization, and hydrolysis [12], the changes in enthalpy and entropy obey a nearly linear relationship when a variable such as binding ligand is varied; this is referred to as entropy-enthalpy compensation [12–17]. The effect is ubiquitous but not universal [18]. The slope of the enthalpy *vs.* entropy plot, referred to as the compensation temperature $T_c$, ranges from about 150K (e.g. for alkane vaporization [19]) to about 300K (most processes). The difficulty in

4

designing high affinity drugs has been attributed to entropy-enthalpy compensation [20–22].

As mentioned above, entropy-enthalpy compensation is not an inevitable consequence of statistical mechanics, particularly along chemical reaction coordinates. Long-lived metastable states due to large barriers, and thus the absence of any significant entropy-enthalpy compensation along the reaction coordinate, are fairly common in condensed matter and biophysics. The diamond phase of carbon is metastable to graphite at standard temperature and pressure, with an enormous conversion barrier; allotropes of boron, polymorphs of silica, and martensite in steel are all metastable phases with prohibitive transition barrier; colloidal systems and emulsions have long-lived metastable phases; long-lived structure with slow, glassy dynamics is common in supercooled liquids; the covalent bonds forming the backbones of DNA, RNA, and proteins are metastable to hydrolysis; several proteins have native, functional states that are metastable, but simply have enormous kinetic unfolding barriers, including alpha-lytic protease, subtilisin, Streptomyces griseus protease B and the aspartic peptidase pepsin [23]. In multimeric systems of chain length $\lesssim 100$ amino acids, native protein structures have been observed to have higher free energy than the amyloid phase, implying that a significant portion of the proteome is conformationally in metastable equilibrium[24]. In contrast to these observations, as a general rule, entropy-enthalpy cancellation does play a critical role in governing the foldability of proteins and resolving the Levinthal paradox [25]. Small barriers in protein folding have been shown to arise due to the locality of interactions and the concomitant loss of entropy in forming stabilizing interactions. If it were not for entropy-enthalpy cancellation as protein chain conformations progressed towards native-like folds, folding barriers would be prohibitively high, proteins could not fold on biological time scales, and life as we know it would not be possible.

The entropy is often obtained from measurements of the enthalpy and free energy by subtraction; one complication however is that errors in enthalpy are often much larger than errors in free energy. In these cases the error in enthalpy can induce a spurious linear relation to the entropy,

5

for example in early measurements for oximation reactions of alkyl thymyl ketones [26], correlated entropy-enthalpy errors were sufficiently large that they could account for the whole measurable effect of entropy-enthalpy compensation, which could then not be definitively proven, regardless of whether or not it existed. Correlated errors generally have an effective compensation temperature equal to the temperature at which the measurements were taken. Compensation exists however when large entropy and enthalpy changes either cancel or compensate each other to yield a relatively small net free energy gain for a given process, regardless of the compensation temperature, and includes cases where the compensation temperature equals the lab temperature [13, 16, 27]. As well, even if the compensation temperature is quite different from the temperature of the experiments, if the correlated scatter is sufficiently large, it can rule out the significance of the effect. Thus there is a need to introduce more rigorous error analysis to judge the significance of any observed entropy-enthalpy compensation. We develop such an analysis in Chapter 2.

## 1.3   Implicit Solvent Models

Proteins fold and function in the crowded environment of the cell. Cytosolic proteins must negotiate a complex milieu which in many ways is significantly different than the environment in the test tube: roughly 15% of water molecules are motionally restricted by protein and membrane surfaces [28]; the surrounding solvent is enriched in ions such as Potassium but depleted in Sodium and Chlorine; osmoprotectants such as trehalose and various amino acids are present in significant concentration; numerous membrane surfaces such as the nucleus, ER, and Golgi impose charged substrates for protein interaction; macromolecular agents such as the microtubules, actin, ribosomes, soluble proteins and RNA occupy roughly 30% ($\approx 300\text{g}/\ell$) of the cellular volume, and modulate stability [29], aggregation propensity [30], and dissociation constants [31, 32].

   Non-cytosolic proteins also fold in environments distinct from the test tube as well as the cytosol, particularly with respect to ionic and redox condi-

tions as well as the chaperone complement. Proteins destined for the plasma membrane or extracellular matrix are trafficked by the secretory pathway through the ER and Golgi [33]. The environments in the ER and cytosol are sufficiently different that the conditions for protein folding are generally mutually exclusive between the two milieux. Folding generally occurs in the lumen of the ER, while function occurs either on the plasma membrane or in the extracellular matrix, which is itself densely occupied by highly charged glycosaminoglycans such as hyaluronan and heparin sulfate—large molecules that may facilitate cellular migration and regulate secreted protein activity. Fibrous proteins such as collagen and fibronectin also occupy the extracellular space, and provide structural rigidity while allowing rapid diffusion of nutrients and signalling metabolites between constituent cells.

The above examples demonstrate the need to correctly account for the effect of the cell environment on protein folding, stability, and function. Accurately accounting for the effects of the cell environment presents a challenge however to both experimental and computational studies. Experimentally, most of what is known about protein folding and stability has resulted from *in vitro* studies at dilute concentrations, and many questions remain as to how well such results apply to a realistic cell environment. Computationally, including explicit solvent along with a realistic concentration of osmolytes in a box of sufficient size to implement periodic boundary conditions outside the range of an electrostatic cutoff typically increases the number of particles in the simulation by a factor on the order of ten or more [34]. While this can be done for small proteins such as Trp-cage [34], investigating larger proteins generally requires coarse-grained models to keep the computational resources required reasonable [35].

Computational studies of crowding on isolated monomeric minimal $\beta$-barrel proteins find that the folding temperature is increased and the folding time decreased [36, 37]. However, molecular crowding has been shown in secretory cells to impair protein folding and lead to aggregate formation in the ER [38]. It has been estimated that increasing the total intracellular protein concentration by 10% can potentially increase the rate of protein misfolding reactions following a nucleation-polymerization mechanism by a factor

of 10 [39]. Consistent with these observations and estimates, another MD folding study of a coarse-grained model of crambin found that the presence of multiple protein copies with a weak inter-protein attractive potential (a more realistic scenario) hindered correct monomeric folding and predisposed the system to aggregation and misfolding [40].

The above considerations motivate the creation of computational models, with which we can account for the cellular environment around a protein in an accurate but less computationally expensive way. Further, while one might naïvely suspect that an explicit approach would yield results that are, in principle, exact, a number of difficulties arise in implementing these models. The "best" form of the interaction between water and other molecules depends on the context of the simulation, and a number of different water models are available with various strengths and drawbacks[41, 42]. Implicit solvent models attempt to address these issues.

Implicit solvent models are attempts to compute the transfer free energy of each configuration of the molecule of interest from a vacuum to a solvent (or, more generally, from one solvent environment to another) through some function of the coordinates of the solute biomolecule.

$$\Delta G_{\text{transfer}} = g(\{\mathbf{R}_i\}) \tag{1.1}$$

where $\{\mathbf{R}_i\}$ is the set of coordinates that defines the position and configuration of the biomolecule. Knowing the transfer free energy (or its derivatives with respect to atom position) for each configuration allows the simulation to update the position of the protein atoms as if the solvent was present. Expressing the transfer free energy in the form of equation 1.1 eliminates the need to explicitly simulate the solvent and hence speeds up the simulation— provided the implicit solvent algorithm is fast enough. A number of models to compute this transfer free energy have been proposed, with the goal of calculating $\Delta G$ with less computational resources than an explicit solvent, while maintaining enough accuracy for the purposes of the simulation. Many such models have been proposed, with various strengths and weaknesses[43].

The most common implicit solvent models split the transfer energy into

two parts: a dielectric contribution and a non-polar contribution. These components are each modelled separately and their contributions summed[44].

$$\Delta G = \Delta G_{\text{nonpolar}} + \Delta G_{\text{dielectric}} \tag{1.2}$$

The dielectric part is typically the most time-consuming part to calculate because of the long-ranged forces involved. The solvent is often assumed to be a continuous dielectric, and the transfer free energy of the biomolecule from vacuum to solvent can then be computed via solving the Poisson-Boltzmann (PB) equation with finite element or finite difference methods[45, 46], which is accurate up to the continuum assumption but relatively costly[47], or by faster but more approximate methods such as the Generalized Born (GB) method[48]. Generalized Born in particular has become popular, and has been implemented and optimized for a variety of forcefields used in molecular dynamics[49, 50].

In the GB model the starting point is the standard Poisson equation for a dielectric:

$$\nabla\left[\epsilon(\mathbf{r})\nabla V(\mathbf{r})\right] = -4\pi\rho(\mathbf{r}) \tag{1.3}$$

where $\epsilon(\mathbf{r})$ is the position-dependent dielectric, $V$ the electric potential, and $\rho(\mathbf{r})$ the density of charge. This can be recast as an equivalent equation for the Green's function:

$$\nabla\left[\epsilon(\mathbf{r})\nabla\mathbf{G}(\mathbf{r}_i, \mathbf{r}_j)\right] = -4\pi\delta(\mathbf{r}_i - \mathbf{r}_j) \tag{1.4}$$

Equation 1.4 has an analytical solution for a model system which consists of a collection of charges at positions $\mathbf{r}_i$ inside a sphere of dielectric $\epsilon_{in}$ of radius $R_{\text{sphere}}$ surrounded by an infinite medium with dielectric $\epsilon_{out}$, with the condition that $\epsilon_{out} \gg \epsilon_{in} \geq 1$:

$$\mathbf{G}(\mathbf{r}_i, \mathbf{r}_j) = \frac{1}{\epsilon_{in}|\mathbf{r}_i - \mathbf{r}_j|} + \mathbf{F}(\mathbf{r}_i, \mathbf{r}_j) \tag{1.5}$$

$$\Delta G_{el} = \frac{1}{2}\sum_{ij}\mathbf{F}(\mathbf{r}_i, \mathbf{r}_j)q_i q_j \tag{1.6}$$

where $\Delta G_{el}$ is the free energy to insert the set of charges $q_i$ into the dielectric sphere at positions $\mathbf{r}_i$ and $\mathbf{F}(\mathbf{r}_i, \mathbf{r}_j)$ is given by

$$\mathbf{F}(\mathbf{r}_i, \mathbf{r}_j) = \left( \frac{1}{\epsilon_{in}} - \frac{1}{\epsilon_{out}} \right) \frac{1}{\sqrt{r_{ij}^2 + \left( R_{\text{sphere}} - \frac{r_i^2}{R_{\text{sphere}}} \right) \left( R_{\text{sphere}} - \frac{r_j^2}{R_{\text{sphere}}} \right)}}$$

(1.7)

The term $R_{\text{sphere}} - \frac{r}{R_{\text{sphere}}}$ is then defined to be the generalized born radius $R$. That this radius is constant is an approximation, as it should in principle include $r$. The free energy to insert a single charged atom is then

$$\Delta G_{el} = -\frac{1}{2} \left( \frac{1}{\epsilon_{in}} - \frac{1}{\epsilon_{out}} \right) \frac{q_i^2}{R_i}$$

(1.8)

(the limiting case of $\mathbf{F}$ when $\mathbf{r}_i = \mathbf{r}_j$) in this case, and this expression is used to calculate the Born radii for each atom.

Equation 1.7 can be modified in various ways, which has given rise to several different methods of calculating Born radii[48, 51, 52]. These methods all share the same essence of calculating a Born radius from a list of known self-interaction values and then using that radius to calculate the interaction between charges.

The nonpolar contribution consists of the free energy to create a cavity in the dielectric and the free energy arising from the van der Waals interaction between the protein and water. The observed linear dependence of the log solubility on the number of $CH_2$ groups and hence chain length, particularly for long chain saturated fatty acids (decanoic acid and longer), and long-chain aliphatic alcohols (1-butanol and longer), can be taken to indicate a free energy change upon transfer to solvent that scales linearly with either volume or surface area. Historically, surface area has been chosen, under the assumption that interactions with the solvent take place at the surface of the molecule in question [53, 54]. Then the free energy difference between an amino acid in water and in a solvent with some osmolyte concentration is, for a given configuration, given in terms of the solvent accessible surface area (SASA) of that configuration by the phenomenological expression: $\Delta G =$

$\gamma \cdot$SASA$+c$, where $\gamma$ is obtained from, $eg$, a tri-peptide experiment [55]. This approximation is severe, and neglects aspects such as the length of the van der Waals interaction[56] and the temperature dependence of the free energy, despite their known importance. The desolvation barrier is another effect of importance which SASA approaches cannot capture[57]. The surface area approach can be expanded such that

$$\Delta G_{\mathrm{nonpolar}} = \sum_i \gamma_i A_i + b \tag{1.9}$$

where $\gamma_i$ and $A_i$ are different regions of the protein. Distinct coefficients can even be assigned to each residue type. These modifications do not address either the temperature dependence or the lack of a desolvation barrier.

In addition to the Generalized Born/Surface Area (GB/SA) approach, there are a number of other approaches which use the general scheme of expressing the transfer free energy as a sum of non-polar and dielectric conbributions. Some attempt more efficient calculation of the dielectric component[58], while others tweak the nonpolar contribution in ways such as adding a term proportional to the volume of the biomolecule[59]. Recent simulation studies have found significant volume contributions to transfer free energies, however [35]. In these studies, model solvents with no enthalpic interaction (hard sphere solvents) still showed significant transfer free energies, due solely to excluded volume. Volume corrections to the surface area model, computed by scaled particle theory or RISM approaches, have been investigated by several authors [60–63]. As well, Baker and colleagues have found that the inclusion of volume terms (computed by scaled particle theory) and dispersion integral terms (computed by Weeks-Chandler-Andersen theory) were essential for an accurate implicit solvent description of atomic-scale nonpolar forces [64].

While the combination of generalized Born electrostatics and surface area nonpolar contributions (GB/SA) is perhaps the most popular choice of implicit solvent, others have been proposed as well. Kovalenko and colleagues[65] have developed an approach based on solving the Ornstein-Zernike (O-Z) equation within the reference site interaction model (RISM)

11

proposed by Chandler[66]. The O-Z equation describes density correlations in a fluid, through the equation

$$h(\mathbf{r}_{12}) = c(\mathbf{r}_{12}) + \rho \int d^3 r_3 \, c(\mathbf{r}_{13}) h(\mathbf{r}_{23}) \qquad (1.10)$$

where $h(\mathbf{r}) = g(\mathbf{r}) - 1$, the so-called total correlation function, and $c(\mathbf{r})$ is the direct correlation function. The O-Z equation defines an expansion in correlation functions between the model's reference sites, which requires a closure relation to then solve[67]. In the RISM-3D approach this closure relation is taken to be

$$g(\mathbf{r}) = \qquad \exp(\chi(\mathbf{r})) \qquad \text{for } \chi(\mathbf{r}) \leq 0 \qquad (1.11)$$

$$g(\mathbf{r}) = \qquad 1 + \chi(\mathbf{r}) \qquad \text{for } \chi(\mathbf{r}) > 0 \qquad (1.12)$$

$$\chi(\mathbf{r}) = \quad -\beta u(\mathbf{r}) + h(\mathbf{r}) - c(\mathbf{r}) \qquad (1.13)$$

where $u(\mathbf{r})$ is the interaction potential between sites.

In principle this approach can be expanded to arbitrary accuracy, but in practice only a few paths can be taken from any given reference site before the computational cost becomes prohibitive. The speed is further improved by only updating the solvent density every few simulation steps, and by using the previous value at each point as an initial guess. This method requires a classical interaction potential between the protein and the solvent and between the solvent molecules themselves, and is limited by the accuracy of this potential.

Another implicit solvent approach that differs fundamentally from GB/SA is the ABSINTH model proposed by Vitalis and Pappu[68]. Here the transfer free energy is not broken up into dielectric and nonpolar components, but rather into a direct mean field interaction term and a screening term. This not only accounts for the solvent-protein interactions, but also the screening of the protein-protein interactions. The protein is broken up into solvation

groups, the solvation free energy of which is expressed as

$$\Delta G_{\text{solv}}^{\text{total}} = \sum_{i=1}^{N_{SG}} \left[ \sum_{k=1}^{n_i} \lambda_k^i \cdot v_k^i \right] \cdot \Delta G_{\text{solv}}^i \qquad (1.14)$$

Here $\Delta G_{\text{solv}}^i$ is the experimental solvation free energy of a particular group of atoms (based on model compounds). The $v_k^i$ are solvation states of the group of atoms measured by the available volume surrounding the atoms, rather than the surface area, and $\lambda_k^i$ are a set of weights to improve the accuracy of the model. The available volume here is measured by considering a sphere of some radius around the atom, and subtracting from that sphere the volume of all other protein atoms which overlap. The remaining volume is thus available to solvent.

Some of the molecular dynamics approaches discussed above implicitly use a solvent without an implicit solvent term appearing in the potential function. As discussed above, Gō models and associative memory hamiltonians, for example, use the known native structures of a protein to define a set of interactions between atoms or residues. Since these structures are observed in the presence of solvent, solvent effects are captured to some extent in these model interactions.

## 1.4  Classical Density Functional Theory

Density functional theory (DFT) is best known in the context of quantum DFT[69], where it is widely used in computational studies of electronic and atomic structure, vibrational spectra, magnetic resonance, and reaction dynamics[70, 71]. The method has also been applied to classical density fields to compute correlation functions and dynamics in liquids[72–77]. In fact the essence of classical density functional theory (cDFT) actually predates the formal development of DFT[78]. Since then, cDFT has been widely used to model liquids around nanostructures[79, 80] and has also been applied to biomolecules[81, 82]. Density functional theories of inter-residue contact probability developed by Plotkin and Onuchic have elucidated the

13

effects of energetic and entropic heterogeneity on protein folding free energy barriers [83, 84], while Wolynes has used DFT in fundamental studies of glass physics and the glass transition [85–91], and also in protein folding[92–94].

Density functional theory relies on the following insight: if the free energy of a system is expressed as a functional of the single particle density, the single particle density that minimizes the functional will be the true single particle density, and the resulting value of the functional will be the true free energy. This is not obvious; the full expression for the free energy includes probabilities over $N$ positions and momenta, where $N$ is the number of particles in the system. Nonetheless a functional of the single particle density can be used to determine the real free energy–a proof of this is shown in Appendix A. The form of the density functional is taken to be

$$
G = \int d^3r\, k_B T \left[\phi(\mathbf{r}) \ln(\delta V \phi(\mathbf{r})) - \phi(\mathbf{r})\right] + \mathcal{V}(\mathbf{r})\phi(\mathbf{r})
$$
$$
+ \Phi[\phi] \tag{1.15}
$$

where $G$ is the free energy of the system, $\phi$ the position dependent density of the solvent or cosolute, $\delta V$ a volume element, $\mathcal{V}$ the external potential, and $\Phi[\phi]$ the free energy functional arising from inter-particle interaction terms. The equilibrium free energy and density is found from equation 1.15 through functional differentiation:

$$
\frac{\delta}{\delta \phi} \left[G - \mu \int d^3r\, \phi(\mathbf{r})\right] = 0 \tag{1.16}
$$

where $\mu$, the Lagrange multiplier associated with the condition that particle number is conserved, is the chemical potential. In the classical DFT literature $G - \Phi$ is known as the ideal free energy, which is the same in each of the cDFT approaches we discuss below, while $\Phi$ is known as the excess free energy. In Appendix A we review the proof that $G$ is a unique function of $\phi$, the single particle density, and hence the function $\phi$ that minimizes $G$ is the equilibrium density.

In traditional approaches to classical DFT[95] $\Phi$ is broken up into two parts: a hard-sphere term and a term arising from any attractive particle-

particle interactions–$\Phi = \Phi_{\text{att}} + \Phi_{\text{hs}}$. The attractive term is taken to be:

$$\Phi_{\text{att}} = \int d^3r \, d^3r' \, U_{\text{att}}(\mathbf{r} - \mathbf{r}')\phi(\mathbf{r})\phi(\mathbf{r}') \tag{1.17}$$

where $U_{\text{att}}(\mathbf{r} - \mathbf{r}')$ is the attractive part of the solvent-solvent potential. This form assumes a uniform fluid; that is, the solvent-solvent correlation function $g(\mathbf{r}, \mathbf{r}') = 1$.

The hard-sphere term in the solvent-solvent interaction functional arises purely from entropic terms and does not have an analytical form. One common approach is to take a weighted density average (WDA), so that the hard-sphere term is written as[95]

$$\Phi_{\text{hs}} = \int d^3r \, \bar{\phi}^\sigma(\mathbf{r})\mathcal{F}(\bar{\phi}^\tau(\mathbf{r})) \tag{1.18}$$

where $\bar{\phi}^\sigma$ and $\bar{\phi}^\tau$ are weighted averages of $\phi$ around position $\mathbf{r}$–e.g. $\bar{\phi}^\sigma = \int d^3r' \, \sigma(\mathbf{r} - \mathbf{r}')\phi(\mathbf{r}')$, and $\bar{\phi}^\tau = \int d^3r' \, \tau(\mathbf{r} - \mathbf{r}')\phi(\mathbf{r}')$. $\mathcal{F}$ is an arbitrary function and $\sigma(\mathbf{r} - \mathbf{r}')$ and $\sigma(\mathbf{r} - \mathbf{r}')$ are local weighting functions, which can have any form but need to go to zero for large $\mathbf{r} - \mathbf{r}'$ in order to be useful. $\mathcal{F}$ can be assigned by requiring that the equation of state for a uniform system in the absence of external potential (such that $\phi(\mathbf{r}) = \rho$, the bulk density) reproduce some known equation of state. For example, to recover the Carnahan-Starling equation of state,

$$P = Nk_BT/V + k_BT\frac{1 + y + y^- y^3}{(1 - y)^3} \tag{1.19}$$

where $y = \pi\rho d^3/6$ and $d$ the hard sphere diameter, $\mathcal{F}$ must be

$$\mathcal{F}(\rho) = k_BT\frac{y(4 - 3y)}{(1 - y)^2} \tag{1.20}$$

so that the pressure from the density functional with a uniform density,

$$P = \frac{\partial G[\rho]}{\partial V} \tag{1.21}$$

15

matches the Carnahan-Starling pressure.

Similarly the forms of the weighting functions $\sigma(\mathbf{r} - \mathbf{r}')$ and $\tau(\mathbf{r} - \mathbf{r}')$ are determined by comparison with other systems. For example, in one dimension the weighting functions

$$\sigma(r) = \frac{1}{2}\delta\left(\frac{d}{2} - |r|\right) \tag{1.22}$$

$$\tau(r) = \frac{1}{d}\Theta\left(\frac{d}{2} - |r|\right) \tag{1.23}$$

where $\Theta$ is the Heaviside function, reproduces the analytic result for one dimensional hard rods of length $d$. This then can be generalized to a three dimensional system, in particular since many systems of interest to early cDFT studies are confined and thus behave as lower dimensional systems on certain length scales[80].

The weighted density average approach can be extended to include gradients of the density–terms such as $1/2 \int d^3r\, k(\nabla\phi)^2$, with $k$ a constant. Oxtaby and colleagues, for example, use such terms to apply classical DFT to crystal growth and nucleation problems[96].

Despite the fact that the correlation function $g(\mathbf{r}, \mathbf{r}')$ was assumed to be uniform in writing $\Phi_{\text{att}}$, WDA models do not assume a uniform density. Rather the correlation function is implicit in the form of the weighted averages. Actually obtaining this correlation function, however, typically involves first solving for the density and thus is non-trivial to obtain from the density functional. In this approach the total correlation function is a prediction of the theory rather than an input.

In general these approaches favour accuracy over speed. By appropriately modelling the solvent-solvent interaction terms and iterating a number of times to find the minimum free energy, cDFT can capture realistic correlation functions and energies, but at a computational cost high enough to be unfeasible for implicit solvent applications in molecular dynamics in large systems.

Kinjo and Takada[81] have applied cDFT to protein systems. Their approach treats both the solvent and the protein with a density field. This

16

approach allows them to study crowding effects in a general way, but treats the protein in an extremely simplified manner in order to reduce it to a manageable density field.

Borgis and colleagues have also applied cDFT to proteins[82], and solvation in particular. Their approach differs somewhat from more traditional cDFT schemes. Rather than break the solvent-solvent interaction potential up into attractive and hard sphere components, Borgis writes the interaction functional $\Phi$ as[97, 98]

$$\Phi[\phi] = \int \int d^3r d^3r' \phi(\mathbf{r})\phi(\mathbf{r}')U(\mathbf{r} - \mathbf{r}')g(\mathbf{r}, \mathbf{r}') \qquad (1.24)$$

where $U$ is the full solvent-solvent interaction potential and $g$ is the correlation function, which is taken from experimental structure factors. Thus here, rather than the solvent-solvent correlation function being a prediction of the theory, it is an input. This approach makes two implicit assumptions: first, that the correlation function is equal to the direct correlation function over the interaction potential, and second, that the two-body and higher *entropic* terms are zero. To see the first assumption, we note that the direct correlation function can be found from the free energy functional by

$$\frac{\delta^2 G}{k_B T \delta\phi(\mathbf{r})\delta\phi(\mathbf{r}')} = c(\mathbf{r}, \mathbf{r}') \qquad (1.25)$$

which, for the interaction functional 1.24, gives $g(\mathbf{r}, \mathbf{r}') \cdot \beta U(\mathbf{r}, \mathbf{r}') = c(\mathbf{r}, \mathbf{r}')$.

To see the second assumption we note that the interaction term is entirely enthalpic; $\partial G/\partial T$ generates only the ideal gas term in this model. Put another way, in the hard sphere limit, in which $U(r) = 0$ for any region in which the correlation function is non-zero, the Borgis interaction functional is identically zero and the model reduces to an ideal gas.

One feature both the Borgis approach and the traditional cDFT approaches share is the need to iterate over the equations resulting from applying equation 1.16 to the respective models to obtain a density which converges to the final self-consistent solution. This iterative nature of the solution creates a cost penalty to applying these methods. In this thesis

we seek a model which avoids iterative solutions in order to develop an algorithm that can be applied many millions of times during a molecular dynamics simulation. Thus, while other approaches start from a relatively accurate functional and seek to simplify it, we will start from the most basic functional and seek to add as little as possible to make it sufficiently accurate.

## 1.5   Molecular Dynamics

One of the principal ways of studying proteins computationally is molecular dynamics (MD)[99]. As in much of condensed matter, the underlying physics of the protein molecule is in principle known: a collection of nuclei and electrons which interact via the electromagnetic force. But performing explicit quantum mechanical calculations on many thousands of atoms at 300K in many configurations is well beyond even the most powerful of computers, and will likely remain so for some time. Thus a series of approximations are made to reduce the problem to something more manageable. The most basic approximation is that the atoms of the protein are treated as classical particles. These particles interact with each other through a series of potentials. In the so-called All Atom approach, each atom of the protein is treated as a point particle, and the atoms interact through potentials such as bond potential, angle potentials, dihedral potential, Coulomb potential, and van der Waals potential; *i.e.* the total potential energy $U$ of the system is given by

$$U = \sum_{ij} V_{ij}^{\text{bond}} + V_{ij}^{\text{Coulomb}} + V_{ij}^{\text{vdW}} \tag{1.26}$$

$$+ \sum_{ijk} V_{ijk}^{\text{angle}} \tag{1.27}$$

$$+ \sum_{ijkl} V_{ijkl}^{\text{dihedral}} \tag{1.28}$$

where each term will be discussed below. Once these terms have been determined, the basic approach of a molecular dynamics program is straight-

forward: at a given time, for each particle in the system, compute the sum of the forces on that particle,

$$\mathbf{F}_i = -\nabla_i U \tag{1.29}$$

then evolve the system using Newton's second law $\mathbf{F} = m\mathbf{a}$ and a discrete timestep $\Delta t$. This is repeated until a simulation of sufficient length to sample the properties of interest has been obtained. The difficulty in practice is the "sufficient length"; a typical value of $\Delta t$ for an all-atom simulation is 2 fs, while the folding time of a protein might be on the order of 10 s[100]. Thus even in this classical approximation many phenomena are out of reach. The timestep $\Delta t$ is determined by requiring a certain level of accuracy in the simulation. Discretizing the integration of Newton's second law introduces error in the system. Conceptually this arises from the possibility that a particle will move through a region in which the potential changes so rapidly that the integration step cannot keep up. This results in a certain amount of "shadow work" done by the integrator on the system[101]. The timestep is thus a compromise between choosing a large value to speed up the simulation while keeping it small enough to maintain stability. In practice, 2 fs is often used as a standard value, with the condition that the bonded interactions involving hydrogens are rigidly constrained (or a so-called "united atom" forcefield is used to eliminate the hydrogen atoms).

The bond potential $V_{\text{bond}}$ is often assumed to be harmonic;

$$V_{\text{bond}}(\mathbf{r}_1 - \mathbf{r}_2) = \frac{1}{2}k(|\mathbf{r}_1 - \mathbf{r}_2| - d_0)^2 \tag{1.30}$$

There are other forms, though; the Morse potential is one such alternative[102]. The bond can also be rigidly constrained. The angle potentials also take a harmonic form,

$$V_{\text{angle}}(\theta_{123}) = \frac{1}{2}k(\theta_{123} - \theta_0)^2 \tag{1.31}$$

In this case there is, in principle, a concern that the form of the potential is not periodic in $\theta$; in practice this concern is dealt with by ensuring that $k$ has a value such that $\theta_{123} - \theta_0$ is never so large as to be an issue. The

dihedral potential (see figure 1.2) cannot in general be harmonic, because it is much softer than the angle potential and there are many situations in which dihedral angles vary through the full $2\pi$. Further there may be several local minima. An example of this is the cis and trans configurations of small molecules. Common forms for the dihedral angle thus include

$$V_{\text{dihedral}}(\theta_{1234}) = k(1 - \cos(n\theta_{1234} - \theta_0)) \tag{1.32}$$

$$V_{\text{dihedral}}(\theta_{1234}) = \sum_{n=1}^{5} C_n (cos\theta)^n \tag{1.33}$$

There are two types of dihedral angles considered. Proper dihedrals occur when four atoms are bonded in a line, as in figure 1.2 panel A. The dihedral angle is then defined as the angle between the plane formed by atoms i, j, and k and the plane formed by atoms j, k, and l. Improper dihedrals pertain to three atoms bonded to a central atom, as in figure 1.2 panel B. Again, the dihedral angle is defined as the angle between the plane formed by atoms i, j, and k, and the plane formed by atoms i, k, and l, but here the interpretation of this angle is somewhat different. While proper dihedrals are used to specify cis and trans configurations, improper dihedrals are used to stiffen groups of planar atoms.

The form of the force-field potential is chosen to account for the fact that interactions between atoms can be strongly non-two-body. Thus the angle and dihedral terms attempt to capture the physics of bonded interactions, which tend towards particular geometric arrangements.

The van der Waals interaction and the Coulomb interaction are termed "non-bonded" interactions. They are present between all pairs of atoms that are in different molecules or in the same molecule but separated by at least a certain number of bonds (typically three). In most cases they have the familiar forms:

$$V_{\text{vdW}}(r) = 4\epsilon \left( \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^6 \right) \tag{1.34}$$

$$V_{\text{Coulomb}}(r) = k_q \frac{q_1 q_2}{r} \tag{1.35}$$

Figure 1.2: **Diagram of dihedral angles.** Panel A shows an example of a proper dihedral angle, while panel B illustrates an improper dihedral.

The van der Waals interaction can, however, be modelled by a 10-12 potential (which may be more appropriate for hydrogen bonding) rather than a 6-12 potential, and some have argued the Buckingham potential is more physically motivated as it attempts to find a form for the repulsive interaction based on observed virial coefficients[103].

This thesis will discuss the derivation of a new implicit solvent model. To put this derivation into context, it is useful to summarize how the parameters in molecular dynamics forcefields are obtained in the first place.

The parameters in equations 1.30-1.35 need to be determined, and there are a number of ways of doing so. Each of them involves attempting to reproduce some aspect of the "real" system. The essential problem, as mentioned above, is that molecules are quantum mechanical entities and we would like to simulate them as classical particles. The CHARMM forcefield[104], for example, uses quantum chemical calculations to generate parameters. The amino acids are broken down into small molecules amenable to a Hartree-Fock computation. These molecules are then placed in the vicinity of a water molecule. The water molecule is placed at a number of different positions to obtain the equilibrium distance to various atoms in the model compound, and the interaction energy at that equilibrium distance; this process is illustrated in figure 1.3. The van der Waals and Coulomb parameters are then set to reproduce these interaction energies and equilibrium distances. Similarly, the bonded parameters are found by varying the distances, angles, and dihedrals of each set of atoms, and finding minimum values and the second

Figure 1.3: **Obtaining the parameters in the CHARMM methodology.** Panel A illustrates a model compound (in this case imidazole) and a nearby water. Panel B shows the interaction energy, calculated in Gaussian[105], for various water-nitrogen distances. The van der Waals and Coulomb parameters are then picked to reproduce the minimum interaction energy and the equilibrium distance.

derivatives of the energy with respect to separation. The bond parameters are set to reproduce the harmonic well each set of atoms rests in.

The OPLS forcefield, on the other hand, compares thermodynamic properties of Monte Carlo simulations of pure liquids of small molecules with experimentally measured bulk properties such as density and heat of vaporization[106]. In each of these forcefields the results for small model compounds are then combined into a force-field for the entire protein.

AMBER is similar to OPLS in that it uses experimental results to determine the force-field parameters, but fits a somewhat different set of properties, such as observed normal modes and vibrational spectra[107].

The all-atom approach is the one most relevant to this thesis, but it is only one of a number of approaches to molecular dynamics. There are a number of coarse-graining methods which can be employed. The simplest are United Atom (UA) forcefields[108]. In these the only coarse-graining is to remove most of the hydrogen atoms, incorporating their effects into modified parameters of the heavy atoms the hydrogens were bonded to–*e.g.*

a methyl group would become a single atom labeled CH$_3$ with an atomic mass of 15 and a radius larger than that of a bare carbon atom.

Coarse graining of atoms can be taken further, subsuming more and more atoms into single particles. Entire side-chains can be approximated as a single particle, or, in the most coarse-grained approached, each residue can be modelled as a single bead[109]. In each of these cases the trade off is for increased simulation speed at the cost of some accuracy.

Coarse-grained approaches require force-fields to be parameterized, and as with all-atom, this can be done in a variety of ways. Typically the parameters are not calculated *ab initio*, but are tuned to reproduce some aspect of the proteins to be studied. The widely used Gō model generates a force-field by starting with the experimentally determined folded structure of the protein, then setting attractive interactions between all pairs of particles in contact in the folded structure (with a equilibrium distance fixed to the observed contact distance) and repulsive interactions between all pairs of particles not in contact[110]. This type of model can be useful for many studies, but obviously cannot be used for things such as structure prediction or investigating intrinsically disordered proteins (proteins that do not adopt a well-defined structure).

More versatile approaches include Associated Memory Hamiltonians. Here, rather than just one folded structure, a database of proteins is used. Each segment of the protein is matched to segments of other proteins that have the similar sequences[111]. The potential energy function for a given protein is then constructed based on an appropriately weighted average of the structures from these similar sequences. In this approach the structure of the protein of interest does not need to be known and intrinsically disordered proteins can be examined. It shares with the Gō model an essentially phenomenological character though.

## 1.6 Computational Methods

As mentioned above, simulations, while providing a level of detail not possible in experiments, are limited in time scale. A protein and accompanying

water molecules, simulated on a high performance cluster of CPUs, might run a nanosecond per hour of wallclock time[112]. Thus simulations of more than a few hundred nanoseconds become prohibitive. The D.E. Shaw company, using a dedicated hardware system and custom software, can simulate on the order of a millisecond[113]. Protein dynamics, however, can occur on the time scales of seconds or longer. Thus while experimental methods determine free energies by looking at populations, computational methods require other approaches.

One of the most relevant questions we can ask of a system is, What is the difference in free energy between two states? Because we cannot easily simulate the systems of interest to biophysics for long enough times to adequately sample each state of interest, we require non-equilibrium techniques to find these differences in free energy. Of particular application to this thesis are thermodynamic integration[114, 115], Bennet acceptance ratio (BAR)[116], the Jarzynski equality[117], and weighted histogram analysis method (WHAM).

Thermodynamic Integration (TI)[114, 115] makes use of the following relation from statistical mechanics: the difference in free energy $\Delta G$ between Hamiltonians $\mathcal{H}_A$ and $\mathcal{H}_B$ is given by

$$\Delta G = \int_0^1 d\lambda \left\langle \frac{\mathrm{d}U}{\mathrm{d}\lambda} \right\rangle \tag{1.36}$$

where the Hamiltonian is parameterized to be a continuous function of $\lambda$ such that $\mathcal{H}(\lambda = 0) = \mathcal{H}_A$ and $\mathcal{H}(\lambda = 1) = \mathcal{H}_B$, $\Delta G$ is the change in free energy, $\langle ... \rangle$ indicates the thermal average, and $U$ is the potential energy of the system. A variety of terms in the forcefield can be set to change with $\lambda$: molecules can be coupled or decoupled from the surrounding solvent molecules, bonds can be formed or broken, and atoms can even be smoothly moved from one element to another in an unphysical process that nonetheless gives meaningful end points. The simulation (or simulations) is then set to obtain sufficient sampling at each value of $\lambda$ to evaluate the integral in equation 1.36.

The Bennet Acceptance Ratio[116] also calculates free energy differences

relative to changes in the forcefield parameters. This approach starts with the following theorem: for any $f(x)$ such that $f(x)/f(-x) = e^{-x}$, the free energy difference between two systems with Hamiltonians $\mathcal{H}_A$ and $\mathcal{H}_B$ is

$$e^{-\beta(\Delta G - C)} = \frac{\langle f(\beta(U_B - U_A - C)) \rangle_A}{\langle f(\beta(U_A - U_B + C)) \rangle_B} \tag{1.37}$$

where $\Delta G$ is the free energy difference between systems with Hamiltonians $\mathcal{H}_A$ and $\mathcal{H}_B$, $U_A$ and $U_B$ are the potential energies of the $\mathcal{H}_A$ and $\mathcal{H}_B$ forcefields respectively, $C$ is an arbitrary constant, and $\langle ... \rangle_A$ and $\langle ... \rangle_B$ are the thermal averages in the $A$ and $B$ ensemble respectively. Equation 1.37 is typically evaluated by running a simulation with Hamiltonian $\mathcal{H}_A$ and a simulation with Hamiltonian $\mathcal{H}_B$. Then in simulation $A$ the quantity $f(\beta(U_B - U_A - C))$ can be computed at each time step and the average over the simulation can be taken to find $\langle f(\beta(U_B - U_A - C)) \rangle_A$. The corresponding quantities for simulation $B$ can be calculated in the same way, and from this the free energy difference is calculated. Bennet showed that the error in the estimation of $\Delta G$ for a given finite sampling is minimized by setting $f(x) = (1 + e^x)^{-1}$ and $C \approx \Delta G$. While we obviously don't know $\Delta G$ *a priori*, this approach allows an iterative scheme to obtain a self-consistent result.

Strictly speaking, equation 1.37 requires that systems $A$ and $B$ occupy the same phase space. This is rarely true in cases of interest, so in practice the full change one wishes to examine is parameterized and broken up into small segments. The free energy to move from state $i$ to state $i + 1$ is calculated along the path and summed to arrive at the total change in free energy. In this way the practical implementation of BAR is very similar to that of TI.

The Jarzynski equality[117] allows the derivation of the change in free energy between two states from many independent simulations determining the work done to move the system between these states. Jarzynski showed that

$$\langle e^{-\beta W} \rangle = e^{-\beta \Delta G} \tag{1.38}$$

where $W$ is the work done in a given simulation, $\langle...\rangle$ is the average over simulations, and $\Delta G$ is the difference in free energies between states. To make use of the Jarzynski equality one performs many simulations in which the system is driven from state A to state B, and the work done in doing so is calculated. One factor complicating this analysis is finding an unbiased estimator for $\langle e^{-\beta W}\rangle$. In particular, given some finite number of simulations $N$, each of gives rise to a work $W_i$, the estimator

$$\sum_{i=1}^{N} e^{-\beta W_i} \tag{1.39}$$

is biased. Unbiased estimators of the change in free energy using the Jarzynski equality have been developed though[118] and the Jarzynski equality is a useful tool in computational biophysics.

The weighted histogram analysis method (WHAM)[119] relies on a technique for estimating the thermal average of an observable in the absence of bias from measurements in the presence of bias. Given $N_{\mathrm{sims}}$ simulations which have different biasing potentials $U_i(z)$ along some reaction coordinate $z$, the probability $P(z)$ of the system being in state $z$ can be calculated by solving self-consistently the following:

$$P(z) = \frac{\sum_i^{N_{\mathrm{sims}}} n_i(z)}{\sum_i^{N_{\mathrm{sims}}} N_i \exp\left(\beta\left[F_i - U_i(z)\right]\right)} \tag{1.40}$$

$$F_i = -k_B T \ln\left(\sum_{\mathrm{bins}} P(z)\exp\left[-\beta U_i(z)\right]\right) \tag{1.41}$$

where $n_i(z)$ is the number of counts in the bin centred on $z$, $N_i$ the total number of counts for simulation $i$, and $F_i$ an energy shift introduced to optimize the estimate of $P(z)$.

The practical implementation of WHAM involves performing many simulations with different biasing potentials that result in a range of values of the reaction coordinate (e.g. the distance between two monomers for a study of binding free energy). The simulations need to be spaced closely enough along the reaction coordinate that the states sampled in different

simulations overlap. Then WHAM can be used to determine the free energy of moving along that reaction coordinate from the knowledge of $P$.

To summarize these methods: Thermodynamic Integration and Bennet Acceptance Ratio are useful in calculating the free energy difference between two states with different forcefield, such as molecules that do not interact with their surroundings in one state but do in the other. Jarzynski allows one to calculate the free energy difference between two states when it is possible to use an external pulling force to move the system from one state to the other, and when this can be done many times to establish an average. WHAM can be used when a path in some coordinate can be defined between the two systems, and simulations at various points along that path can be performed.

It is important to note here that all these methods calculate the *difference* in free energy between two systems. The absolute free energy of a system is a quantity we cannot calculate, nor would it be of particular interest if we could. Only relative free energies matter.

## 1.7 Aims of this Thesis

This thesis investigates protein-solvent interactions using classical density functional theory. In Chapter 2 we examine a variety of experimental data on the transfer free energy to move a protein into various solutions, and in particular the difference in transfer free energy between moving the protein from vacuum to pure water and moving the protein from vacuum to water plus a cosolute. We introduce a new way of looking at the uncertainty in measuring the enthalpy and entropy of such transfers and ascertain whether entropy-enthalpy compensation is a real effect or an experimental artifact. In chapter 4 we introduce classical density functional theory in the context of transfer free energies and prove that the free energy of a bath of particles is independent of external potential under certain conditions. In chapter 3 we apply cDFT to cosulutes and show that even in a very approximate form the theory still produces useful insights. We also re-examine the issue of entropy-enthalpy compensation through the lens of cDFT. Finally in chapter

5 we develop a cDFT implicit solvent model. We finish with a discussion of future directions for this work.

# Chapter 2

# Experimental Analysis: Entropy-Enthalpy Compensation and Cosolutes

As mentioned in Section 1.2, entropy-enthalpy compensation is a phenomenon in which changes in entropy and enthalpy upon some perturbation largely cancel, leaving a change in free energy that is much smaller in magnitude than the changes in entropy and enthalpy. In this chapter we investigate how broadly this effect of entropy-enthalpy compensation applies to macromolecular systems, by analyzing the experimentally-derived enthalpy and entropy of transferring two-state proteins from water, perhaps with buffers and at some pH which need not be 7, into the same solution but in the presence of various *cosolutes*. These cosolutes can be osmolytes, denaturants, crowders, or other proteins; *i.e.* we place no restriction on the size or on how relatively favorable or unfavorable the interactions are with the protein[35].

While the rest of the thesis examines protein-solvent interactions using theoretical and computational tools, in this chapter we examine experimental data on protein-solvent interactions. In addition to addressing entropy-enthalpy compensation, a topic of interest in and of itself, our analysis in this chapter will motivate our later development of an implicit solvent theory that accounts for entropy as well as enthalpy. Many presently available

implicit solvent theories typically do not account for entropy, and are purely enthalpic. We will see in this chapter that entropy and enthalpy enter the transfer free energy on equal footing, and must both be accounted for to create an accurate theory of implicit solvation. We will also return to the data in this chapter in Section 3.5 to show how it can inform our classical DFT approach to solvation in a direct way.

In what follows, we begin by introducing various thermodynamic equations that define the two-state model in Section 2.1. In Section 2.2.1 we introduce a Monte Carlo procedure for estimating the experimental uncertainty of thermodynamic quantities obtained from calorimetry assays. In Section 2.2.2 we show that entropy-enthalpy compensation occurs for the transfer of a diverse set of two-state proteins to various solvents. This may be the most general class of systems that have been observed to obey compensation. While it may be intuitive that a given protein and solvent series may compensate, it is not obvious that there would exist compensation across *both* solvents and proteins. For example, the excluded volume component of transfer is generally non-compensated and different across protein-solvent systems. We use the method derived in Section 2.2.1 to confirm that entropy-enthalpy compensation is a significant effect. In Section 2.2.2 we plot the experimental data at lab temperature, which exhibit definitive entropy-enthalpy compensation across a diverse set of proteins and cosolutes. In this section we emphasize the importance of accounting for the (often neglected) concentration dependence of the heat capacity change upon unfolding. Finally we conclude in Section 2.3.

## 2.1 Methods and Theory

### 2.1.1 Thermodynamic Equations for Protein Unfolding

Two-state models in protein folding have a long and rich history, and have empirical validity for many proteins [120–122]; various aspects of two-state folding, including applications to protein denaturation, protein stability, and the prediction of so-called $m$-values, are described elsewhere [25, 121–132].

Here, we adopt the two-state model for a set of proteins that either have been shown previously to satisfy the van't Hoff two-state criterion[132] or that have comparably small residuals when fit to a two-state model.

The changes in enthalpy $\Delta H$, entropy $\Delta S$, and free energy $\Delta G$ upon unfolding can be obtained if the change of heat capacity upon unfolding $\Delta C_p = C_{pu} - C_{pn}$ is measured. Here $C_{pu}$ and $C_{pn}$ are the unfolded and native state heat capacities respectively, which may be temperature-dependent. A temperature-independent unfolding heat capacity is often used as a good approximation [133, 134], while others have considered a linear temperature-dependence of the unfolding heat capacity [7, 135]. Here, we adopt the most general temperature dependence of the unfolding heat capacity, following the method used in Wintrode *et. al.* [136], wherein the folded heat capacity $C_{pn}$ is observed to obey a linear temperature-dependence, and the unfolded heat capacity obeys a non-linear temperature-dependence determined by the heat capacities of the amino acid constituents. Specifically, the heat capacity of the unfolded state $C_{pu}$ is given by

$$ C_{pu} = (N - N_{gly} - 1)C_p(\text{bb}) + C_p(\text{N/C term}) + \sum_{i=1}^{N} C_p(R_i) \,. \qquad (2.1) $$

Here $N$ is the chain length, $N_{gly}$ is the number of glycine residues in the polypeptide chain, $C_p(bb)$ is the heat capacity of the peptide backbone, $C_p(\text{N/C term})$ is the heat capacity of the N- and C- termini, and $C_p(R_i)$ is the heat capacity of the side chain corresponding to the $i$th amino acid (glycine is included in this sum). Values for $C_p(bb)$, $C_p(\text{N/C term})$, and $C_p(R_i)$ have been obtained by Makhatadze and Privalov [137] for temperatures of 5, 25, 50, 75, 100, and 125° C. For the proteins and cosolutes we consider in Section 2.2.1, we use the values in reference ([137]) to interpolate $C_{pu}(T)$ from 5 to 125° C with a cubic spline.

With $\Delta C_p(T)$ determined numerically, $\Delta H$, $\Delta S$, and $\Delta G$ can be calcu-

lated from

$$\Delta H = \Delta H_f + \int_{T_f}^{T} \Delta C_p(T') dT' \tag{2.2}$$

$$\Delta S = \Delta S_f + \int_{T_f}^{T} \frac{\Delta C_p(T')}{T'} dT' \tag{2.3}$$

$$\Delta G = \Delta H - T\Delta S \tag{2.4}$$

The reference temperature $T_f$ is taken to be the temperature at which the unfolding free energy is zero: $\Delta H(T_f) = T_f \Delta S(T_f)$. The unfolding heat capacity is given by $\Delta C_p(T) = C_{pu}(T) - C_p^n(T)$, with $C_{pu}$ and $C_p^n$ described above. The non-linearity in the heat capacity is fixed in the model by the composition of the protein. The linear temperature-dependence of the native heat capacity is determined empirically from the fit to the data for each protein-cosolute system. Thus when using this model to fit data, the free parameters in the heat capacity are the unfolding heat capacity at the transition temperature, $\Delta C_{pf}$, and the linear coefficient to the temperature dependence of the heat capacity of the native state, $\Delta C'_{pn}$. $\Delta C_p(T)$ is parameterized as $\Delta C_p = \Delta C_{pf} + C_{pu}(T) - C_{pu}(T_f) - \Delta C'_{pn}(T - T_f)$, where $C_{pu}(T)$ has the non-linear $T$ dependence in equation (2.1).

In the approximation that $\Delta C_p$ is a linear function of temperature: $\Delta C_p = \Delta C_{pf} + \Delta C'_p(T - T_f)$, where $\Delta C'_p = \partial \Delta C_p / \partial T$, Equations (2.2)-(2.4) become

$$\Delta H = \Delta H_f + \Delta C_{pf}(T - T_f) + \frac{\Delta C'_p}{2}(T - T_f)^2 \tag{2.5}$$

$$\Delta S = \Delta S_f + \Delta C_{pf} \ln\left(\frac{T}{T_f}\right) + \Delta C'_p \left[T - T_f - T_f \ln\left(\frac{T}{T_f}\right)\right] \tag{2.6}$$

$$\Delta G = \Delta H - T\Delta S \tag{2.7}$$

The expressions for $\Delta H$, $\Delta S$, and $\Delta G$ in the limiting case of a $T$-independent unfolding heat capacity may be obtained by setting $\Delta C'_p = 0$ in (2.5)-(2.7);

*e.g.* the unfolding free energy is

$$\Delta G = \left(1 - \frac{T}{T_f}\right)\Delta H_f + \left[T - T_f - T\ln\left(\frac{T}{T_f}\right)\right]\Delta C_{pf} \qquad (2.8)$$

In Section 2.2.1, we examine the effect these approximations have on the parameters obtained from experimental data.

The probability $p_u$ for the system to be unfolded in the two-state model is

$$p_u = 1/(1 + e^{\beta\Delta G}) \qquad (2.9)$$

with $\Delta G$ given in equation (2.4). Equation (2.9) can be equivalently written as $\Delta G = -k_B T \ln K_u$ where $K_u$ is the folding equilibrium constant, given in the two-state model by $K_u = p_u/(1 - p_u)$.

The total heat capacity in the two-state model is given (for example by differentiating $\langle H \rangle = H_n(1 - p_u) + H_u p_u$ with respect to $T$) by

$$C_p = C_{pu} - \Delta C_p + \frac{\Delta C_p}{1 + e^{\beta\Delta G}} + \frac{\Delta H^2}{4k_B T^2}\text{sech}^2\left(\frac{\Delta G}{2k_B T}\right) \qquad (2.10)$$

where $C_{pu}$, $\Delta C_p$, $\Delta G$, and $\Delta H$ are all temperature-dependent.

A plot of the Gibbs free energy *vs* temperature obtained from empirical data may be fit to Equation (2.4) to obtain values of $\Delta H_f$, $\Delta S_f$, $\Delta C_{pf}$, and $\Delta C'_p$. Similarly, a plot of the fraction of unfolded protein *vs* temperature may be fit to equation (2.9), or a plot of excess heat capacity may be fit to equation (2.10) to obtain values of these parameters. In all cases, once $\Delta H_f$, $\Delta S_f$, $\Delta C_{pf}$, and $\Delta C'_p$ are obtained, equations (2.2) and (2.3) can be used to obtain $\Delta H(T)$ and $\Delta S(T)$ at various temperatures. In Section 2.2.1 we will compare the best-fit values for the three models of $\Delta C_p$ described above, *i.e.* Equations (2.2)-(2.4), (2.5)-(2.7), and the temperature-independent $\Delta C_p$ model (*cf.* Equation (2.8)). This procedure can be performed at various cosolute concentrations, providing experimental data is available. Then the *changes* in unfolding enthalpy $\delta\Delta H(T,c)$ and entropy $\delta\Delta S(T,c)$ upon transfer from a solution of cosolute concentration 0 to one of concentration c at a given temperature $T$ can be determined.

We also define the change in the midpoint parameters at each respective cosolute concentration, nonzero and zero, upon transfer. :

$$\delta\Delta H_f(c) \equiv \Delta H(T_f(c),c) - \Delta H(T_f(0),0) \qquad (2.11)$$

$$\delta\Delta S_f(c) \equiv \Delta S(T_f(c),c) - \Delta S(T_f(0),0) \qquad (2.12)$$

$$\delta\Delta C_{pf}(c) \equiv \Delta C_p(T_f(c),c) - \Delta C_p(T_f(0),0) \qquad (2.13)$$

In what follows we will often drop the explicit concentration dependence in writing various thermodynamic equalities when it is unambiguous.

## 2.2 Results

### 2.2.1 Monte Carlo Method to Determine Statistical Errors

To analyze the uncertainty involved in fitting the data, we perform a Monte Carlo procedure. We fit a given data set, such as $C_p(T)$, $\Delta G(T)$, or $p_u(T)$, to equations (2.10), (2.4), or (2.9) respectively. Using the root mean square of the residuals for a given fit, we then generate a large number of sample data sets by drawing each point from a normal distribution with a mean equal to the value of the best fit curve at that point and a standard deviation equal to the root mean square of the residual. Each of these generated sample data sets is then fitted and new fit parameters are obtained, thus generating a *distribution* of values for $\Delta H_f$, $\Delta S_f$, and, depending on the model, either $\Delta C_p$, $\Delta C_{pf}$ and $\Delta C'_p$, or $\Delta C_{pf}$ and $\Delta C'_{pn}$. We can fit the different models for $\Delta C_p$ described in Section 2.1.1 to compare the parameters extracted. The uncertainty in the thermodynamic parameters could also be obtained by examining the covariance matrix of the fit parameters, but the Monte Carlo method we use allows us to extrapolate uncertainties to other regimes (as we will do in section 2.2.2) without truncating any moments in obtaining the variance.

As an example of fitting the stability $\Delta G(T)$, we have used experimental measurements by Zweifel and Barrick of the thermal denaturation of notch ankyrin in various concentrations of urea[138]. The best fit to the 0M data

34

in figure 4a for reference ([138]) (plotted as green circles in Panel A of Figure 2.1) yields a root mean square of the residuals of 0.383 kJ/mol, so the square of this becomes the variance of the normal random distribution centered around the value of the best fit curve. From this we generate 1000 sample data sets, each of which is fit to either Equation (2.4), (2.7) or (2.8), depending on the model. In Table 2.1, we compare the parameters obtained with the three different models of $\Delta C_p$ described above. We see that the parameters are consistent with each other, and with the tabulated value in reference ([138]). We analyze the variances and correlations of this data; this is reported in Table 2.2. We see that all parameters are generally strongly correlated or strongly anti-correlated. Figures 2.1A,B show that the three models give similar curves even when extrapolating to high temperatures, though the variance in $\Delta G$ at high $T$ is significantly smaller for the $T$-independent $\Delta C_p$ model than the other two models considered.

We perform the same analysis to compare the three heat capacity models for the stability $\Delta G$ *vs* $T$ data for hisactophilin given in reference ([6]). The comparison between the parameters that the three models give for fitting the same data are given in Table 2.1. The midpoint parameters $\Delta H_f$ and $\Delta S_f$ for the different models again all agree within the uncertainties obtained from the Monte Carlo procedure.

Figure 2.1 panels C and D plot data from reference ([6]) for the stability $\Delta G$ vs $T$ for hisactophilin in 0 M and 1 M urea. The data for 1 M urea includes both hot denaturation and cold denaturation regions. Comparing panels C and D of Figure 2.1 we can see that the *model* variance is much less for the 1M urea data, in that all three models predict similar curves, presumably as a result of having a larger range of $\Delta G(T)$ data to fit. The large uncertainty in the non-linear $T$ dependent model for the 1 M urea data is likely caused by fitting a more flexible model to a limited set of data.

Interestingly, there is a change in sign of the curvature at low temperatures in Panel C of Figure 2.1 for the non-linear $T$-dependent model. This effect is caused by a change in sign of $\Delta C_p$ at around 310 K. A similar effect is seen at high temperatures for some generated sets of data in the non-linear $T$-dependent model at 1 M urea (Figure 2.1D).

We have performed the same analysis on data measuring the fraction of unfolded protein as a function of temperature. The data examined is for histidine-containing phosphocarrier protein (HPr), from reference ([139]), and for Arc Repressor, from reference ([140]). The midpoint parameters $\Delta H_f$ and $\Delta S_f$ are again in agreement between the three models (see Table 2.1) but the way that the different heat capacity models extrapolate quantities such as the stability and the enthalpy is markedly different, as it was for hisactophilin— see Figure 2.2 for HPr.

We have performed the same analysis on heat capacity $vs$ $T$ data for $\alpha$-lactalbumin, from reference ([7]). The midpoint parameters $\Delta H_f$ and $\Delta S_f$ are in agreement between the models (Table 2.1). Again, however, the way that the models extrapolate stability and enthalpy is very different (Figure 2.2).

In all the cases we have examined, the values of the unfolding entropy and enthalpy at the transition midpoint are robust across all three models. Further, the value $\Delta C_{pf}$ agrees within uncertainty for all the cases we looked at between the linear $T$-dependent $\Delta C_p$ model and the non-linear $T$-dependent $\Delta C_p$ model. Fitting experimental data to a temperature-independent $\Delta C_p$ model will be sufficient, if only midpoint parameters are required and the accuracy of the unfolding heat capacity $\Delta C_{pf}$ is not particularly important. However Figures 2.1 and 2.2 indicate that such data is prone to significant extrapolation errors.

### 2.2.2 Transfer Entropy and Enthalpy for Various Proteins and Solvents

**Transfer Entropy and Enthalpy at the Transition Midpoint**

The above analysis indicates that the thermodynamic parameters obtained by fitting experimental data are most accurately determined near the transition midpoint. We thus now examine the entropy and enthalpy of transfer for various proteins at their transition midpoints, from water to water plus various cosolutes, $cf.$ Equations (2.11) and (2.12).

For a number of proteins, thermodynamic data exists for more than one

Figure 2.1: **Stability and enthalpy as a function of temperature** for notch ankyrin and hisactophilin. The green circles are experimental data from ref ([138]) for notch ankyrin (panels A and B) and ref ([6]) for hisactophilin (panels C and D). The blue lines are fits for the $T$-independent $\Delta C_p$ model, the red lines are fits for the linear $T$-dependent model, and the black lines are fits for the non-linear $T$-dependent model (*cf.* Section 2.1.1). The solid lines arise from the best fit parameters for each model, while the dashed lines represent one standard deviation away, determined by the Monte Carlo procedure described in Section 2.2.1. (Panel A) Stability for notch ankyrin *vs* temperature in buffer. (Panel B) Enthalpy for notch ankyrin in buffer. (Panel C) Stability for hisactophilin in buffer. (Panel D) Stability for hisactophilin in buffer with 1M urea. The insets in Panels A, C, and D show the correlation of the midpoint enthalpy and entropy, in which the models are represented by the same colors as above. All data in the insets lies on top of the red scatter points; the blue and black points have been displaced for clarity. 1000 Monte Carlo instances have been generated for the inset plots. Bars indicate one standard deviation.

Figure 2.2: **Analyzing unfolding fraction data and heat capacity data.** Panels A and B show fraction of unfolded population data and heat capacity data from refs ([139]) and ([7]) respectively (green circles), along with best fit curves. Panels C and D show the corresponding stability as a function of temperature, and panels E and F show the corresponding enthalpy as a function of temperature. In all panels the blue lines are fits for the $T$-independent $\Delta C_p$ model, the red lines are fits for the linear $T$-dependent model, and the back lines are fits for the non-linear $T$-dependent model. The solid lines arise from the best fit parameters for each model while the dashed lines represent one standard deviation away, determined by the Monte Carlo procedure described in Section 2.2.1.

Table 2.1: **Thermodynamic parameters for several proteins** used in our analysis, and comparison to literature data where available: $\alpha$-Lactalbumin, from heat capacity data in ref. ([7]), Arc Repressor, from fraction of unfolded protein data in ref. ([140]), Creatine Kinase from heat capacity data in ref. ([141]), Hisactophilin, from stability data in ref. ([6]), Histidine-containing phosphocarrier protein (HPr), from fraction of unfolded protein data in ref. ([139]), Notch Ankyrin, from stability data in ref. ([138]), and RNase A from heat capacity data in ref. ([142]). Values obtained from the three models of the temperature dependence of $\Delta C_p$ are given, as well as the values obtained from the appropriate reference where available. The reference value in ref. ([138]) assumed a temperature-independent $\Delta C_p$, the value of $\Delta H_f$ from ref. ([6]) was obtained by integrating $C_p$ up to $T_f$, and the values from ref. ([139]) were obtained assuming a temperature-independent $\Delta C_p$.

| Protein | Model | $\Delta H_f$ (kJ/mol) | $\Delta S_f$ (kJ/mol/K) | $\Delta C_{pf}$ (kJ/mol/K) | $\Delta C'_p$ (kJ/mol/K$^2$) | $\Delta C'_{pn}$ (kJ/mol/K$^2$) |
|---|---|---|---|---|---|---|
| $\alpha$-Lactalbumin | $T$ independent | $297 \pm 1$ | $0.876 \pm 0.003$ | $4.38 \pm 0.07$ | – | – |
| | $T$ linear | $304 \pm 1$ | $0.896 \pm 0.004$ | $2.88 \pm 0.26$ | $-0.101 \pm 0.017$ | – |
| | $T$ non-linear | $304 \pm 1$ | $0.895 \pm 0.004$ | $2.93 \pm 0.26$ | – | $0.065 \pm 0.017$ |
| | Reference Value[7] | 310 | – | 5.3 | -0.05 | – |
| Arc Repressor | $T$ independent | $139 \pm 3$ | $0.454 \pm 0.011$ | $0.536 \pm 0.5$ | – | – |
| | $T$ linear | $118 \pm 7$ | $0.385 \pm 0.02$ | $1.45 \pm 1.5$ | $1.18 \pm 0.44$ | – |
| | $T$ non-linear | $117 \pm 7$ | $0.384 \pm 0.024$ | $1.49 \pm 1.4$ | – | $1.15 \pm 0.4$ |
| Creatine Kinase* | $T$ independent | $780 \pm 3$ | $2.37 \pm 0.01$ | $25 \pm 2$ | – | – |
| | $T$ linear | $733 \pm 3$ | $2.24 \pm 0.01$ | $66 \pm 4$ | $-73 \pm 3$ | – |
| Hisactophilin | $T$ independent | $215 \pm 7.4$ | $0.659 \pm 0.022$ | $6.15 \pm 0.75$ | – | – |
| | $T$ linear | $218 \pm 5.9$ | $0.669 \pm 0.018$ | $12.7 \pm 1.7$ | $-0.610 \pm 0.17$ | – |
| | $T$ non-linear | $218 \pm 17$ | $0.667 \pm 0.051$ | $12.8 \pm 3.3$ | – | $0.784 \pm 0.28$ |
| | Reference Value[6] | 226 | – | – | – | – |
| HPr | $T$ independent | $312 \pm 2$ | $0.929 \pm 0.006$ | $5.27 \pm 0.7$ | – | – |
| | $T$ linear | $314 \pm 2$ | $0.935 \pm 0.006$ | $4.41 \pm 0.85$ | $0.359 \pm 0.150$ | – |
| | $T$ non-linear | $314 \pm 2$ | $0.935 \pm 0.006$ | $4.44 \pm 0.80$ | – | $0.367 \pm 0.14$ |
| | Reference Values[139] | 316 | 0.941 | 6.0 | – | |
| Notch Ankyrin | $T$ independent | $593 \pm 9$ | $1.86 \pm 0.03$ | $15.1 \pm 0.5$ | – | – |
| | $T$ linear | $602 \pm 30$ | $1.89 \pm 0.1$ | $16.2 \pm 3.5$ | $-0.045 \pm 0.14$ | – |
| | $T$ non-linear | $601 \pm 33$ | $1.89 \pm 0.1$ | $15.7 \pm 3.5$ | – | $0.114 \pm 0.14$ |
| | Reference Value[138] | – | – | 15.1 | – | – |
| RNase A* | $T$ independent | $496 \pm 1$ | $1.477 \pm 0.002$ | $14.3 \pm 0.1$ | – | – |
| | $T$ linear | $468 \pm 1$ | $1.396 \pm 0.003$ | $22.6 \pm 0.2$ | $0.89 \pm 0.02$ | – |
| | Reference Values[142]$^\ddagger$ | 515 | 1.52 | – | – | – |
| | Reference Values[142]$^\S$ | 479 | 1.42 | – | – | – |

– Not applicable for the respective model.
$^*$ Literature data had background heat capacity subtracted for these proteins, so the non-linear temperature-dependent model could not be applied.
$^\ddagger$ Literature values obtained from differential scanning calorimetry. $^\S$ Literature values obtained from spectroscopy measurements.

Table 2.2: **Comparison of the variance and covariance of fits to** $\Delta G$ *vs* $T$ data for Notch Ankyrin from ref. ([138]), for the three models of the temperature dependence of $\Delta C_p$ discussed. For each model, a matrix is given in which the diagonal elements are the relative deviations for that quantity, and the off-diagonal elements are the correlation coefficients for the two quantities. Relative deviation for *e.g.* $\Delta H_f$ is defined as $(\langle \Delta H_f^2 \rangle - \langle \Delta H_f \rangle^2)^{1/2}/\langle \Delta H_f \rangle$, where averages are over the Monte Carlo generated data. In all models the entropy and enthalpy of unfolding are highly correlated.

| | | $\Delta H_f$ | $\Delta S_f$ | $\Delta C_p$ | |
|---|---|---|---|---|---|
| $T$-independent $\Delta C_p$ | $\Delta H_f$ | 0.016 | 0.9997 | 0.989 | |
| | $\Delta S_f$ | | 0.017 | 0.989 | |
| | $\Delta C_p$ | | | 0.033 | |
| | | $\Delta H_f$ | $\Delta S_f$ | $\Delta C_{pf}$ | $\Delta C_p'$ |
| $T$-linear $\Delta C_p$ | $\Delta H_f$ | 0.055 | 0.99997 | 0.987 | -0.960 |
| | $\Delta S_f$ | | 0.056 | 0.988 | -0.959 |
| | $\Delta C_{pf}$ | | | 0.22 | -0.991 |
| | $\Delta C_p'$ | | | | 3.0 |
| | | $\Delta H_f$ | $\Delta S_f$ | $\Delta C_{pf}$ | $\Delta C_{pn}'$ |
| $T$-non-linear $\Delta C_p$ | $\Delta H_f$ | 0.055 | 0.99997 | 0.987 | -0.958 |
| | $\Delta S_f$ | | 0.056 | 0.988 | -0.957 |
| | $\Delta C_{pf}$ | | | 0.22 | -0.990 |
| | $\Delta C_{pn}'$ | | | | 1.20 |

cosolute; as well, for a number of cosolutes thermodynamic data exists for more than one protein. These commonalities and differences enable useful comparisons.

Because $\Delta H_f = T_f \Delta S_f$, the folding temperature $T_f$ is unchanged upon transfer to cosolute if $\delta \Delta H_f / \Delta H_f^0 = \delta \Delta S_f / \Delta S_f^0$ where $\Delta H_f^0$ and $\Delta S_f^0$ are the mid-point enthalpy and entropy of unfolding at cosolute concentration $c = 0$, respectively, and $\delta \Delta H_f$ and $\delta \Delta S_f$ are defined in Equations (2.11) and (2.12). Thus, on a plot with $\delta \Delta H_f / \Delta H_f^0$ and $\delta \Delta S_f / \Delta S_f^0$ on the ordinate and abscissa, a line of slope unity would constitute no change in folding

temperature due to the cosolute. Further, since

$$T_f^0 + \delta T_f = \frac{\Delta H_f^0 + \delta \Delta H_f}{\Delta S_f^0 + \delta \Delta S_f} \approx T_f^0 + T_f^0 \left( \frac{\delta \Delta H_f}{\Delta H_f^0} - \frac{\delta \Delta S_f}{\Delta S_f^0} \right)$$

we conclude that the distance of each point from the line $y = x$ is, to first order, the relative change in the folding temperature, $\delta T_f / T_f^0$, as a result of the transfer.

We may further interpret the deviation from the line $y = x$ on a plot of $\delta \Delta S_f / \Delta S_0$ *vs* $\delta \Delta H_f / \Delta H_0$ as the free energy of transfer at fixed temperature $T_f^0$, $\delta \Delta G(T_f^0)$, divided by the unfolding enthalpy $\Delta H_f^0$. That is, $\delta \Delta H_f / \Delta H_f^0 - \delta \Delta S_f / \Delta S_f^0 = \delta \Delta G(T_f^0) / \Delta H_f^0$.

In Figures 2.4 and 2.5, we divide both $\delta \Delta H_f / \Delta H_f$ and $\delta \Delta S_f / \Delta S_f$ by the concentration $c$ of the cosolute, to compare cosolutes solutions having different concentrations. Thus the axes in Figures 2.4 and 2.5 can be thought of as a decomposition of $m$-values[143] into enthalpic and entropic components, each normalized by the corresponding unfolding enthalpy or entropy in the absence of cosolute.

Linear regression to the data in Figures 2.4 and 2.5, when taken together gives a slope of $0.99 \pm 0.04$. The statistical test outlined in Krugg *et. al.* [26] requires that the slope of the best fit $\delta \Delta H_f / \Delta H_f^0$ *vs.* $\delta \Delta S_f / \Delta S_f^0$ line be more than $2\sigma$ away from the harmonic mean of the temperatures at which the experiments were performed; thus for Figure 2.4 this requires that the slope be $2\sigma$ away from unity. Figure 2.4 fails this test. Nevertheless, we show that the results in Figure 2.4 are in fact statistically significant, given the magnitude of the experimental errors. We now describe a treatment of the statistical significance of entropy-enthalpy compensation that is valid when the slope of the $\Delta H$-$\Delta S$ plot is near unity.

The Monte Carlo method in Section 2.2.1 can be applied to data at various concentrations of cosolutes to assess the significance of the linear relationship between enthalpy and entropy observed in Figures 2.4-2.5. Fitting each experimental data set to the appropriate equation as described in Section 2.2.1 results in a set of best fit parameters, as well as a set of

residuals from the best fit. For the following protein/cosolute systems—$\alpha$-lactalbumin in ethanol[7], arc repressor in KCl[140], creatine kinase in glycerol[141], hisactophilin in urea[6], histidine containing phosphocarrier in urea[139] notch ankyrin in urea[138], and RNase A in urea[142] — 1000 data sets were generated as described in Section 2.2.1, for each of several concentrations of cosolute. Thus 1000 values for $\delta\Delta H_f$ and $\delta\Delta S_f$ were generated. We plot the results of this procedure as scatter points in Figures 2.4-2.6.

If we take the average of the extent of the scatter for these six data points as an estimate for the experimental uncertainty, and apply it to all other data in Figure 2.4, we can assess the significance of the apparent linear relationship in the plot. The analysis rests on the assumption that the deviations from the line of slope unity, $\delta T_f = 0$, are much smaller than the deviations from zero of either $\delta\Delta H_f$ or $\delta\Delta S_f$, *i.e.*, that the data are essentially distributed along the diagonal. We may then consider the data as transformed to a coordinate system that is rotated $\pi/4$ counterclockwise, and translated so that the origin coincides with the mean of the data. Then the data consists approximately of points distributed along the abscissa all having zero ordinate. If the variance of this data is large compared to what would be expected from the experimental error as derived from the above Monte-Carlo method, then the result of entropy-enthalpy compensation is significant.

To assess the significance we use a bootstrapping method, to avoid requiring an assumption as to the distribution of the points along the $y = x$ line. From the 48 data points in Figures 2.4 and 2.5 we perform random sampling with replacement to generate new sets of data (also with $n = 48$). We find the standard deviation of each of these generated sets and thus obtain a distribution in $\sigma$. We obtain another distribution in $\sigma$ by sampling from the Monte Carlo generated deviations. From the 6000 total generated points (1000 for each of 6 proteins) we sample with replacement to obtain many sets of 1000 deviations. The standard deviations of these sets forms another distribution in $\sigma$. Distributions formed with this procedure are plotted in figure 2.3. The overlap of these two distributions then provides the significance–the likelihood that the scatter in the experimental data arises

from the fit uncertainty. For the data in Figures 2.4 and 2.5 we obtain $p < 10^{-6}$. With near certainty, these results illustrate entropy-enthalpy compensation rather than experimental error.

For almost all the points in Figures 2.4 and 2.5 the magnitude of the enthalpy change is larger than the magnitude of the entropy change. None of the systems we examined showed a destabilizing cosolute with a change in entropy of unfolding larger than the change in enthalpy of unfolding, and only a few systems showed a stabilizing cosolute with a change in entropy of unfolding larger than the change in enthalpy of unfolding. Thus for the most part enthalpy drives the change in stability, while entropy tries to catch up and partially compensates.

Figure 2.5 also shows the transfer enthalpy and entropy for simulation results by O'Brien *et al.* on Cold shock protein and protein L [144] (open black symbols in Figure 2.5). Here, thermodynamic parameters were extracted from fits to simulated heat capacity curves, for the transfer of the above proteins to either urea or TMAO. A surface-area based Tanford transfer model [145] was used to model the cosolute solution. We have not found experimental values for these thermodynamic parameters in the literature; the values in Figure 2.5 are thus predictions as a consequence of both the simulation method for generating unfolded ensembles, and the Tanford transfer model, which are subject to experimental test.

**Importance of the Cosolute- and Temperature-Dependence of the Unfolding Heat Capacity**

The unfolding heat capacity $\Delta C_p(T, c)$ is generally both temperature- and concentration-dependent. We define the change in unfolding heat capacity upon transfer, $\delta \Delta C_{pf}$, in Equation (2.13). While the quantities $\delta \Delta H_f$ and $\delta \Delta S_f$ in Equations (2.11) and (2.12) are independent of $\delta \Delta C_{pf}$, since they are always evaluated at the respective transition temperatures, the thermodynamics for transfer at fixed temperature (e.g. lab temperature) does depend on $\delta \Delta C_{pf}$.

A number of the works cited here obtained a concentration-independent

43

$\Delta C_p$ however, by equating $\Delta C_p$ to the slope of $\Delta H_f$ vs. $T_f$ data for various cosolute concentrations. This assumes that $\Delta C_p$ is constant with varying cosolute concentration, and hence $\delta \Delta C_{pf} = 0$. These proteins/cosolutes are thus indicated in Table 2.3 by "$0^\dagger$" in the column for $\delta \Delta C_{pf}$. This assumption may be sufficient if $\delta \Delta H_f$, $\delta \Delta S_f$, or $\Delta C_p(c = 0)$ is the quantity of interest, however when Equations (2.2) and (2.3) are used to evaluate thermodynamic parameters at lab temperature, this assumption produces unacceptably large errors. Examples of the change in thermodynamic values obtained by setting $\delta \Delta C_{pf} = 0$ are shown in Figure 2.6, for the transfer of Hisactophilin and RNase A to urea: neglecting $\delta \Delta C_p$ changes the resulting value of $\delta \Delta H$ by $\approx 80$ kJ/mol for Hisactophilin and $\approx 40$ kJ/mol for RNase A. This is to be compared with the error introduced by neglecting the temperature dependence of $\Delta C_p$, which was $\approx 30$ kJ/mol for Hisactophilin and $\approx 10$ kJ/mol for RNase A.

In Table 2.3 we have made a note of where the $\delta \Delta C_p = 0$ assumption has been made, and we have not plotted the corresponding $\delta \Delta H_{\text{lab}}$ and $T \delta \Delta S_{\text{lab}}$ data in Figure 2.6. Note that the potential for large errors due to $\delta \Delta C_p$ is irrelevant when comparing data at the respective folding temperatures, i.e. for the quantities in Equations (2.11) and (2.12), and Figures 2.4 and 2.5.

Figure 2.7 plots the concentration-dependence of $\Delta C_{pf}$ for several protein-cosolute systems, obtained by using the non-linear temperature-dependent model in equation 2.1. The values plotted do not change significantly if the linear temperature-dependent model is used (see for example Table 2.1, which shows that the values obtained from the two models are comparable). Some proteins have a $\Delta C_{pf}$ showing weak concentration-dependence (e.g. Barstar in GdmHCl, Acylphosphatase in urea), while for others, $\Delta C_{pf}$ shows significant concentration-dependence (RNase A in urea, $\alpha$-Lactalbumin in ethanol). Furthermore, $\Delta C_{pf}$ need not even be monotonic in $T$; the non-monotonic behaviour exhibited by RNase A in urea is well beyond what can be explained by the experimental uncertainty. Factoring in the inter-model uncertainty does not change this; as panel B shows the non-monotonic behaviour is present in both the T-independent and linear T-dependent models. A concentration-independent heat capacity is often obtained from

linear fits of the unfolding enthalpy *vs.* melting temperature for various osmolyte concentrations. For the proteins in Figure 2.7 that have strong $c$-dependence, this would be a recipe prone to large errors.

### Transfer Entropy and Enthalpy at Lab Temperature

For 19 proteins and cosolutes that we had investigated, the concentration dependence of $\Delta C_{pf}$ is known. For these proteins, we have obtained the transfer enthalpy of unfolding $\delta \Delta H(T = 25°C)$ and the transfer entropy of unfolding $\delta \Delta S(T = 25°C)$ at lab temperature; values are tabulated in Table 2.3. For 12 protein-cosolute systems, thermodynamic parameters at the folding temperatures corresponding to different cosolute concentrations were tabulated in the literature. For these systems, a temperature-independent $\Delta C_p$ model was invariably used to obtain the tabulated values. We thus had to also assume a temperature-independent $\Delta C_p$ model in order to extrapolate the thermodynamic values to lab temperature. We show below however that this procedure may be prone to large errors.

Seven references contained plotted data, which we had fitted to obtain thermodynamic parameters. For 5 of these protein-cosolute systems, the non-linear temperature-dependent $\Delta C_p$ model was used to extrapolate to lab temperature. Two of these systems had baselines subtracted in the published data, so a linear temperature-dependent $\Delta C_p$ model was used to extrapolate to lab temperature. All of these 7 protein-cosolute systems show scatter due to our Monte Carlo procedure that is indicated in Figure 2.6 (though for $\alpha$-Lactalbumin in ethanol and RNase A in urea the scatter is small).

The extent of the scatter in Figure 2.6 makes it clear that for any of the three methods of obtaining $\delta \Delta H$ and $\delta \Delta S$ (heat capacity, fraction unfolded, or unfolding free energy, with the method for each protein indicated in Table 2.3), the uncertainties are highly correlated and can be quite large. There is very little scatter orthogonal to the lines of constant stability on the $\delta \Delta H$-$T\delta \Delta S$ plot. The scatter *along* the equi-stability line is significantly larger however; for arc repressor in particular, the scatter is large enough to

render the sign of $\delta\Delta H$ and $T\delta\Delta S$ uncertain. The scatter in the data for RNase transfer to urea is quite small on the other hand, even though the scatter in the data for creatine kinase in glycerol, also from heat capacity measurements, is large. The average standard deviation along the diagonal was 21 kJ/mol, compared with an average of error of 0.32 kJ/mol perpendicular to the diagonal, corresponding to the change in the unfolding free energy upon transfer $\delta\Delta G$.

We apply the same procedure described in Section 2.2.2 to evaluate the significance of the entropy-enthalpy compensation here. In this case there are 23 data points in Figure 2.6, and the data themselves have a sample standard deviation of about $s \approx 65$ kJ/mol, whereas the mean Monte-Carlo standard deviation applied to each data point is about $\sigma \approx 19$ kJ/mol. Bootstrapping with the same procedure described in section 2.2.2 rejects the hypothesis that the scatter arises from the uncertainty in fitting with a significance $p = 5 \times 10^{-5}$. The result in Figure 2.6 thus also illustrates entropy-enthalpy compensation rather than experimental error.

One caveat of the significance is that the experimental uncertainty only included the fit uncertainty, and not the model to model uncertainty. At lab temperature the model to model uncertainty is approximately a factor of two larger than the fit uncertainty. To assess the significance of entropy-enthalpy compensation with model to model uncertainty factored in, we perform another bootstrapping procedure. This time we create a distribution in $\sigma$ by sampling many sets of 6 standard deviations from the 6 deviations found from the model to model uncertainty. This distribution is then compared to the distribution arising from bootstrapping from the experimental lab temperature data. With this consideration the significance from the bootstrapping method drops to $p = 0.085$.
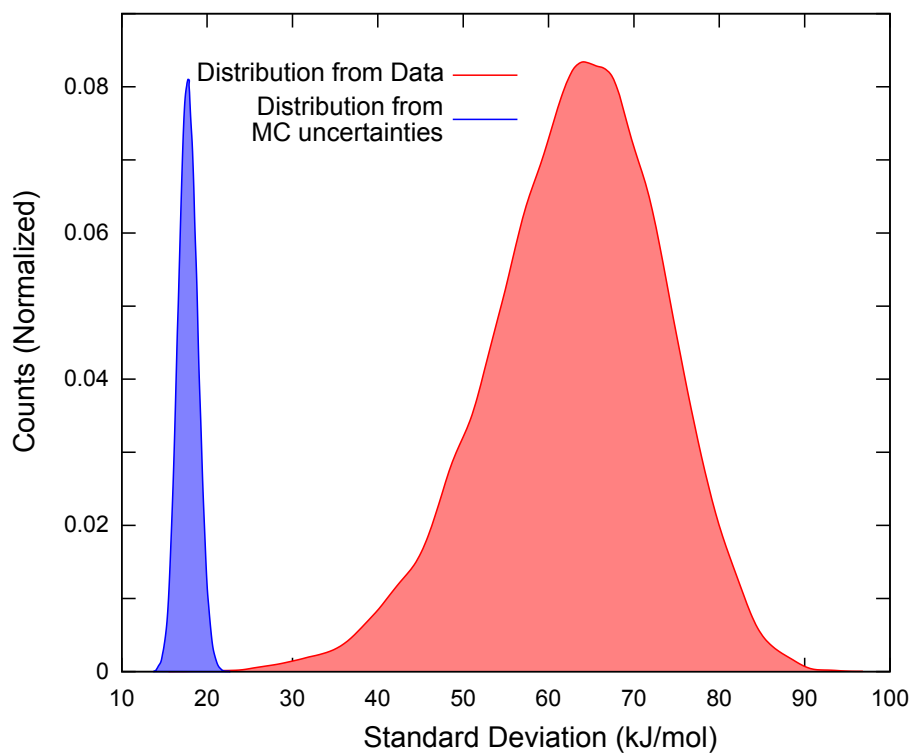
Figure 2.3: Distribution in $\sigma$ arising from the bootstrapping method described in section 2.2.2. The overlap of the two distributions represents the statistical significance of entropy-enthalpy compensation–in this case $p = 5 \times 10^{-5}$.

Table 2.3: **Empirical values for thermodynamic unfolding parameters** upon transfer of various proteins to various solvents. For each protein-cosolute system, the values are listed at the concentration corresponding to (100g/l). The literature reference from which the values were obtained, along with the corresponding figure or table in that reference, is listed in the last column. For systems in which the parameters were obtained through fitting a curve to the data in the reference work, the equation from this work used for the fitting is also listed below the table. Systems without an equation listed were those in which $\Delta H_f$ and $\Delta S_f$ values were directly available.

| Protein | Cosolute | pH | $\Delta H_f^{0**}$ | $\Delta S_f^0$ | $T_f$ | $\Delta C_{pf}$ | $\frac{\delta\Delta H_f}{\Delta H_f^0}$ | $\frac{\delta\Delta S_f}{\Delta S_f^0}$ | $\delta\Delta C_{pf}$ | $\delta T_f$ | $\delta\Delta H_{\mathrm{lab}}$ | $T\delta\Delta S_{\mathrm{lab}}$ | Ref/Fig/Tab[§] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\alpha$-CTgen[$\alpha$] | Trehalose | 2.5 | 403 | 1.27 | 317 | 4.4 | 6.7e-3 | 9.7e-4 | 0[†] | 1.8 | 2.69[†] | 0.37[†] | [146] Tab. 1 |
| $\alpha$-lactalbumin[*] | Ethanol | 8 | 310 | 0.916 | 338 | 5.3 | -8.2e-3 | -5.9e-3 | 0.10 | 19 | 1.11 | 1.67 | [7] Fig. 3[c] |
| $\alpha$-lactalbumin | TMAO | 7 | 209 | 0.663 | 315 | 6.5 | 0.27 | 0.25 | -0.133 | 5.0 | 22.5 | 17.4 | [147] Tab. 1 |
| Acylphosphatase | Urea | 5.5 | 351 | 1.06 | 331 | 6.2 | -0.11 | -0.10 | -0.037 | -3.7 | -4.17 | 1.07 | [148] Tab. 2 |
| Arc Repressor[*] | KCl | 4 | 141 | 0.463 | 304 | 1.0 | 0.25 | 0.14 | 5.20 | 29 | -89.1 | -102 | [140] Fig. 7[a] |
| Barstar | GdnHCl | 8 | 292 | 0.849 | 343 | 6.2 | -0.43 | -0.41 | -1.9 | -12 | -16.2 | -5.64 | [149] Tab. 2 |
| Cro Protein | Urea | 6 | 195 | 0.591 | 330 | 3.8 | -0.22 | -0.20 | 0.09 | -8.2 | -7.47 | -3.45 | [150] Tab. 2 |
| Cytochrome c | Sorbitol | 2 | 226 | 0.740 | 305 | 5.2 | 0.12 | 0.09 | -1.7 | 8.4 | 13.6 | 7.28 | [151] Tab. 1 |
| Cytochrome c | Trehalose | 7 | 161 | 0.502 | 321 | N/A | -3.4e-2 | -3.7e-2 | 0[†] | 24 | -5.38[†] | -4.78[†] | [146] Tab. 1 |
| Creatine Kinase[*] | glycerol | 8.05 | 782 | 2.38 | 329 | 92 | 2.0e-2 | 1.9e-2 | 2.46 | 0.32 | -93.1 | -88.4 | [141] Fig. 4[c] |
| De Novo $\alpha$ B | GdnHCl | 7.3 | 103 | 0.300 | 343 | 2.3 | -0.53 | -0.52 | 0.21 | -7.2 | 0.564 | 4.52 | [152] Tab. 2 |
| De Novo $\alpha$ C | GdnHCl | 7.3 | 153 | 0.441 | 347 | 2.7 | -0.48 | -0.47 | -0.07 | -6.5 | 3.94 | 10.5 | [152] Tab. 2 |
| Hisactophilin[*] | Urea | 5.8 | 215 | 0.658 | 327 | 6.1 | -0.89 | -0.85 | 3.63 | -87 | -142 | -123 | [6] Fig. 6[b] |
| Hexokinase | Glucose | 8 | 700 | 2.19 | 320 | 30 | 0.51 | 0.46 | 1.5 | 11 | -30.0 | -58.1 | [153] Tab. 2 |
| HPr[$\beta$] | Urea | 7 | 315 | 0.935 | 336 | 4.4 | -0.16 | -0.16 | -0.0054 | -27.1 | 36.4 | 39.7 | [139] Fig. 2[a] |
| Lectin (Pea) | Urea | 7.2 | 1130 | 3.25 | 347 | 22 | -1.6e-2 | -1.1e-2 | 0.74 | -1.8 | -7.12 | -3.40 | [154] Tab. 2 |
| Lysozyme | DMSO | 2.5 | 535 | 1.58 | 339 | 7.8 | -4.4e-3 | -3.4e-3 | 0[†] | -034 | 1.00[†] | 1.50[†] | [155] Tab. 2 |
| Lysozyme | Trehalose | 7 | 397 | 1.20 | 331 | N/A | -1.5e-2 | -2.2e-2 | 0[†] | 2.4 | -5.99[†] | -7.09[†] | [146] Tab. 1 |
| Lysozyme | TMAO | 6 | 535 | 1.50 | 357 | 6.8 | 7.5e-2 | 6.5e-2 | 0.089 | 3.3 | 12.3 | 5.59 | [147] Tab. 1 |
| Notch Ankyrin[*] | Urea | 8 | 592 | 1.86 | 318 | 15 | -0.45 | -0.43 | -1.26 | -11 | -63.2 | -43.8 | [138] Fig. 4[b] |
| RNase A | $\beta$-hydrox[$\gamma$] | 5.5 | 364 | 1.09 | 334 | 4.4 | 2.8e-2 | 2.4e-2 | 0.194 | 1.3 | -3.16 | -4.36 | [156] Tab. 1 |
| RNase A | Betaine | 5.5 | 364 | 1.09 | 334 | 4.4 | 4.2e-2 | 4.1e-2 | 0.214 | 0.32 | 5.39 | 3.42 | [156] Tab. 1 |
| RNase A | Betaine | 6.0 | 364 | 1.09 | 334 | 0 | 5.3e-2 | 5.0e-2 | 0[†] | 0.95 | 0[†] | -0.291[†] | [157] Figs. 2,4[d] |
| RNase A | Trehalose | 7 | 385 | 1.20 | 321 | 4.7 | 1.4e-2 | 9.1e-3 | 0[†] | 1.6 | 5.51[†] | 3.38[†] | [146] Tab. 1 |
| RNase A | Glycine | 6.0 | 364 | 1.09 | 334 | 0 | -2.0e-2 | -2.7e-2 | 0[†] | 2.4 | 0[†] | -0.659[†] | [157] Figs. 2,4[d] |
| RNase A | Sarcosine | 6.0 | 364 | 1.09 | 334 | 0 | 1.1e-2 | -4.0e-4 | 0[†] | 3.8 | 0[†] | -1.20[†] | [157] Figs. 2,4[d] |

| Protein | Cosolute | pH | $\Delta H_f^0$ | $\Delta S_f^0$ | $T_f$ | $\Delta C_{pf}$ | $\frac{\delta\Delta H_f}{\Delta H_f^0}$ | $\frac{\delta\Delta S_f}{\Delta S_f^0}$ | $\delta\Delta C_{pf}$ | $\delta T_f$ | $\delta\Delta H_{\text{lab}}$ | $T\delta\Delta S_{\text{lab}}$ | Ref/Fig/Tab[§] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RNase A | TMAO | 7 | 490 | 1.46 | 336 | 5.2 | 7.3e-2 | 6.2e-2 | 0.022 | 3.5 | 16.1 | 9.35 | [147] Tab. 1 |
| RNase A[Cal] | GdnHCl | 7 | 515 | 1.53 | 337 | 11 | -0.23 | -0.22 | 0[†] | -4.3 | -0.246[†] | 10.8[†] | [142] Tab. 1 |
| RNase A[UV] | GdnHCl | 7 | 452 | 1.35 | 335 | 6.3 | -0.16 | -0.14 | 0[†] | -7.8 | 0.44[†] | 12.3[†] | [142] Tab. 1 |
| RNase A[Cal] | Methylurea | 7 | 515 | 1.53 | 337 | 7.1 | -5.3e-2 | -4.5e-2 | 0[†] | -2.8 | -0.508[†] | 4.19[†] | [142] Tab. 1 |
| RNase A[UV] | Methylurea | 7 | 452 | 1.35 | 335 | 4.8 | -4.4e-2 | -3.4e-2 | 0[†] | -3.5 | -0.419[†] | 4.28[†] | [142] Tab. 1 |
| RNase A[Cal] | Dimethylurea[δ] | 7 | 515 | 1.53 | 337 | 3.7 | -2.8e-2 | -1.8e-2 | 0[†] | -3.4 | -0.08[†] | 5.04[†] | [142] Tab. 1 |
| RNase A[UV] | Dimethylurea[δ] | 7 | 452 | 1.35 | 335 | 1.2 | -8.7e-2 | 3.2e-2 | 0[†] | -39 | 0.79[†] | 5.64[†] | [142] Tab. 1 |
| RNase A[Cal] | Ethylurea | 7 | 515 | 1.53 | 337 | 3.6 | -4.0e-2 | -2.7e-2 | 0[†] | -4.5 | 0.13[†] | 6.02[†] | [142] Tab. 1 |
| RNase A[UV] | Ethylurea | 7 | 452 | 1.35 | 335 | 3.9 | -4.2e-2 | -2.7e-2 | 0[†] | -5.2 | 1.05[†] | 8.25[†] | [142] Tab. 1 |
| RNase A[Cal] | Butylurea | 7 | 515 | 1.53 | 337 | 6.0 | -0.21 | -0.17 | 0[†] | -16 | 4.3[†] | 23.4[†] | [142] Tab. 1 |
| RNase A[UV] | Butylurea | 7 | 452 | 1.35 | 335 | 7.3 | -0.26 | -0.21 | 0[†] | -21 | 7.4[†] | 30.3[†] | [142] Tab. 1 |
| RNase A[*] | Urea | 7 | 501 | 1.49 | 336 | 14.6 | -4.0e-2 | -2.9e-2 | -1.5 | -3.8 | 43.0 | 46.6 | [142] Fig. 1[c] |
| Ssh10b | GdnHCl | 6.8 | 307 | 0.840 | 365 | N/A | -0.18 | -0.14 | 0[†] | -17 | -54.4[†] | -34.3[†] | [158] Tab. 2 |
| Tryps inh[ε] | Trehalose | 7 | 236 | 0.711 | 332 | 1.1 | 8.8e-3 | 3.4e-3 | 0[†] | 1.8 | 2.08[†] | 0.74[†] | [146] Tab. 1 |
| Ubiquitin | PVP | 5.4 | 100 | 0.265 | 377 | 1.5 | -0.2 | -0.2 | -0.1 | 0 | -26.0 | -14.8 | [159] Tab. 3 |
| Ubiquitin | Ficoll | 5.4 | 100 | 0.265 | 377 | 1.5 | -0.5 | -0.6 | -0.5 | 94 | -36.1 | -44.5 | [159] Tab. 3 |
| Ubiquitin | BSA | 5.4 | 100 | 0.265 | 377 | 1.5 | 0.5 | 0.5 | 1.6 | 0 | -310 | -295 | [159] Tab. 3 |
| Ubiquitin | Lysozyme | 5.4 | 100 | 0.265 | 377 | 1.5 | -0.1 | -0.2 | 0.1 | 47 | -88.1 | -107 | [159] Tab. 3 |
| Ubiquitin | NaCl | 2 | 203 | 0.617 | 329 | 3.0 | 0.547 | 0.383 | 0[†] | 39 | 111[†] | 70.4[†] | [160] Tab. 1 |
| Ubiquitin | CaCl$_2$ | 2 | 203 | 0.617 | 329 | 4.6 | 3.21 | 2.67 | 0[†] | 48 | 652[†] | 490[†] | [160] Tab. 1 |
| Ubiquitin | MgCl$_2$ | 2 | 203 | 0.617 | 329 | 4.4 | 1.87 | 1.49 | 0[†] | 50 | 379[†] | 273[†] | [160] Tab. 1 |
| Ubiquitin | GdmCl | 2 | 203 | 0.617 | 329 | 4.7 | 0.248 | 0.201 | 0[†] | 13 | 50.3[†] | 37.0[†] | [160] Tab. 1 |

[§]Literature reference and corresponding tabulated value, or figure used to extract the values listed here.

**Values listed for $\Delta H_f^0$ and $\Delta S_f^0$ are the unfolding enthalpy and entropy at $T_f$, 0M cosolute concentration, and pH indicated in the corresponding column. Throughout this table, values for enthalpy, entropy, and heat capacity are given in kJ/mol, kJ/mol/K, and kJ/mol/K respectively.

[a]Fit to Equations 2.4 and 2.9 [b]Fit to Equation 2.4. [c]Fit to Equation 2.10. [d] Used $\Delta H(T_f) = T_f\Delta S(T_f)$ to obtain unfolding entropy N/A: $\Delta C_p$ data not available. [†]$\Delta C_p$ measured for one concentration and assumed constant with respect to concentration in these references. The values for $\delta\Delta H_{\text{lab}}$ and $T\delta\Delta S_{\text{lab}}$ for these systems are thus likely inaccurate.

[α]α-chymotrypsinogen. [β]Histidine-containing phosphocarrier protein. [γ] β-hydroxyectoine. [δ] N-N' Dimethylurea. [ε]Trypsin Inhibitor.

[Cal] Measurements taken with calorimetry; [UV] Measurements taken with UV absorption;

[*] Monte-Carlo error analysis is performed on this protein in Figures 2.4-2.6.

## 2.3　Concluding Remarks

In this chapter we have observed and analyzed significant entropy-enthalpy compensation across both diverse proteins and diverse cosolute solutions, by performing a rigorous thermodynamic analyis of calorimetric and spectroscopic data, which included Monte-Carlo error estimates and a comparison across different models of the temperature-dependence of the unfolding heat capacity. Uncertainties in enthalpy and entropy, while much larger than the uncertainty in free energetic stability, do not rule out significant entropy-enthalpy compensation as a general phenomenon in protein transfer. The accuracy of the temperature-dependence and concentration-dependence of the unfolding heat capacity is not important near the folding transition, but is important if we are interested for example in the stability at lab temperature.

We can consider several possible scenarios for stabilizing and destabilizing cosolutes. A stabilizing cosolute, for example, could act in two different ways: it could increase the change in enthalpy upon unfolding, stabilizing the protein, while increasing the entropy of unfolding by a lesser amount. Or it could decrease the change in entropy upon unfolding, while decreasing the enthalpy by a lesser amount. We refer to the former case as an enthalpically stabilizing cosolute, and the latter case as an entropically stabilizing cosolute. Conversely a cosolute could decrease the change in enthalpy upon unfolding, destabilizing the protein, while decreasing the entropy of unfolding by a lesser amount; we refer to this as an enthalpically destabilizing cosolute. An entropically destabilizing cosolute would then be one that increases the entropy of unfolding while increasing the enthalpy of unfolding by a lesser amount. A cosolute could also in principle be stabilizing in both enthalpy and entropy (or destabilizing in both), which we refer to as uncompensated stabilization (or destabilization). These regions are illustrated in figure 2.8.

Each of these regions could in principle be populated, but examining the data in figures 2.4 and 2.5 shows an interesting pattern. In the systems we looked at, the majority of cosolutes were enthalpically stabilizing

Figure 2.4: **Entropy-enthalpy compensation for protein unfolding** transfer enthalpy $\delta\Delta H_f$ and unfolding transfer entropy $\delta\Delta S_f$, both evaluated at the folding midpoint and suitably normalized as described below. The legend, listed from the upper right data point to the lower left data point, indicates the protein, cosolute, pH, and corresponding source of the experimental data. Cosolutes above and to the left of the diagonal are destabilizing as noted; cosolutes below and to the right of the diagonal are stabilizing. Abscissa/ordinate are the transfer enthalpy/entropy normalized by the unfolding enthalpy/entropy in the absence of solute, per 100g/L of cosolute, i.e. $\delta\Delta H_f/(\Delta H_f \cdot c)$ *vs* $\delta\Delta S_f/(\Delta S_f \cdot c)$. Also plotted here are the Monte Carlo-generated scatter points for Arc Repressor in KCl (red circle), Notch Ankyrin in urea (green circle), and Hisactophilin in urea (blue circle). Bars on each of these three points show the standard deviation in the direction of the scatter. The scatter here does not substantially reduce the significance of the linear compensating trend. See also Table 2.3, which gives thermodynamic parameters for the proteins we study here.

Figure 2.5: **Further illustration of entropy-enthalpy compensation** for various proteins and solvents. The notation here is the same as in Figure 2.4, but the scale of the plot is significantly smaller. Scatter as a result of uncertainty for Creatine Kinase in glycerol (blue circle) is shown, along with bars to indicate the standard deviation. Scatter was calculated for RNase A in urea (cyan circle) and $\alpha$-lac in ethanol (mustard circle), but the scatter is smaller than the data point appearing on this plot. Black open symbols correspond to simulation data using the Tanford transfer model taken from O'Brien *et. al.* (ref. [144]). Legend labels are ordered from upper right data point to lower left data point.

Figure 2.6: **Entropy-enthalpy compensation for the transfer of various proteins to various solvents** is also seen by plotting $\delta\Delta H$ *vs.* $T\delta\Delta S$ at lab temperature ($25°$ C). The points cluster close to the $\delta\Delta H = T\delta\Delta S$ line; the deviation from that line (horizontal or vertical) represents the absolute change in stability upon transfer at $25°$ C. Points above the line correspond to destabilizing cosolutes, points below the line correspond to stabilizing solutes. Scatter points representing the range of uncertainty in obtaining the enthalpy and entropy are also shown, as determined by the Monte Carlo method described in Section 2.2.1. The scatter is highly correlated with a magnitude that in some cases is large enough to change the sign of $\delta\Delta H$ and $T\delta\Delta S$. The compensation is statistically significant however—see Section 2.2.2. In the case of $\alpha$-Lactalbumin in ethanol, the scatter was smaller than the symbol. The cyan circle with black outline indicates the Hisactophilin data assuming $\Delta C_p$ is independent of the concentration of urea, but has non-linear temperature-dependence obtained by fitting to Equation (2.4); the cyan circle with black square outline is the value obtained assuming $\Delta C_p$ is independent of temperature, but still accounting for the concentration-dependence. These approximations both introduce significant error: $\approx 80$ kJ/mol and 30 kJ/mol respectively. Similarly, assuming a concentration-independent unfolding heat capacity introduces an error of $\approx 40$ kJ/mol for RNase A in urea (circled blue cross) and a temperature-independent $\Delta C_p$ introduces an error of $\approx 10$ kJ/mol. Legend labels are ordered from upper right data point to lower left data point.

53

Figure 2.7: **A) Concentration-dependence of the heat capacity** for several protein-cosolute systems. For RNase A in urea, hisactophilin in urea, and $\alpha$-Lactalbumin in ethanol, error bars were determined from the Monte Carlo method described in Section 2.2.1. Error bars are not present for acylphosphatase or barstar because the corresponding literature data were not available for application of the Monte-Carlo method. The x-axis was normalized to facilitate comparison across proteins. B) heat capacity vs concentration for the T-independent and linear T-dependent models of $\Delta C_p$ for RNase A in urea.

or enthalpically destabilizing. Few were entropically stabilizing, and none entropically destabilizing. None could be confidently assigned as uncompensated stabilizers or destabilizers based on the error bars shown. The picture that emerges here is one in which enthalpy plays the dominant role in stabilization or destabilization, with entropy partially compensating the effect of enthalpy. This is consistent with recent results by Senske *et al.*[161] which found that even some macromolecular crowders, which are typically assumed to act primarily through entropy, and in fact enthalpic stabilizers.

Early results by Ben-Naim [162], Grunwald [163], Karplus [164], and Lee [165] have analyzed the invariable entropy-enthalpy compensation that occurs during solvent reorganization around a solute due to solvent-solvent interactions. In these theories, cavity creation results in a singular solute-solvent potential and is non-compensating, the limiting case being the free energy of inserting a non-interacting, hard-sphere solute. This issue is un-

Figure 2.8: **A diagram illustrating the possible categories of cosolute** discussed in section 2.3

likely to be a factor in the transfer scheme wherein a solute (protein) is transferred from pure buffer to solution containing cosolute: volume is indeed lost to buffer and cosolute upon transfer at constant pressure, but is also gained to the pure buffer system. A systematic analysis of entropy-enthalpy compensation in protein transfer using density functional theory to capture the effects of solvation is an interesting topic for future work.[166, 167]

# Chapter 3

# Classical Density Functional Theory and Protein-Cosolute Interactions

The problem we consider in this chapter is that of calculating the free energy change upon moving a solute such as a protein from a pure water environment and inserting it into a water and cosolute environment. Figure 3.1 illustrates the problem we are considering in the context of the Tanford transfer model for protein folding [168]. The cycle depicted here implies that if we know properties (such as the unfolding free energy) of the protein of interest in water, and we know the transfer free energy of each state of the protein from water to water and cosolute, we can find the corresponding properties of the protein in water and cosolute. This can be implemented either as an implicit solvent model, which we will discuss in Section 3.4, or in a post-processing way as in Ref. [145], which we will use to empirically fit our DFT model in Section 3.3.

The organization of this chapter is as follows. We begin in Section 3.1 by investigating the expected behavior of the surface and volume contributions to the transfer free energy in a heuristic model. In Section 3.2 we derive the principal equations for the DFT model of the transfer free energy. In Section 3.2.2 - Section 3.4, we consider several examples of how the DFT model can

Figure 3.1: **A diagram of the Tanford transfer model**, for a transfer process going from a pure water environment to one of water and cosolutes. Knowledge of the free energy of unfolding $\Delta G_{\text{wat}}^{u \to f}$ in the absence of cosolutes can be combined with the transfer free energies of the folded $(\Delta G_{w \to o}^{\text{fold}})$ and unfolded $(\Delta G_{v \to s}^{\text{un}})$ states to obtain the free energy of unfolding in the presence of cosolutes $\Delta G_{\text{cos}}^{u \to f}$.

be applied, making connections with the model developed in Section 3.1. We finally conclude and give our outlook on future directions for this approach.

## 3.1 Volume and Area Terms in the Transfer Free Energy

### 3.1.1 Volume Considerations

To appreciate the terms that we expect in an expression for the transfer free energy, we initially consider both volume and surface area effects in a more qualitative way. We consider the difference in volume occupied by the folded and unfolded states, or more precisely the expanded and collapsed states of a polymer, to obtain the corresponding free energy difference in the

presence of a bath of "hard-sphere" cosolutes. There are thus no surface interactions to consider, and we seek to estimate the magnitude of the volume effect; we also ignore for the time being the change in internal free energy as the polymer collapses. The free energy change upon collapse of a protein or polymer then arises from the change in entropy of the cosolutes, due to the change in available phase space. For hard-sphere cosolutes, the volume occupied by the expanded polymer will be larger than that of the collapsed polymer. The same considerations apply to a collapsed *vs.* expanded protein; unfolded states of proteins are generally found to be expanded relative to the folded state [169]. In what follows, let $r_a$ be the mean amino acid radius, $r_o$ the cosolute radius, and $N_p$ the number of amino acids in the polymer or protein. Treating the unfolded protein crudely as a meandering cylindrical tube (see Figure 3.2a inset), the volume is approximately $\pi(r_a + r_o)^2(2r_a N_p + 2r_o)$, which is that of a cylinder of radius $r_a + r_o$ and length $2N_p r_a + 2r_o$. The volume of the collapsed globule, or folded protein, can be modelled as a sphere of radius $R_p + r_o$, where $R_p$ is the protein radius as probed by a zero-radius cosolute particle, i.e. the collapsed volume is $(4/3)\pi(R_p + r_o)^3$. When $r_o = 0$, the unfolded and folded volumes must be equal, giving $R_p^3 = (3/2)N_p r_a^3$. The change in available volume for cosolutes $\Delta V(r_o)$ upon polymer collapse is thus positive, and is plotted in Figure 3.2 as a function of cosolute radius $r_o$, for a chain of length $N_p = 70$.

We can compare the results of the above simple model to data taken from simulations of a $C_\alpha$ Gō model of cold-shock protein (PDB 2L15), with 70 amino acids, generated with the GROMACS molecular dynamics package. The Gō potential was generated using a shadow map for the native contacts [170] by the SMOG@ctbp server [171]. The simulated free energy surface has a double-well structure with well-defined folded ($f$) and unfolded ($u$) ensemble as observed in $C_\alpha$ Gō models for other single domain proteins [172]. We take conformational snapshots in each ensemble and measure the volume using a variable probe radius with the program VOIDOO [173]. The average volume change $\Delta V = \langle V_u \rangle - \langle V_f \rangle$ for a given probe radius is plotted in Figure 3.2a. The theory and simulation data compare quite well given the simplicity of the model.

Figure 3.2: a) **The change in volume upon collapse** $\Delta V(r_o) = V_u - V_f$, as a function of cosolute radius $r_o$, for a polymer chain of length $N_p = 70$ residues and with $r_a = 6$ Å. The magnitude of the change in volume monotonically increases as $r_o$ increases. Also plotted are the average $\Delta V = \langle V_u \rangle - \langle V_f \rangle$ values of simulation trajectories of Cold-Shock Protein ($N = 70$, PDB 2L15) against probe radius. (Inset) Schematic of collapsed/folded and unfolded polymer. Folded polymer has radius $R_p$; unfolded polymer has tube radius $r_a$ and length $N_p r_a$. b) Minus the change in free energy upon collapse as a function of cosolute radius $r_o$, for both constant packing fraction $\eta$ and constant concentration $\rho$. The value of $\rho$ was set to $1M$, and the value of $\eta$ was set so that the free energy change would be equal to that at constant $\rho$ at a typical cosolute radius of 3.1 Å. This gave a packing fraction $\eta \approx 0.075$.

59

We now consider the free energy as a function of either uniform density $\rho$ or packing fraction $\eta$ of the cosolutes. Given a large effective box with volume $V_{box}$ containing a given protein, the packing fraction of cosolutes $\eta$ (i.e. the volume density) is given by

$$\eta = \frac{\frac{4}{3}\pi r_o^3 N_o}{V_{box} - V_{prot}(r_o)} \approx \frac{\frac{4}{3}\pi r_o^3 N_o}{V_{box}} = \frac{4}{3}\pi r_o^3 \cdot \rho \,,$$

where $\rho$ is the number density. So, at a fixed packing fraction the number of cosolutes $N_o$ scales as $r_o^{-3}$.

To estimate the volume contributions to the free energy change upon collapse, $\Delta G_V(r_o)$, as a function of cosolute radius but at either fixed density or packing fraction, we use the ideal gas approximation for the osmotic pressure $p_{osm} = \rho k_B T$ to obtain

$$\Delta G_V(r_o) = p_{osm}\Delta V(r_o) = \rho k_B T \Delta V(r_o) = \frac{\eta k_B T \Delta V(r_o)}{\frac{4}{3}\pi r_o^3} \qquad (3.1)$$

where $\Delta V(r_o)$ is obtained from the model above.

A plot of the magnitude of the free energy change upon collapse as a function of cosolute radius, here exclusively due to the increase in entropy of cosolute particles, is shown in Figure 3.2b. Based on these considerations we can estimate the volume-like contribution for typical cosolute sizes and concentrations. Taking TMAO (Trimethylamine N-oxide) as an example, we expect the cosolute radius to be about 2 Å, from the water oxygen-TMAO nitrogen radial distribution function[174]. Given this radius and a concentration of 300 g/L, for a protein of length $N_p = 70$ we estimate a volume contribution to the free energy of $\approx 4k_B T$. The free energy of unfolding is linear in protein length, so a larger protein of $N_p = 300$ has an estimated $\Delta G \approx 17k_B T$.

### 3.1.2 Surface Considerations

The presence of cosolutes in solution can make the effective solvent more repulsive to protein resulting in stabilization, or more attractive to the protein

resulting in denaturation. What effect is observed depends on the energy $\epsilon$ of cosolute-protein binding and also the concentration $c$ (or equivalently the chemical potential $\mu$) of the cosolute.

The energy $\epsilon$ of binding of the cosolute is actually the difference in internal free energy of binding between cosolute and water, since for example water may have some attraction to the polymer, and also a cosolute may supplant more than one water molecule in the process of binding.

Previous treatments of transfer free energy analysis as a condensation problem onto the surface of the protein have been undertaken primarily in the context of protein denaturation and the prediction of $m$-values [175–177]. The process of condensation of a cosolute to a surface is equivalent to the well-known statistical mechanical problem of Langmuir's isotherm [178], for which the partition function $\mathfrak{Z}$ in the $(T, \mu)$ ensemble for a substrate with $M$ absorbing sites is given by $\left(1 + \mathrm{e}^{-\beta(\epsilon-\mu)}\right)^{M}$. The mean covering ratio $f$ is then given by

$$f = \frac{kT}{M} \frac{\partial \log \mathfrak{Z}}{\partial \mu} = \frac{1}{1 + \mathrm{e}^{\beta(\epsilon-\mu)}} \,, \tag{3.2}$$

and the mean energy of condensation on the surface is $Mf\epsilon$. Here we neglect interactions between cosolutes when bound. The Helmholtz free energy in this model is given by

$$G = -pV + fM\mu = -k_BT \log(\mathfrak{Z}) + fM\mu$$

with $T, \mu$ partition function $\mathfrak{Z}$ as given above.

We can relate the Langmuir isotherm to the free energy of a protein surface by assuming that each cosolute occupies an area $a_0 \approx \pi r_o^2$ on the protein surface, so that we can write $M = A/a_0$, where $A$ is the protein's solvent accessible surface area in a given conformation. The change in free energy $G_A$ upon condensation becomes

$$G_A = -k_BT \frac{A}{a_0} \log\left(1 + \mathrm{e}^{-\beta(\epsilon-\mu)}\right) + f\frac{A}{a_o}\mu \tag{3.3}$$

If the concentration of unbound cosolute is dilute, an ideal gas approximation suffices for the chemical potential: $\mu = kT \log\left(\rho/\rho_Q\right)$, where $\rho_Q$ is a

reference concentration (typically taken to be 1M). The quantity $e^{-\beta\epsilon}/\rho_Q$ is typically treated as an equilibrium constant in the literature [176, 177]. We consider both dilute and non-dilute limits below. The protein's exposed surface area is obtained from the volume given in Section 3.1.1 by $A = \partial V/\partial r_o$, so the collapsed exposed area is $4\pi (R_p + r_o)^2$ and the expanded (random coil) exposed area is $2\pi (r_o + r_a) [(2N_p + 1) r_a + 3r_o]$.

### 3.1.3  Combined Surface/Volume Model for the Transfer Free Energy

We can now write the total free energy of collapse $\Delta G$ arising from cosolutes by combining the volume and surface area terms in equations (3.1) and (3.3). We can also remove the ideal gas assumption by expressing $\Delta G$ in terms of the Carnahan-Starling (CS) approximations to the pressure and chemical potential: [179]

$$
\begin{aligned}
p &= \rho k_B T \frac{1 + \eta + \eta^2 - \eta^3}{(1 - \eta)^3} \\
\mu &= k_B T \log(\rho/\rho_Q) + k_B T \frac{8\eta - 9\eta^2 + 3\eta^3}{(1 - \eta)^3} \, .
\end{aligned}
\tag{3.4}
$$

where $\eta$ is the volume fraction (the volume per molecule times $\rho$). Then the free energy becomes:

$$
\Delta G = p\Delta V + \left( \frac{k_B T}{\pi r_o^2} \log (1 - f) + \frac{f\mu}{\pi r_o^2} \right) \Delta A
\tag{3.5}
$$

with $f$ given in (3.2) and $p$ and $\mu$ given in (3.4), and where

$$
\Delta V(r_o) = \frac{4}{3}\pi \left( \left( \frac{3N_p}{2} \right)^{1/3} r_a + r_o \right)^3 - 2\pi(r_a + r_o)^2(N_p r_a + r_o)
$$

$$
\Delta A(r_o) = 4\pi \left( \left( \frac{3N_p}{2} \right)^{1/3} r_a + r_o \right)^2 - 2\pi(r_a + r_o)[(2N_p + 1)r_a + 3r_o]
$$

are the volume and surface area change upon folding (or collapse).

We plot equation (3.5) in Figure 3.3 as a function of cosolute radius $r_o$,

for condensation energies $\epsilon = 2k_BT$ and $\epsilon = -k_BT$. To assess the limits of the ideal gas model, we have also plotted the ideal gas results in Figure 3.3. For repulsive cosolute-protein interactions, both surface and volume terms stabilize the folded or collapsed state (Figure 3.3). The free energy change upon collapse is monotonically decreasing (increasing in magnitude) from zero, and more strongly favoring collapse as cosolute radius is increased. Non-ideal excluded volume effects in the cosolute pressure and chemical potential enhance the stabilizing effect. For attractive cosolute-protein interactions, the situation is more complex. At small values of cosolute radius $r_o$, the collapsed phase is destabilized by cosolute-protein binding, which favors expansion. As $r_o$ increases, the volume change upon collapse increases, which begins to entropically favor collapse. The osmotic pressure initially increases modestly, additionally favoring collapse. However the chemical potential also increases modestly, driving condensation of cosolute and favoring expansion. These two effects nearly cancel each other rendering the real and ideal gas curves nearly coincident up to $r_o \approx 4\text{Å}$. The sigmoidal dependence of covering fraction $f$ in equation (3.2) on chemical potential $\mu$ results in a sudden condensation of cosolute onto the protein around $r_o \approx 5\text{Å}$, which induces the system to favor expansion at these radii. While the number of condensed cosolutes is bounded, the osmotic pressure is not, and eventually collapse is favored once again through volume terms. The cosolute radius $r_o$ can only increase until $\eta \approx 0.6$ (near crystal packing densities), giving a cutoff of $r_o^{(cut)} \approx (3\eta/4\pi\rho)^{1/3}$, or about 6.2Å for 1M concentration.

In the limit that the cosolute is dilute, $\rho e^{-\beta\epsilon}/\rho_Q \ll 1$ and we can expand the logarithm in equation (3.5) to obtain an area contribution to the free energy of $-\rho k_B T A e^{-\beta\epsilon}/a_0\rho_Q$, so that the free energy change upon unfolding becomes

$$\Delta G = \rho k_B T \left( \Delta V - At e^{-\beta\epsilon} \right) . \tag{3.6}$$

Here we have used the fact that $(a_0\rho_Q)^{-1}$ has units of length and can be thus be interpreted physically as a thickness $t$ over which the surface interaction acts.

Having looked at these preliminary volume and surface considerations,

Figure 3.3: **Total free energy change $\Delta G$ upon collapse** in units of $k_B T$, as a function of cosolute radius $r_o$. Values of packing fraction $\eta$ corresponding to the values of $r_o$ on the x-axis are shown above the plot. Curves are taken from equation (3.5) which combines surface area and volume terms. Here the polymer length $N_p = 70$, the cosolute concentration $\rho = 1$M, and $r_a = 6$Å. Red curves show $\Delta G$ upon collapse for a repulsive cosolute with interaction energy $+2k_B T$, i.e. a crowding particle. Blue curves show $\Delta G$ upon collapse for an attractive cosolute with interaction energy $-k_B T$, i.e. a weak denaturant. Plotted are both the model with ideal gas (IG, dashed) and Carnahan-Starling (C-S, solid) pressure and chemical potential.

we now turn to a classical density functional theory formulation, which provides a more complete understanding of the transfer free energy, and as well, reduces to equation (3.6) in the appropriate limits.

## 3.2 The Density Functional Theory Formulation

We now consider a density functional formulation of the problem of transfer free energy. In what follows, we will assume that the intra-protein energy of a given configuration of a protein is in principle known and the net interac-

tion between any given site on the protein and either the cosolute or water is in principle known. We then wish to calculate $\Delta G$, the free energy of transferring the protein from water to an cosolute solution, or, equivalently, of transferring the cosolutes from an aqueous solution to one containing the protein (see Figure 3.1). In short, we wish to consider the effect that the presence of cosolutes has on the free energy of the protein.

The uniqueness of the Kohn-Sham density functional may be extended to finite temperatures, so that the free energy of the protein-solvent system is uniquely expressed as a functional of the single particle density $\phi(\boldsymbol{r})$ [180]. We thus seek an expression for the free energy of the cosolutes and water in an arbitrary external potential. For our purposes in obtaining a transfer free energy, we will treat a given protein configuration, with atom positions $\{\boldsymbol{R}_i\}$, as the source of the external potential. We write the free energy in the standard way [80]:

$$
\begin{aligned}
G(\{\boldsymbol{R}_i\}) = \int d^3r \; k_B T(-S_o(\phi_o(\boldsymbol{r})) - S_w(\phi_w(\boldsymbol{r}))) + \mathcal{V}_o(\boldsymbol{r})\phi_o(\boldsymbol{r}) + \mathcal{V}_w(\boldsymbol{r})\phi_w(\boldsymbol{r}) \\
+ \Phi_o[\phi_o] + \Phi_w[\phi_w] + \Phi_{ow}[\phi_o, \phi_w] \quad\quad (3.7)
\end{aligned}
$$

Here $\phi_j$ is the density function for the cosolutes ($o$) or water ($w$), and $\mathcal{V}_j$ the external potential on the respective species. The entropy density for each species can be written as

$$
S_o(\boldsymbol{r}) + S_w(\boldsymbol{r}) = -\phi_o(\boldsymbol{r}) \log \left[\lambda_o^3 \phi_o(\boldsymbol{r})\right] - \phi_w(\boldsymbol{r}) \log \left[\lambda_w^3 \phi_w(\boldsymbol{r})\right] \quad\quad (3.8)
$$

where $\lambda_o$ and $\lambda_w$ are constants with units of length, analogous to thermal wavelengths. The terms $\Phi_o$, $\Phi_w$, and $\Phi_{ow}$ are the multi-particle correlation contributions to the free energy for the respective species. For example, the two particle correlation part of $\Phi_o$ would have the form

$$
\Phi_o^{(2)}[\phi_o] = \int \int d^3r_1 d^3r_2 \; \phi_o(\boldsymbol{r}_1)\phi_o(\boldsymbol{r}_2)U_{oo}(\boldsymbol{r}_1 - \boldsymbol{r}_2)g(\boldsymbol{r}_1, \boldsymbol{r}_2|\mathcal{V}) \quad\quad (3.9)
$$

where $U_{oo}$ is the interaction potential between two cosolutes and $g$ the two-particle correlation function. The full multi-particle function is not known

65

exactly, and so, as in electronic DFT, while equation (3.7) is exact in principle, approximations must be made to use it in practice [181].

We now make two key assumptions. The first is that the water and cosolute densities are completely correlated, such that all vacua are occupied by either water or cosolute. Thus $N_w v_w + N_o v_o = V$, where $v_i$ is the volume of an individual water or cosolute molecule, and $V$ the total volume. Dividing this by $V v_w$ and allowing the local density of a given species to vary gives

$$\phi_w(\boldsymbol{r}) + f \phi_o(\boldsymbol{r}) = \rho_w \tag{3.10}$$

where $f = v_o/v_w$ and $\rho_w = 1/v_w$ (the factor of $f$ allows for the cosolute molecule to be a different size than the water molecule). Equation (3.10) is not valid in the interior of the protein, so we split our system up into two regions: a hard wall region $V_{hw}$ in which $\phi_w = \phi_o = 0$, and the rest of the system, which has a volume $V$ identical to the volume of the cosolute-water bath prior to the insertion of the protein, and in which Equation (3.10) is valid. We further take $V_{hw}$ to be the same as the change in volume of the aqueous system the protein was removed from in the transfer process (see Figure 3.1), so that the total system of water, protein, and cosolute-water solution does not change volume during the transfer process.

With the approximation of equation (3.10) we can write

$$\mathcal{V}_o(\boldsymbol{r})\phi_o(\boldsymbol{r}) + \mathcal{V}_w(\boldsymbol{r})\phi_w(\boldsymbol{r}) = \Delta\mathcal{V}(\boldsymbol{r})\phi_o(\boldsymbol{r}) + \mathcal{V}_w(\boldsymbol{r})\rho_w \tag{3.11}$$

$$\Phi_o[\phi_o] + \Phi_w[\phi_w] + \Phi_{ow}[\phi_o, \phi_w] = \Phi_t[\phi_o] \tag{3.12}$$

where $\Delta\mathcal{V}(\boldsymbol{r}) = \mathcal{V}_o(\boldsymbol{r}) - f\mathcal{V}_w(\boldsymbol{r})$.

The second approximation in our treatment is that the cosolute number density is much less than that of water. Using this approximation along with the one given in Equation (3.10), the entropy in Equation (3.8) becomes

$$-S_o(\boldsymbol{r}) - S_w(\boldsymbol{r}) = \phi_o(\boldsymbol{r}) \log \left[\lambda_o^3 \phi_o(\boldsymbol{r})\right] + (\rho_w - f\phi_o(\boldsymbol{r})) \log \left[\lambda_w^3 (\rho_w - f\phi_o(\boldsymbol{r}))\right] \tag{3.13}$$

$$\approx \phi_o(\boldsymbol{r}) \log \left[\lambda_o^3 \phi_o(\boldsymbol{r})\right] - f\phi_o(\boldsymbol{r}) + \rho_w \log \left[\lambda_w^3 \rho_w\right] - f\phi_o(\boldsymbol{r}) \log \left[\lambda_w^3 \rho_w\right]$$

In this way we express each part of equation (3.7) in terms of cosolute density and constant terms. The free energy functional may then be written as

$$G = \int d^3r \, k_B T(\phi_o(\boldsymbol{r}) \log[\lambda_o \phi_o(\boldsymbol{r})] - (\gamma + 1)\phi_o(\boldsymbol{r})) + \Delta \mathcal{V}(\boldsymbol{r})\phi_o(\boldsymbol{r})$$
$$+ V\rho_w \log \lambda_w^3 \rho_w + \mathcal{U}\rho_w + \Phi_t[\phi_o] \tag{3.14}$$

where $\mathcal{U} \equiv \int d^3r \mathcal{V}(\boldsymbol{r})$, and $\gamma + 1 \equiv f(1 + \log(\lambda_w^3 \rho_w))$. Since $V$ is the volume of the system, the term $V\rho_w$ is equal to $V/v_w = N'_w$, the total number of water molecules in a system of pure water of volume $V$.

Thus, dropping the subscripts, letting $\mathcal{V} \equiv \Delta\mathcal{V}$, and ignoring any position independent terms, we can write the free energy as

$$G = \int d^3r \, k_B T \left( \phi(\boldsymbol{r}) \log \lambda^3 \phi(\boldsymbol{r}) - \phi(\boldsymbol{r}) \right) + k_B T \gamma \phi(\boldsymbol{r}) + \mathcal{V}(\boldsymbol{r})\phi(\boldsymbol{r})$$
$$+ \Phi[\phi] \tag{3.15}$$

where $\Phi[\phi]$ is the functional containing the multi-particle correlation part of the free energy, and $\lambda \equiv \lambda_o$ is a constant with units of length analogous to the thermal wavelength, whose value will be shown to be unimportant. For now we will formally manipulate $\Phi$ without making assumptions about its form. We can find the density that minimizes the free energy by use of the Euler-Lagrange equations, with the constraint that the cosolute density when integrated over the total volume is the total number of cosolutes:

$$\int_V d^3r \, \phi(\boldsymbol{r}) = N_o \,. \tag{3.16}$$

We thus write

$$\frac{\delta}{\delta\phi} \left[ G - \mu_o \left( \int_V d^3r \, \phi(\boldsymbol{r}) - N_o \right) \right] = 0$$
$$\text{or} \quad k_B T \log \lambda^3 \phi(\boldsymbol{r}) + \mathcal{V}(\boldsymbol{r}) - k_B T\gamma + \frac{\delta\Phi}{\delta\phi} - \mu_o = 0 \tag{3.17}$$

where $\mu_o$ is the Lagrange multiplier corresponding to the constraint in equation (3.16). Physically, we can interpret equation (3.17) as a statement that

$\frac{\delta G}{\delta \phi}$ is equal to the chemical potential $\mu_o$, and thus must be a constant value at all points in space. Solving this for the density field gives

$$\phi(\boldsymbol{r}) = e^{\gamma} \lambda^{-3} e^{-\beta(\mathcal{V}(\boldsymbol{r}) + \Phi' - \mu_o)} \tag{3.18}$$

where $\Phi' \equiv \frac{\delta \Phi}{\delta \phi}$.

To obtain $\mu_o$ from equation (3.18), we use the constraint on the total number of particles in equation (3.16) which yields

$$e^{\beta \mu_o} = \frac{e^{\gamma} \lambda^3 N_o}{\int_V d^3 r \, e^{-\beta(\mathcal{V}(\boldsymbol{r}) + \Phi')}} \,. \tag{3.19}$$

From here we can obtain the transfer free energy, which is given by the free energy of the cosolute bath in the presence of the external protein potential, $\mathcal{V}(\boldsymbol{r})$, minus the free energy of the cosolute bath without the protein potential $(\mathcal{V}(\boldsymbol{r}) = 0)$. We thus have

$$
\begin{aligned}
\Delta G =& \Delta \mu_o N_o \\
=& - k_B T N_o \log \left( \frac{e^{\gamma} \lambda^{-3}}{N_o} \int_V d^3 r \, e^{-\beta(\mathcal{V}(\boldsymbol{r}) + \Phi'_f(\boldsymbol{r}))} \right) \\
& + k_B T N_o \log \left( \frac{e^{\gamma} \lambda^{-3}}{N_o} \int_V d^3 r \, e^{-\beta \Phi'_i} \right)
\end{aligned} \tag{3.20}
$$

where the volume $V$ integrated over is the volume outside of hard-wall volume of the protein, and is the same in the initial and final systems. The difference $\Delta G$ is independent of $\lambda$ and $\gamma$.

The bath in the initial state is homogeneous and isotropic, so $\Phi'_i$ in equation (3.20) is independent of position. Thus it may be factored out of the integral,

$$\int_V d^3 r \, e^{-\beta \Phi'_i} = V e^{-\beta \Phi'_i}$$

so that

$$\Delta G = -k_B T N_o \log \left( \frac{1}{V} \int_V d^3 r \, e^{-\beta(\mathcal{V}(\boldsymbol{r}) + \Delta \Phi')} \right) \tag{3.21}$$

where $\Delta \Phi' = \Phi'_f(\boldsymbol{r}) - \Phi'_i$. The expression in equation (3.21) consists of the

logarithm of the integral of a Boltzmann weight for the effective potential $\mathcal{V}(\boldsymbol{r}) + \Delta\Phi'(\boldsymbol{r})$. Here $\mathcal{V}(\boldsymbol{r})$ and $\Delta\Phi'(\boldsymbol{r})$ enter on equal footing. Recall that $\mathcal{V}$ is the protein-cosolute potential, treating the protein as an external source. $\Phi'$ is the functional derivative of the multi-particle part of the free energy. If we use the two-particle cosolute contribution from equation (3.9), we obtain

$$\Delta\Phi_o^{(2)\prime} = \left.\frac{\delta\Phi_o^{(2)}}{\delta\phi_o(\boldsymbol{r})}\right|_{\mathcal{V}} - \left.\frac{\delta\Phi_o^{(2)}}{\delta\phi_o(\boldsymbol{r})}\right|_{\mathcal{V}=0} = \int d^3r' \left[\phi_{of}(\boldsymbol{r}')g(\boldsymbol{r},\boldsymbol{r}'|\mathcal{V}) - \phi_{oi}(\boldsymbol{r}')g(\boldsymbol{r},\boldsymbol{r}'|\mathcal{V}=0)\right] U_{oo}(\boldsymbol{r},\boldsymbol{r}'),$$
$$(3.22)$$

which gives the difference of two terms in the presence and absence of the external protein potential, where each term corresponds to the equilibrium-averaged interaction energy between cosolutes, up to pair correlations. Thus the term $\Delta\Phi'$ in Equation (3.21) can be interpreted as the change in energy due to redistribution of the environment in response to the change in external potential.

We can recast equation (3.21) into a form that will be somewhat more useful later:

$$\Delta G = -k_B T N_o \log\left(1 + \frac{1}{V}\int_V d^3r \left[e^{-\beta(\mathcal{V}(\boldsymbol{r})+\Delta\Phi')} - 1\right]\right) \qquad (3.23)$$

which has the advantage that when $\mathcal{V}$ and $\Delta\Phi'$ are both zero, the integrand is also zero, and thus the integral can be taken over all space.

In equation (3.23) we can take the limit $V \to \infty$, with $N_o/V = \rho$ fixed. Then, assuming that the region over which the integrand in equation (3.23) is non-zero is finite, we can expand the logarithm to first order to obtain

$$\Delta G = -k_B T N_o \frac{1}{V}\int d^3r \left(e^{-\beta(\mathcal{V}(\boldsymbol{r})+\Delta\Phi')} - 1\right) \qquad (3.24)$$

which has the form

$$\Delta G = p_{id}\Delta V_{\text{eff}}$$

where $p_{id} = N_o k_B T/V$ is the ideal gas osmotic pressure, and $V_{\text{eff}} = \int d^3r \left[1 - e^{-\beta(\mathcal{V}(\boldsymbol{r})+\Delta\Phi')}\right]$ is an effective change in volume. In the dilute limit, the osmotic pressure $p = p_{id}$; then $V_{\text{eff}}$ may be interpreted as the change in volume available to

the cosolutes.

We now need to address $\Delta\Phi'$ to progress further. The obvious first approximation is to set $\Delta\Phi' = 0$; we will see below that this approximation can in fact go quite a long way, depending on the solvent. This is consistent with the observations in Figure 3.3 where the ideal gas approximation, which neglects cosolute-cosolute correlations, holds for typical molecular radii at 1M concentration. It is worth noting that this is not ignoring the cosolute-cosolute, cosolute-water, and water-water correlations completely; it is merely assuming that they are the same in the initial and final baths. Making this approximation, we have

$$\Delta G = -k_B T N_o \log\left(1 + \frac{1}{V}\int_V d^3 r \left[\mathrm{e}^{-\beta\mathcal{V}(\boldsymbol{r})} - 1\right]\right) \qquad (3.25)$$

Equation (3.25) represents an approximation to the transfer free energy that, while severe, nonetheless takes into account both the change in energy and change in entropy of the cosolute bath.

### 3.2.1 Validation Tests in Model Solvents

As a test of the density functional theory, we have used equation (3.25) to calculate the transfer free energy of several small molecules into model cosolutes. To simplify the simulations, we looked at transfer from vacuum to a van der Waals gas of cosolutes, which were taken to be single atoms interacting through a VDW potential. The density of the cosolutes was set to 1M. The molecules we transferred were the side chains of alanine and valine, with C-$\beta$ capped with a hydrogen to replace the backbone (*ie*, the molecules were methane and propane). The coordinates were taken from an existing protein structure file, and the angle and bond parameters were generated with the GROMACS utility pdb2gmx. The charges were set to zero for all atoms, and the interaction was purely van der Waals. We list the VDW parameters in Table 3.1. Figure 3.4 shows the interaction potential for the two different cosolutes we used. The transfer energies were calculated both with equation (3.25) and by simulating the transfer

Table 3.1: **van der Waals parameters for the atoms used** in the simulation test of the DFT, as taken from the CHARMM parameter set. Cos2 is a relatively attractive spherical cosolute, while the potential of Cos1 is dominated by steric repulsion. The interaction is parameterized as $V(r) = 4\epsilon\left[(\sigma/r)^{12} - (\sigma/r)^6\right]$.

| Atom | $\sigma$ (Å) | $\epsilon$ (kJ/mol) |
|---|---|---|
| Ala C-$\beta$ | 0.36705 | 0.33472 |
| Ala H | 0.23520 | 0.092048 |
| Val C-$\beta$ | 0.40536 | 0.08368 |
| Val C-$\gamma$ | 0.36705 | 0.33472 |
| Val H | 0.23520 | 0.092048 |
| Cos1 | 0.40536 | 0.08368 |
| Cos2 | 0.36705 | 0.33472 |

Table 3.2: **Comparison of test cases between density functional theory (DFT) and thermodynamic integration (TI)**

| Molecule/cosolute | DFT $\Delta G$ (kJ/mol) | TI $\Delta G$ (kJ/mol) |
|---|---|---|
| Ala/Cos1 | $0.188 \pm 0.002$ | $0.187 \pm 0.002$ |
| Val/Cos1 | $0.255 \pm 0.004$ | $0.261 \pm 0.004$ |
| Ala/Cos2 | $0.055 \pm 0.002$ | $0.059 \pm 0.003$ |
| Val/Cos2 | $-0.018 \pm 0.004$ | $-0.011 \pm 0.004$ |

in GROMACS and using Thermodynamic Integration (TI) [182–184]. The results are summarized in Table 3.2, and show excellent agreement between TI and DFT. This is notable since the result was obtained neglecting the inter-particle correlations, and at 1M the pressure of the cosolutes was $\approx 1.5$ that of the ideal gas pressure, which indicates that the cosolute-cosolute interactions were significant.

### 3.2.2 Connecting DFT to Previous Surface/Volume Models

We now take a simplified model of a protein potential to compare with the results obtained previously in Section 3.1 for the solvent contribution to the change in free energy upon protein collapse. In this model we will consider the protein to have an excluded volume of $V_{prot}$; that is, within

Figure 3.4: **Comparison of cosolute potential functions** for the test cases parameterized in Table 3.1. Cos2 is significantly more attractive than Cos1, which is reflected in the transfer free energies in Table 3.2

that volume the potential is infinite. From the discussion in Sections 3.1.1-3.2 concerning excluded volume, we saw that the changes in volume treated there are volumes from which cosolutes are excluded. We also consider the protein to have a surface region of thickness $t$ that exerts a potential on the cosolutes of depth $\epsilon$; this region is sufficiently thin that we can approximate its volume as $V_{surface} \approx tA$. If we use this model in the expression for the free energy in the limit of large system size (equation (3.25) ) then we obtain a free energy upon transfer of

$$\Delta G = \rho k_B T \left( V_{prot} + (1 - e^{-\beta \epsilon}) tA \right) . \tag{3.26}$$

The DFT transfer free energy with this simplified model provides a natural split between the volume contribution $p_{id}V_{prot}$ and the surface area contribution $p_{id}(1 - e^{-\beta\epsilon})tA$. Thus the DFT result, in the appropriate model, naturally generates the free energy contributions derived in Section 3.1 from more bespoke considerations. Specifically, if we take the total volume of the protein upon insertion to be $V = V_{prot} + tA$, then equation (3.26) is identical to equation (3.6). The simplified DFT model here reduces to our earlier considerations and helps give a physical interpretation of the quantity $\rho_Q$ as it pertains to the protein surface.

We can also see that, in order to obtain a SASA approximation in which $\Delta G$ is independent of temperature, one would have to assume that the cosolute-protein binding energy $\epsilon \ll k_B T$, and that volume terms were either negligible compared to surface terms, or they were proportional to them. We find below that $\epsilon \approx k_B T$ in order to obtain empirically-derived transfer free energies to TMAO, which does not satisfy the above inequality. As well, we can use the tube model from Section 3.1 for protein volume and surface area to estimate the relative contributions of volume and area: for a cosolute of radius $r_o = 2.5$ Å and a protein with $N_p = 70$, $V/tA = 0.62$ in the unfolded state, and $V/tA = 0.77$ in the folded state. The volume here is by no means negligible.

We thus expect on general grounds that the transfer free energy will be dependent on temperature. One way of looking at the simplified limit for the transfer free energy in equation (3.26) is as a derivation of a new phenomenological form for the transfer energy, containing both temperature and volume dependence:

$$\Delta G = \gamma_1 k_B T (V_{solute}) + \gamma_2 k_B T (\text{ASA}) e^{-\beta\epsilon} , \qquad (3.27)$$

where one can now fit the parameters $\gamma_1$, $\gamma_2$, and $\epsilon$, to empirical data.

## 3.3 Empirically Deriving DFT Transfer Free Energy Parameters

The potential $\mathcal{V}(\boldsymbol{r})$ in equation (3.25) is an effective potential given by $\mathcal{V}_o(\boldsymbol{r}) - f\mathcal{V}_w(\boldsymbol{r})$. Obtaining $f$ and $\mathcal{V}_w$ *ab initio* may be nontrivial, so we examine some model systems, and compare with empirical methods. To begin with, we will assume that the potential takes the form of a sum of terms from each particle in the protein, where a particle may be an atom in an all-atom model, or a bead modeling an amino acid in a coarse-grained approach:

$$\mathcal{V}(\boldsymbol{r}) = \sum_{i=1}^{N_p} v_i^{\text{eff}}(\boldsymbol{r} - \boldsymbol{R}_i) \,.$$

Here $N_p$ is the number of particles in the protein, and $\boldsymbol{R}_i$ the position of the $i$th particle.

We consider a model consisting of backbone $C_\alpha$ atoms and coarse-grained side-chain beads, which then form the particles for our potential. We make the assumption that the protein-cosolute potentials have a 6-12 form:

$$v_i(r) = 4\epsilon_i \left[ \left( \frac{\sigma_i}{r} \right)^{12} - \left( \frac{\sigma_i}{r} \right)^6 \right] \,,$$

and we wish to determine the potential parameters $\sigma_i$, $\epsilon_i$ for each amino acid that reproduce the transfer energies found experimentally when DFT is applied using the above potential. As a starting point, we examine those used by Auton and Bolen [145].

Two constraint equations are required for each amino acid. For the first equation, we note that the beads representing the various amino acid side chains have residue radii $r_{oi}$ that may be obtained from measured partial molar volumes [185]. We can then apply a constraint to the above 6-12 parameters $\sigma_i$, $\epsilon_i$ by requiring that at a distance $r_{oi}$ from the residue centre,

$$v_i(r_{oi}) = 0.6 \ \text{kcal} \cdot \text{mol}^{-1} \,. \tag{3.28}$$

To obtain the remaining equation determining the parameters $\sigma_i$, $\epsilon_i$, we

require that the DFT transfer free energy, as computed by the dilute limit of equation (3.25) for the single particle representing an amino acid side chain, should be equal to the experimental value as given in reference [145], specifically for transfer into a solution of 1M TMAO. This involves computing the integral over the cosolute-accessible volume in the expression

$$\rho k_B T \int d^3r \left(1 - e^{-\beta v_i(r)}\right) \tag{3.29}$$

and setting the result to the empirical value of $\delta g_i$ for each amino acid.

The sum of the transfer free energies of each amino acid in a Gly-X-Gly tripeptide is often used to approximate the conformationally-averaged transfer free energy for a protein.[145] Here we consider the tripeptide transfer free energies. The integral in expression (3.29) then involves integration over a solid angle $\Omega_i$ determined by the fraction of solid angle available to the side chain in the tripeptide *vs.* that for the isolated residue, i.e.

$$\Omega_i = \frac{A^i_{tri}}{A^i_{iso}} 4\pi$$

The potential $v_i$ is then fully determined from equation (3.28) along with

$$\Omega_i \rho k_B T \int_0^\infty dr \, r^2 \left(1 - e^{-\beta v_i(r)}\right) = \delta g_i . \tag{3.30}$$

We can now construct potentials for each amino acid transfer free energy given in reference [186]. The parameters derived from doing so are listed in Table 3.3. The backbone-cosolute interaction was parameterized as $v_{BB}(r) = C/r^{12}$, as this better represented its strongly repulsive character. The value of $C$ obtained by fitting to $\delta g_{BB}$ was $C = 7.510 \times 10^7$ kcal·Å$^{12}$.

In this context, the DFT formulation provides a way of using the information from tri-peptide experiments in a way that captures both energetic and entropic effects. The parameters just obtained can be used to determine the change in the transfer free energies for isolated residues as temperature changes. The experimental transfer free energies $\delta g_i$ are predicted to increase as temperature increases, with the new values at $T = 310$K given in

Table 3.3: **Parameter values yielding transfer free energies** $\delta g$ **to 1M TMAO for amino acid side chains and backbone at 300K, and the predicted** $\delta g$ **at 310K .**

| Type | $r_o$ (Å) [a] | $\delta g$ (cal/mol) [b] | $\sigma$ (Å) [c] | $\epsilon$ (kcal/mol) [d] | $\delta g(T = 310K)$ (cal/mol) [e] |
|------|------|------|------|------|------|
| Ala | 2.52 | -14.64 | 3.517 | 0.6286 | -12.65 |
| Arg | 3.28 | -109.3 | 4.088 | 1.022 | -104.0 |
| Asn | 2.74 | 55.69 | 4.564 | 0.0483 | 58.06 |
| Asp | 2.79 | -66.67 | 3.627 | 1.055 | -63.31 |
| Gln | 3.01 | 41.41 | 4.397 | 0.1710 | 44.57 |
| Glu | 2.96 | -83.25 | 3.799 | 0.9973 | -78.88 |
| His | 3.04 | 42.07 | 4.428 | 0.1707 | 45.28 |
| Ile | 3.09 | -25.43 | 4.084 | 0.5692 | -21.59 |
| Leu | 3.09 | 11.6 | 4.246 | 0.3405 | 15.15 |
| Lys | 3.18 | -110.23 | 3.968 | 1.126 | -104.7 |
| Met | 3.09 | -7.65 | 4.154 | 0.4538 | -3.791 |
| Phe | 3.18 | -9.32 | 4.237 | 0.4587 | -5.397 |
| Pro | 2.78 | -137.7 | 3.457 | 1.987 | -133.5 |
| Ser | 2.59 | -39.04 | 3.4905 | 0.8849 | -36.45 |
| Thr | 2.81 | 3.75 | 3.9312 | 0.3889 | 6.41 |
| Trp | 3.39 | -152.9 | 4.157 | 1.150 | -146.5 |
| Tyr | 3.23 | -114.3 | 4.020 | 1.103 | -109.2 |
| Val | 2.93 | -1.02 | 4.021 | 0.4238 | 1.78 |
| BB [f] | 2.25 | 90.0 | - | - | 92.7 |

[a]Distance where the cosolute-amino acid potential is taken to be 0.6 kcal·mol$^{-1}$
[b]Empirical transfer free energies to 1M TMAO
[c]van der Waals size parameter
[d]van der Waals well depth
[e] predicted transfer free energies at $T = 310$K
[f]Backbone is parameterized for TMAO by a purely repulsive potential (see text)

Table 3.3. Increasing temperature by $0.03k_B T$ increases the transfer free energy by $\approx 0.6k_B T$ for a 100 residue protein. This change is not large, but the relative temperature change is also small. The transfer entropy is significant: $d(\delta g)/dT \approx 20k_B$.

## 3.4   Using DFT for Implicit Solvent Models

The DFT methodology has been applied to the problem of solvation to calculate fluid correlation functions, solvation free energies, and reorganization energy in charge transfer [167, 187]. The use of time-dependent density functional theory has been well-established to understand solvation dynamics in single-component solvents [188] as well as selective solvation in binary mixtures [189, 190]. The methodology has also been applied to the connect static and dynamic approaches to the glass transition by Kirkpatrick and Wolynes [86]. The DFT methodology as described above may also be be applied to the problem of finding the effective forces for molecular dynamics simulation in an implicit solvent, which we briefly describe here.

We again write the external potential due to solute-solvent interactions as

$$\mathcal{V}(\boldsymbol{r}) = \sum_j v_j(|\boldsymbol{R}_j - \boldsymbol{r}|)$$

we can write the force on the $i$th particle from the transfer free energy in equation (3.24) (neglecting solvent inter-particle correlations) as

$$
\begin{aligned}
\mathbf{F}_i &= \nabla_{R_i} \left[ k_B T \rho \int d^3 r \, \left( 1 - \mathrm{e}^{-\beta \sum_j v_j(|\boldsymbol{R}_j - \boldsymbol{r}|)} \right) \right] \\
&= k_B T \rho \beta \int d^3 r \, \mathrm{e}^{-\beta \sum_j v_j(|\boldsymbol{R}_j - \boldsymbol{r}|)} \nabla_{R_i} v_i(|\boldsymbol{R}_i - \boldsymbol{r}|) \\
&= \rho \int d^3 r \, \mathrm{e}^{-\beta \sum_j v_j(|(\boldsymbol{R}_j - \boldsymbol{R}_i) - \boldsymbol{r}|)} \nabla v_i(r) \quad\quad (3.31)
\end{aligned}
$$

We immediately see that the integrand is non-zero only when $\nabla v_i(r)$ is non-zero, so that if there is an effective cutoff $r_c$ such that $v_i(r) \approx 0$ for $r > r_c$, then the integral in equation (3.31) only needs to be taken in the region $r < r_c$. This is a generalization of the result obtained by Götzelmann *et al*[191], who have shown that for a hard sphere potential, only the solvent density at the surface of the spheres was relevant to the calculation of depletion forces. Here we extend this analysis to arbitrary potentials.

Consider a particle with a spherically symmetric $v_i(r)$, as assumed above.

The net force on this particle when isolated is zero. When a second particle exerting potential $v_j(r)$ on the cosolutes is brought near, the net force on the first due to the solvent is a result of the now asymmetric solvent density. We note here we are treating the indirect force rather than the direct force between the particles, which can be calculated by direct application of the interparticle potential. The region of asymmetric solvent density constitutes a restricted volume to be integrated over in equation (3.31), as only the region of overlap between the two spheres defined by the cutoff in potential around $\boldsymbol{R}_i$ and $\boldsymbol{R}_j$ contributes to the net force (see e.g. Figure 3.5b below). In addition, the solvent field in this overlap region will maintain cylindrical symmetry about the axis joining the two particles, which means that the force will be along this axis as well. This suggests that the force on particle $i$ can be written as

$$\boldsymbol{F}_i = \sum_{|\boldsymbol{R}_{ij}|<2r_c} F_{ij}(|\boldsymbol{R}_{ij}|)\hat{\boldsymbol{R}}_{ij} .$$

Here $\hat{\boldsymbol{R}}_{ij}$ is the unit vector from particle $j$ to particle $i$, and $F_{ij}$ is a scalar function of the interparticle distance $|\boldsymbol{R}_{ij}| \equiv |\boldsymbol{R}_i - \boldsymbol{R}_j|$, which is determined by the overlap integral in equation (3.31), and which could in principle be pre-computed and tabulated to speed up execution.

### 3.4.1 Depletion and Impeded-Solvation Interactions in an Implicit Solvent Model

We can use equation (3.31) to investigate the forces due to solvent on colloidal particles. In what follows, we imagine the "solvent" to be simplified cosolutes within an implicit solvent bath. This subject has been well-studied (see e.g. refs. [192–196] ); our goal here is simply to show that the DFT transfer free energy provides a natural way of calculating depletion forces as well as transfer energies, and that even the approximated form in equation (3.31) yields non-trivial results for the depletion force.

We investigate a model consisting of two spheres that interact only by a hard wall potential of radius $r_s$. Each sphere also interacts with a bath of cosolutes through a 6-12 (van der Waals) potential: $V(r) =$

78

$4\epsilon \left( (\sigma/r)^{12} - (\sigma/r)^6 \right)$, with $\sigma = r_s + r_o$. With this model we examine the force as a function of the sphere separation $d$. Any force between the spheres is entirely due to cosolute-mediated effects.

When the solute particles are far apart, they dress themselves with cosolute solvation shells because of the attractive solute-cosolute potential. As we imagine moving the two solute particles closer together, eventually the repulsive region of one solute particle overlaps with the attractive region of the other solute particle, and vice versa. This situation is unfavorable for the solute particles, and the energy may be lowered by moving them further apart; hence there is a repulsive force at these distances (see Figure 3.5). As the solute particles continue to approach each other, the above repulsive region encroaches on the regions of space where the van der Waals potential is deeper. A larger amount of potentially favorable binding energy is removed per distance travelled, and the repulsive force due to "impeded-solvation" increases. The repulsive force is maximal when the solute separation $d$ is roughly $2\sigma$. For separations $d < 2\sigma$, the repulsive regions of the two solute spheres begin to overlap. This reduces the volume excluded, or more precisely repulsive to, cosolutes. This reduced excluded volume results in an attractive force which is the traditional depletion force. Eventually the depletion force becomes stronger than the above impeded-solvation force, and the net force becomes attractive. We note that such effects would not be present in standard GB/SA models of implicit solvation.

In general, direct inter-particle interactions must be superimposed on the above scenario. Which force dominates at a given separation will then depend on the values for $r_s$, $r_o$, and $\epsilon$, along with the strength of the direct interaction. The above-described repulsive effect has been observed before in hard-sphere solutes using the Derjaguin approximation to obtain an effective surface tension [191]. Here we see that the effect arises naturally from the presence of an attractive potential in the density functional theory.

Figure 3.5: **Solvent-induced force on a pair of "hard-wall" spheres** as a function of the separation distance, as obtained from equation (3.31). Spheres interact with cosolutes through a LJ potential (see text). The only parameters that determine the force are thus $\sigma$ and $\epsilon$, which appear in the LJ potential that enters into the DFT expression for the force. Each curve in the figure corresponds to a given well-depth $\epsilon$ in the sphere-cosolute potential. The depletion force is dominant at small separation, but there is a region in which the spheres are mutually repulsive due to lost attraction or "impeded-solvation" to the solvent.

Figure 3.6: **Schematic renderings of the solute spheres** in Figure 3.5 at several distances. a) The sphere-cosolute interaction is through a LJ potential, which is negative beyond a distance $\sigma = r_s + r_o$ (shown as the green region), and positive and repulsive for $d < \sigma$ (red region). The direct sphere-sphere interaction is only through a hard-wall potential of radius $r_s$. The cosolutes have radius $r_o$. (b) Sphere configuration when distance $d = 2\sigma$. An cosolute can just fit between the spheres at this distance- the LJ potential is zero in this configuration if the cosolute (dashed sphere) is centered directly between the solute particles. Such separations have positive force between the solutes in Figure 3.5a, due to "impeded-solvation": the repulsive interaction between one sphere-cosolute pair removes some of the attractive region from the other sphere-cosolute pair (region shown in magenta). At the separation shown in (c), the solvent-induced force between the spheres is now attractive; the volume of the removed attractive region now varies weakly with separation, and bringing the spheres closer together gains free energy by removing the depletion zone highlighted in blue.

## 3.5 Transfer Enthalpy and Entropy from Density Functional Theory

Having developed a DFT for cosolute-protein interactions, we can return to the experimental results analyzed in Chapter 2 and ask what constraints the observed data places on our simple models.

The transfer enthalpy $\delta H$ may be found from the free energy (equation (3.25) in Section 3.2) through $\partial(\beta \delta G)/\partial \beta$:

$$\delta H = \int d^3 r \, \mathcal{V}(\boldsymbol{r})\phi(\boldsymbol{r}) = \tilde{c} \int d^3 r \, \mathcal{V}(\boldsymbol{r}) e^{-\beta \mathcal{V}(\boldsymbol{r})} , \tag{3.32}$$

where $\phi(\boldsymbol{r})$ is given by

$$\phi(\boldsymbol{r}) = \frac{N e^{-\beta \mathcal{V}(\boldsymbol{r})}}{\int d^3 r' \, e^{-\beta \mathcal{V}(\boldsymbol{r}')}} \equiv \tilde{c} e^{-\beta \mathcal{V}(\boldsymbol{r})} , \tag{3.33}$$

The transfer entropy may be found from $-\partial \delta G/\partial T$ or directly from $\delta S = \beta(\delta H - \delta G)$:

$$\delta S = \beta \int d^3 r \, \mathcal{V}(\boldsymbol{r})\phi(\boldsymbol{r}) + N \log \left( \frac{1}{V} \int d^3 r \, e^{-\beta \mathcal{V}(\boldsymbol{r})} \right) . \tag{3.34}$$

We can now take a simple, heuristic model wherein we assume that the protein occupies a volume $V_n$ in the native state, and over that volume exerts an average potential $\epsilon_n$ on the cosolutes, relative to that of water. Likewise, the protein occupies an average volume $V_u$ in the unfolded ensemble, and exerts an average potential $\epsilon_u$ on the cosolutes. Then (3.32) and (3.34) reduce to:

$$\delta H_n = c \, V_n \epsilon_n e^{-\beta \epsilon_n} \tag{3.35}$$

$$\delta S_n = c \, V_n \left[ \beta \epsilon_n e^{-\beta \epsilon_n} + \left( e^{-\beta \epsilon_n} - 1 \right) \right] \tag{3.36}$$

with similar expressions holding for the transfer enthalpy and entropy of the unfolded state in terms of $V_u$ and $\epsilon_u$. This then gives us model predictions for $\delta \Delta H = \delta H_u - \delta H_n$ and $\delta \Delta S = \delta S_u - \delta S_n$.

The observation of entropy-enthalpy compensation means that the deviations from the line $\beta\delta\Delta G = 0$ are not large compared to the scale of $\beta\delta\Delta H$ or $\delta\Delta S$. The dimensionless quantity

$$\beta\delta\Delta G = -cV_u \left(\mathrm{e}^{-\beta\epsilon_u} - 1\right) + cV_n \left(\mathrm{e}^{-\beta\epsilon_n} - 1\right) \qquad (3.37)$$

defines a hypersurface as a function of the dimensionless variables $cV_u, cV_n, \beta\epsilon_u, \beta\epsilon_n$. We can estimate $cV_n$ for a 1M solution and a typical steric native volume for a two-state protein, $V_n \approx 2.5 \times 10^4 \text{\AA}^3$. Then $cV_n \approx 15$, perhaps larger depending on the size of the protein or cosolute.

Then the inequality $|\beta\delta\Delta G| < a$, where $a$ is a constant of order unity, defines a region in the space $\beta\epsilon_u, \beta\epsilon_n, cV_u$ bounded by the surfaces $\beta\delta\Delta G = a$ and $\beta\delta\Delta G = -a$. In Figure 3.7, we analyze what requirements entropy-enthalpy compensation imposes on the parameters of the model, by taking the two proteins that had the minimal and maximal free energy change upon transfer to 100g/l of cosolute. These systems correspond to Arc repressor in KCl, with $\beta\delta\Delta G = 6.1$ at the transition temperature in water, to RNase A in Butylurea, with $\beta\delta\Delta G = -8.1$ at the transition temperature in water.

The two surfaces $\beta\delta\Delta G = 6.1$ and $\beta\delta\Delta G = -8.1$ in Figure 3.7 are quite close: for Arc repressor for example, a difference in unfolding interaction energy of $\Delta\epsilon_u \approx 0.2k_BT$ can move points from the stabilizing surface to the destabilizing surface. We also observe that as $\epsilon_n$ and $\epsilon_u$ become increasingly negative, all systems, for both stabilizing and destabilizing cosolutes, converge to $V_u \approx V_n$. Thus, if the effective interaction potentials are net attractive, even small changes in unfolded volume can yield dramatically different values of unfolding transfer enthalpy $\delta\Delta H$ for example. We found that both Lysozyme/DMSO and RNase A/Butylurea systems decreased the volume experienced by the cosolutes upon unfolding, for most of the range of the energetic parameters. Strongly stabilizing cosolutes such as KCl showed significant increases in volume upon unfolding, if the interaction energies were small. This particular observation is consistent with previous findings for the volume increase upon unfolding of Trp-cage miniprotein in model hard-sphere cosolutes[35].
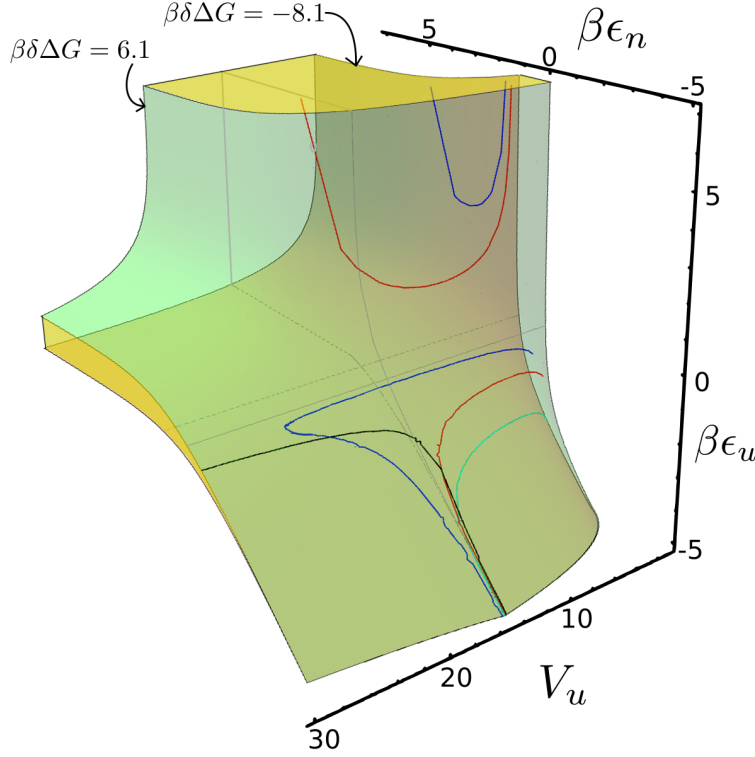
Figure 3.7: **Allowed volume in parameter space for the proteins we investigated, as predicted by the classical density functional theory** of protein transfer free energy [166] for the model described in Section 3.5. The two bounding surfaces shown bracket the observed transfer free energy transfer values for the proteins we considered. The transfer free energy is taken for cosolutes with concentration 100g/ml, at the transition midpoint in water, $\delta\Delta G(T_f^o)$. The upper surface in green corresponds to a transfer free energy of $\beta\delta\Delta G = 6.1$, observed for the transfer of Arc repressor to KCl. The lower surface in orange corrsponds to a transfer free energy of $\beta\delta\Delta G = -8.1$, observed for the transfer of RNase A to Butylurea. The volume bounded by the surfaces, where parameters characterizing all other proteins reside, is shaded yellow. The Arc repressor system is also constrained by the transfer enthalpy $\beta\delta\Delta H_f = 13.9$, which further restricts parameter values to lie on the black curve within the green surface. The RNase A system, with $\beta\delta\Delta H = -42$, is constrained to lie on the cyan curve. The transfer of Lysozyme to DMSO corresponds to the two red curves. The lower red curve corresponding to negative interaction energies obeys $V_u < V_n$, while the upper red curve crosses the plane $V_u = V_n$ for some energetic parameters. The transfer of RNase A to N-N' Dimethylurea corresponds to the two blue curves. Both red and blue curves lie between the surfaces. The intersection of the surfaces with the plane $cV_u = cV_n$, and with the plane $\epsilon_u = 0$ are shown as grey solid or dashed lines for the upper and lower surfaces respectively.

In fact, rigorous inequalities for the volume change upon unfolding as well as the protein-cosolutes effective interaction energies may be shown for the case of a stabilizing cosolute with negative change in unfolding enthalpy upon transfer ($\delta\Delta G > 0$ and $\delta\Delta H < 0$) and for a destabilizing cosolute with positive change in unfolding enthalpy upon transfer ($\delta\Delta G < 0$ and $\delta\Delta H > 0$). Examples of the former protein cosolute system are ubiquitin in ficoll or lysozyme, RNase A in glycine, and cytochrome c in trehalose. Interestingly, we did not find a protein that fit into the latter class among the proteins we investigated. Consider first the case in which $\beta\delta\Delta G > 0$ and $\beta\delta\Delta H < 0$, and assume that the protein-cosolute energies are negative with respect to water ($\epsilon_u < 0$, $\epsilon_n < 0$). From equations (3.35) and (3.37), letting $u \equiv -\epsilon_u$ and $n \equiv -\epsilon_n$, we have $V_u u e^u > V_n n e^n$ and $V_u(e^u - 1) < V_n(e^n - 1)$. Eliminating $V_u/V_n$ from these inequalities yields $ue^u/(e^u - 1) > ne^n/(e^n - 1)$. Since the function $f(x) = xe^x/(e^x - 1)$ is monotonically increasing for positive $x$, this inequality directly shows that $|\epsilon_u| > |\epsilon_n|$. Therefore $(e^n - 1)/(e^u - 1) < 1$ and $V_u/V_n < 1$. Thus the effective volume decreases upon unfolding. Similarly, for the case in which $\beta\delta\Delta H > 0$ and $\beta\delta\Delta G < 0$, $|\epsilon_n| > |\epsilon_u|$ and $V_u > V_n$.

Inequalities may also be obtained in the limit of strong entropy-enthalpy compensation, where $|\delta\Delta H| \gg 1$ and $|\delta\Delta G| \approx 0$, still assuming $\epsilon_n, \epsilon_u < 0$. Then equation (3.37) gives $V_u/V_n \approx (e^n - 1)/(e^u - 1)$. If $\delta\Delta H$ is large and positive, equation (3.35) gives $V_u/V_n \ll ne^n/ue^u$. Together these yield $|\epsilon_n| > |\epsilon_u|$ and $V_u > V_n$. If $\delta\Delta H$ is large and negative, $|\epsilon_u| > |\epsilon_n|$ and $V_n > V_u$.

For the more realistic case of $\epsilon_n, \epsilon_u > 0$, if $\beta\delta\Delta H < 0$, $\beta\delta\Delta G > 0$, and the interaction energies are larger than $k_B T$ ($\beta\epsilon_n, \beta\epsilon_u > 1$), then $\epsilon_u > \epsilon_n$; if $\beta\delta\Delta H > 0$, $\beta\delta\Delta G < 0$, and $\beta\epsilon_n, \beta\epsilon_u > 1$, then $\epsilon_n > \epsilon_u$.

Further analysis of the parameter space for the simple model introduced here, as well as more realistic models that include both bulk and a surface terms in the free energy functional, are a topic for future research.

## 3.6 Conclusions

In this chapter we have explored the application of the density functional framework to protein transfer free energies. We have focused primarily on conceptual questions, such as the role of solvent excluded volume, the temperature dependence of transfer free energies, and how the density functional theory (DFT) would reduce to a Volume + SASA model of transfer free energy.

We compared the DFT results with those from a simplified model that treated the protein as a tube with a given volume and surface area, on which cosolutes could condense. The DFT contains contributions from both enthalpy and entropy, so it allows for the calculation of the temperature-dependence of the transfer free energy.

A further development of the theory presented here which accounts for interparticle correlations while maintaining computational efficiency is an important topic for future research. As well, the calculation of transfer free energies was implemented here for a model system with simplified potentials that were parameterized to experimental values. One could extend this by implementing the theory using more realistic potential models, and all-atom representations of a protein or peptide. The various approximations involved in these potentials and models could then be tested and the limits of their validity determined through comparisons with experiment and simulation. The DFT framework may also provide a method to obtain computationally efficient but still accurate implicit solvent models for molecular dynamics simulation, a subject of immense practical importance. In general, the framework of density functional theory can provide a powerful tool to explore aspects of solvation in the context of protein folding, and can do so in a systematic way.

# Chapter 4

# Theoretical Considerations in Classical Density Functional Theory

## 4.1 Defining the Transfer Free Energy

Having shown the usefulness of classical density functional theory in treating cosolutes, we now turn to water. While cosolutes tend to be relatively dilute, such that the approximation that $\Delta\Phi' = 0$ used in Equation 3.25 is acceptable (because the cosolute-cosolute energy is not changed by the presence of an external potential), water is not dilute and this approximation is unlikely to be very good. Thus a suitable approach to determining $\Phi'$ must be found.

To begin, though, we broaden our scope to consider the transfer free energy in a more general sense. The usefulness of the transfer free energy (also known as the free energy of solvation in the context of polar solvents) extends beyond protein simulations–every molecule of interest, whether a biological molecule, a functional inorganic molecule, or a member of some nanostructure, operates in a background of solvent molecules that impact the function of the molecule one is looking at, but whose behaviour, in and of itself, is not of interest.
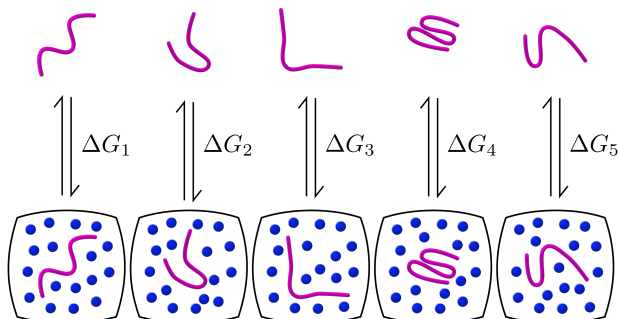
Figure 4.1: **Diagram of the transfer free energy.** The top row represents the states we can simulate cheaply—those without the solvent present. The bottom row represents the states we actually want to access—those with the solvent present. The vertical transitions constitute the transfer process. The transfer process can be broken up into two terms: the change in free energy of the solvent, $\Delta G_{\text{solvent}-\text{solvent}}$, and the change in free energy of the environment, $\Delta G_{\text{solute}-\text{solvent}}$

The transfer free energy can be implemented in computational studies in several different ways. The implicit solvent models discussed in Chapter 1 are typically implemented as an extra force at each time step of the simulation–*i.e.* the total force on the $i$th particle will be $\mathbf{F}_i = \mathbf{F}_i^{\text{forcefield}} + \nabla_i \Delta G$, where $\mathbf{F}_i^{\text{forcefield}}$ is the term arising from the explicit particles in the simulation and $\Delta G$ is the transfer free energy. Alternatively the transfer free energy can be implemented in a post-processing way. This makes use of the Tanford transfer model[168] to modify the weights and energies of the observed simulation states. This approach requires that the initial simulation offer sufficient sampling of the final states[197], and thus is limited to relatively small changes in the environment, such as a change from a protein in water to a protein in water with a small concentration of urea[144].

We need to take care to clarify what the transfer free energy consists of, and which parts of it we would like to calculate. Consider an initial state which consists of an isotropic bath of solvent molecules and, isolated from this bath, a solute molecule. This system has free energy $G_i$. The final state

will be the solute molecule dissolved in the bath of solvent molecules, and this state will have free energy $G_f$. Then the total change in free energy $\Delta G_{\text{total}}$ will have three terms: a term arising from the solute-solvent interaction, a term arising from the solvent-solvent interactions (accounting for both the change in entropy as the ensemble shifts and the change in solvent-solvent potential), and a term arising from the solute-solute interactions (again both in entropy and enthalpy). That is,

$$\Delta G_{\text{total}} = \Delta G_{\text{solute}-\text{solvent}} + \Delta G_{\text{solvent}-\text{solvent}} + \Delta G_{\text{solute}-\text{solute}} \quad\quad (4.1)$$

Experimental measures of the transfer free energy measure $\Delta G_{\text{total}}$, while typically implicit solvent models are concerned with calculating the first two terms, $\Delta G_{\text{solute}-\text{solvent}} + \Delta G_{\text{solvent}-\text{solvent}}$, which we will refer to as $\Delta G_{\text{solvent}}$. The reason for this is that if during the course of a simulation $\Delta G_{\text{solvent}}$ is implemented correctly, $\Delta G_{\text{solute}-\text{solute}}$ will then fall out naturally as the system adjusts; *e.g.* the swelling of a polymer in response to a solvent that makes surface exposure favourable. Also note that $\Delta G_{\text{total}}$ and $\Delta G_{\text{solvent}}$ are equal in the limit of a small rigid solute molecule, since then neither the solute energy nor its configurational space change.

We now wish to consider the general case of $\Delta G_{\text{solvent}}$ within the framework of classical density functional theory. To do this we take the solvent to be in an NPT ensemble and the solute molecule to be a fixed source of external potential. We can then take the solvent-solute term $\Delta G_{\text{solute}-\text{solvent}}$ as the minimum work required to insert this fixed potential, and the solvent-solvent term $\Delta G_{\text{solvent}-\text{solvent}}$ as the change in free energy of the solvent;

$$\Delta G_{\text{solvent}-\text{solvent}} = N\Delta\mu$$

where $N$ is the number of solvent molecules and $\Delta\mu$ is the change in chemical potential of the solvent molecules. The chemical potential of the solvent molecules is not the same in general as the chemical potential of the solute molecule, an important point we will return to. We will show below that $\Delta\mu$ is identically zero in the NPT ensemble.

As a simple example to motivate our conclusions, consider the free energy change to transfer an idealized solute into an ideal gas that initially has volume $V_i$ and pressure $P_i$. The idealized solute has a volume $v_0$ over which it produces a step potential of $U_0$. The configuration integral of the partition function of the gas after insertion is proportional to

$$\mathcal{Z} \propto [V_f - v_0(1 - e^{-\beta U_0})]^N$$

where $V_f$ is the final volume of the solvent + solute. $V_f$ is greater than $V_i$ if $U_0 > 0$ and less than $V_i$ if $U_0 < 0$; it is determined from the constancy of the pressure, $P_f = P_i$, where

$$P_i = \frac{N k_B T}{V_i}$$

The final pressure may be found from

$$P = - \left. \frac{\partial G}{\partial V} \right|_{V_f} = \frac{N k_B T}{V_f + v_0(e^{-\beta U_0} - 1)}$$

so that

$$V_f = V_i - v_0(e^{-\beta U_0} - 1)$$

Thus the change in free energy of the solvent is

$$N \Delta \mu = G_f - G_i = N k_B T \log \left[ \frac{V_f - v_0(1 - e^{-\beta U_0})}{V_i} \right] = 0$$

That is, the change in free energy of the ideal gas upon inserting the region of potential $U_0$ is 0, regardless of $U_0$ and $v_o$. The transfer free energy is then equal to the change in free energy of the environment–in this case $P\Delta V = -Pv_0(1 - e^{-\beta U_0})$.

## 4.2   The DFT Formulation

We consider the solute molecule being transferred as a fixed potential that acts on the solvent, and we ask what the difference in free energy of the

solvent is in going from the system without that potential to that with the external potential. So we write the free energy as a functional of the solvent density $\varphi$:

$$G = \int_{V_{sys}} d\mathbf{r}\, \mathcal{V}(\mathbf{r})\varphi(\mathbf{r}) + k_B T \left(\varphi(\mathbf{r}) \log \varphi(\mathbf{r}) - \varphi(\mathbf{r})\right)$$
$$+ \, \Phi[\varphi] \tag{4.2}$$

Here $\mathcal{V}(\mathbf{r})$ is the potential the solvent feels due to the solute, the integral is taken over the volume of the system $V_{sys}$, $\Phi$ contains all non-ideal terms in the free energy (*ie*, all two and higher particle correlation terms) and $k_B T$ is Boltzmann's constant times the temperature. While in earlier chapters we approximated $\Phi[\phi]$, here we will leave it as is, and manipulate it formally without specifying the form it takes.

While in Chapter 3 we considered the cosolute contribution to the transfer free energy in an NVT ensemble (while later allowing $V \to \infty$), in this chapter we consider the more general case of a transfer from vacuum to an arbitrary solvent bath, in an NPT ensemble.

Minimizing equation (4.2) subject to a fixed number of solvent particles amounts to minimizing the function

$$\mathcal{L} = G + \mu \left( N - \int_{V_{sys}} d\mathbf{r}\, \varphi(\mathbf{r}) \right)$$

where the Lagrange multiplier $\mu$ is the chemical potential. The Euler-Lagrange equation becomes:

$$\frac{\delta \mathcal{L}}{\delta \varphi} = \mathcal{V} + k_B T \log \varphi(\mathbf{r}) + \frac{\delta \Phi}{\delta \varphi}(\mathbf{r}) - \mu = 0 \tag{4.3}$$

We now define solvent redistribution energy density $\Phi' \equiv \frac{\delta \Phi}{\delta \varphi}$ and note that $\Phi'(\mathbf{r})$ can be thought of physically as the energy arising from solvent-solvent interactions; if we consider only the two-particle enthalpic component of $\Phi'$ we get

$$\Phi'(\mathbf{r}) = \int d\mathbf{r}'\, U(\mathbf{r} - \mathbf{r}')g(\mathbf{r}, \mathbf{r}')\varphi(\mathbf{r}')$$

where $U(\mathbf{r} - \mathbf{r}')$ is the solvent-solvent interaction potential and $g(\mathbf{r}, \mathbf{r}')$ the two-particle total correlation function.

Equation (4.3) gives an (implicit) expression for $\varphi(\mathbf{r})$,

$$\varphi(\mathbf{r}) = e^{\beta\mu} e^{-\beta(\mathcal{V}(\mathbf{r}) + \Phi'(\mathbf{r}))} \tag{4.4}$$

both sides of which can be integrated over the system volume to give

$$N = e^{\beta\mu} \int_{V_{sys}} d\mathbf{r} \, e^{-\beta(\mathcal{V}(\mathbf{r}) + \Phi'(\mathbf{r}))}$$

or

$$\mu = -k_B T \log\left(\frac{1}{N} \int_{V_{sys}} d\mathbf{r} \, e^{-\beta(\mathcal{V}(\mathbf{r}) + \Phi'(\mathbf{r}))}\right) \tag{4.5}$$

For the transfer problem we find the difference in free energy $N\Delta\mu$ between the initial case in which the external potential on the solvent is zero everywhere, and the final state. We thus write

$$\Delta\mu = -k_B T \log\left(\frac{\int_{V_f} d\mathbf{r} \, e^{-\beta(\mathcal{V}(\mathbf{r}) + \Phi'_f(\mathbf{r}))}}{V_i e^{\Phi'_i}}\right) \tag{4.6}$$

where $V_f$ and $V_i$ are the final and initial system volumes respectively, and $\Phi'_f$ and $\Phi'_i$ are likewise the final and initial solvent redistribution energy densities. In the denominator of equation (4.6) we have used the fact that, in the absence of an external potential, the system is homogeneous and isotropic, so any property of the solvent must be independent of position $\mathbf{r}$.

We now consider the conditions under which the transfer is made. We will assume that there is some region $\Omega_\infty$ in the solvent that is sufficiently far from the solute potential in the final system such that

$$\mathcal{V}(\mathbf{r} \in \Omega_\infty) = 0 \tag{4.7}$$

$$\varphi_f(\mathbf{r} \in \Omega_\infty) = \varphi_i \tag{4.8}$$

$$\Phi'_f(\mathbf{r} \in \Omega_\infty) = \Phi'_i \tag{4.9}$$

Figure 4.2: **Density, potential, and $\Phi'$ of water** as a function of distance from a spherical van der Waals potential. The interaction function $\Phi'(\mathbf{r})$ decays to zero on the same length scale as $\mathcal{V}(\mathbf{r})$.

Using equation (5.3) in equation (4.8) gives

$$\mathrm{e}^{\beta\mu_i}\mathrm{e}^{-\beta\Phi'_i} = \mathrm{e}^{\beta\mu_f}\mathrm{e}^{-\beta\left(\mathcal{V}(\mathbf{r}\in\Omega_\infty)+\Phi'_f(\mathbf{r}\in\Omega_\infty)\right)} \ . \tag{4.10}$$

Using (4.7) and (4.9) in equation (4.10) gives $\mathrm{e}^{\beta\mu_i} = \mathrm{e}^{\beta\mu_f}$, or

$$\Delta\mu = 0 \ , \tag{4.11}$$

i.e. the change in free energy of the solvent upon transfer is zero. Thus the solvent transfer free energy is given entirely by the solvent-solute interaction, which is simply the work done by the external potential to insert itself into the solvent bath.

$$\Delta G_{\mathrm{solvent}} = \Delta G_{\mathrm{solvent-solute}} \tag{4.12}$$

93

Equations (4.11) and (4.12) together are the principal result of this chapter.

We now consider a set of conditions that are sufficient for the existence of a region $\Omega_\infty$ that satisfies equations (4.7)-(4.9):

1. The transfer of the potential source (the solute) into the system occurs at constant solvent particle number $N$, pressure $P$, and temperature $T$.

2. The external potential acts on a finite, enclosed region of the system. This implies a fixed solute molecule transferred to a fixed position, so that we are addressing the solvation contribution to the transfer free energy.

3. Define lengths $\ell_C$ and $\ell_G$ as follows: At distances $r > \ell_C$, the direct correlation function $C(r)$ satisfies $|C(r) - 1| < \epsilon_C$, where $\epsilon_C$ is a system-dependent constant. Similarly, at distances $r > \ell_G$, the total correlation function $G(r)$ satisfies $|G(r) - 1| < \epsilon_G$, where again $\epsilon_G$ is a system dependent constant. We thus require that the direct correlation length $\ell_C$ and the total correlation length $\ell_G$ of the solvent molecules are both finite in the initial and final systems.

4. The system is sufficiently large that there exists a finite region $\Omega_\infty$ that is everywhere farther than $\ell_C + \ell_G$ from the region over which the potential acts.

The above conditions ensure that identical initial and final pressures yield identical values of $\varphi$ and $\Phi'$ far from the solute. In Figure 4.2 we plot the density relative to the equilibrium density $\phi_f / \phi_i$, the external potential $\mathcal{V}$, and the difference in the solvent-solvent interaction terms $\Phi'_f - \Phi'_i$ as a function of distance from a spherical van der Waals potential for a simulation of TIP3P water at 300K. The simulation was performed in GROMACS with the CHARMM27 forcefield (see Appendix B for specific parameters used). One feature that is immediately apparent is that the length scale over which $\Phi'_f - \Phi'_i$ is non-zero is very small–slightly larger than 1 nm. This is true even for the long-range electrostatic force; due to screening by the dielectric

medium of water (enhanced by the salt present in real biological systems) the effective range of the force will be small. One atom per $nm^3$ corresponds to a concentration of around 1 M, so we can conclude that up to that point a solute is well described by the infinite dilution limit.

## 4.3 Discussion

To return to the simple example in section 4.1, we consider the ideal gas form of the DFT transfer free energy, in which $\Phi' = 0$. Then

$$\phi = \frac{N}{V} e^{-\beta(\mathcal{V})} \tag{4.13}$$

and the work done to insert a step potential $U_0$ over a volume $V_0$ is then

$$\Delta G_{\text{solvent}} = \int d^3r \int_o^{U_0} d\mathcal{V} \frac{N}{V} e^{-\beta(\mathcal{V})} \tag{4.14}$$

$$= \frac{kTN}{V} V_0 (1 - e^{-\beta U_0}) \tag{4.15}$$

which is simply the ideal gas pressure times the change in total volume. We discussed in chapter 3 how this result also comes about by considering the constant volume system and letting $V \to \infty$. Thus in the ideal gas approximation, the solvent-solute transfer free energy is equal to $P\Delta V$, which is also the result given from the NVT ensemble in the limit of infinite V. As we will discuss below, these equalities are not true when we move away from the ideal gas.

The result in equation 4.11 is significant for two reasons. The first is that, to our knowledge, it is a general proof in the context of density functional theory. That the solvent-solvent term in the transfer free energy is zero was shown by Yu and Karpluss using different methods for a specific case[164]. Their result made use of the hyper-netted chain closure relation, which is only one possible closure relation of the Ornstein-Zernike equation. To our knowledge the identity $\Delta\mu = 0$ has not been extended to all closure relations, and thus our result here, which is independent of closure rela-

tion, is significant. Indeed, we have shown that the free energy of a bath of particles in the NPT ensemble is independent of the external potential so long as that potential is finite in range. Further, while our definition of the terms in the transfer free energy here makes it clear that $\Delta\mu$ corresponds to Karpluss and Yu's solvent re-organization term, this has not been appreciated in other literature on classical density functional theory. Ramirez, Mareschal, and Borgis [198], for example, do not make use of such a result in their discussion of the transfer free energy.

The other reason this result is significant is practical: it is a useful theorem in developing an implicit solvent model within the cDFT framework. This is the subject we turn to next.

# Chapter 5

# Including Solvent-Solvent Interactions in DFT

## 5.1 A Simple Form for the Solvent-Solvent Interaction

In developing an implicit solvent model, the overall goal is to develop a way of computing the solvent forces on protein atoms quickly; the main reason one uses such models in practice is to speed up simulations. To that end we require that the final equations for the force involve only the protein coordinates and constants—not the solvent density. That said, to address solvent forces in a DFT approach we obviously need to consider the solvent density—but we consider approximations to the free energy functional that give the solvent density a simple enough form that we can evaluate the solvent forces on a given atom as a single spatial integral over scalar functions. The challenge in the context of this work is to express the quantity $\Phi'$ (see section 3.2 equations 3.17 and 3.18) in a way that allows for both accuracy and speed.

After the assumption that the solvent-solvent interaction energy does not change upon transfer, which we have made in Chapter 3, the next simplest approximation mathematically is to assume that the solvent-solvent

interactions take the form of a delta function: $\Phi(r, r') = f(\phi)\delta(r - r')$. That is, limiting the free energy function to terms no higher than two-body,

$$G = \int d^3r\, kT \left( \phi(\mathbf{r}) \ln[\lambda^3 \phi(\mathbf{r})] - \phi(\mathbf{r}) \right) + \mathcal{V}(\mathbf{r})\phi(\mathbf{r}) + \int d^3r' \frac{\gamma}{2}\delta(\mathbf{r} - \mathbf{r}')\phi(\mathbf{r})\phi(\mathbf{r}')$$
(5.1)

This density functional can be minimized in the usual way to give

$$\frac{\delta G}{\delta \phi} = kT \ln[\lambda^3 \phi(\mathbf{r})] + \mathcal{V}(\mathbf{r}) + \gamma\phi(\mathbf{r}) = \mu,$$
(5.2)

which can be re-arranged to give

$$\lambda^3 \phi = e^{\beta\mu}e^{-\beta\mathcal{V}}e^{-\beta\gamma\phi}$$
(5.3)

where it is understood that $\phi$ and $\mathcal{V}$ depend on position. Equation 5.3 can be solved for the density to give

$$\phi = \frac{1}{\beta\gamma}W\left[\frac{\beta\gamma e^{\beta\mu}e^{-\beta\mathcal{V}}}{\lambda^3}\right]$$
(5.4)

where W is the Lambert-W function.

At this point we can make use of the main result from Chapter 4: in the NPT ensemble the chemical potential is independent of external potential, so long as the potential is finite in range. Thus in Equation 5.2 we can set $\mathcal{V} = 0$ everywhere, and note that this implies $\phi$ must be a uniform constant $\rho = N/V$ so that

$$\mu = kT \ln[\lambda^3 \rho] + \gamma\rho$$
(5.5)

Inserting this into Equation 5.4 then gives

$$\phi(\mathbf{r}) = \frac{1}{\beta\gamma}W\left[\beta\gamma\rho e^{\beta\gamma\rho}e^{-\beta\mathcal{V}(\mathbf{r})}\right]$$
(5.6)

If we take the usual assumption that the potential $\mathcal{V}$ has the form

$$\mathcal{V}(\mathbf{r}) = \sum_i v_i(\mathbf{r} - \mathbf{R}_i)$$

where $\mathbf{R}_i$ is the position of the $i$th atom in the protein, and $v_i$ is the potential energy at $\mathbf{r}$ due to atom $i$, then we can calculate the force on the $i$th atom from equation 5.6 via

$$\mathbf{F}_i = \int d^3r\, \phi(\mathbf{r}) \frac{\partial v_i}{\partial \mathbf{R}_i} \tag{5.7}$$

We can also calculate the transfer free energy through

$$\Delta G = \int d^3r \int_0^{\mathcal{V}_{\text{fin}}} d\mathcal{V}\, \phi(\mathbf{r}; \mathcal{V}) \tag{5.8}$$

The integral over $\mathcal{V}$ evaluates analytically using properties of the Lambert-$W$ function to give

$$\Delta G = -\frac{1}{\gamma} \int d^3r\, \mathrm{W}\left(\beta\gamma\rho e^{\beta\gamma\rho - \beta\mathcal{V}(\mathbf{r})}\right) + \frac{1}{2}\mathrm{W}\left(\beta\gamma\rho e^{\beta\gamma\rho - \beta\mathcal{V}(\mathbf{r})}\right)^2$$
$$+ V\left(\beta\rho + \frac{1}{2}\gamma(\beta\rho)^2\right) \tag{5.9}$$

Going back to section 4.3, it is worth noting here that equation 5.9 is not equal to $P\Delta V$ (when $\gamma \neq 0$). For this form of the interaction potential, the pressure can be calculated to be

$$P = \frac{k_B T}{V} + \frac{\gamma}{\rho^2} \tag{5.10}$$

and the change in volume is

$$\Delta V = \int d^3r \left[\frac{1}{\rho\beta\gamma}\mathrm{W}\left(\beta\gamma\rho e^{\beta\gamma\rho - \beta\mathcal{V}(\mathbf{r})}\right) - 1\right] \tag{5.11}$$

The non-interacting solvent case has the somewhat fortuitous result that the NPT ensemble transfer free energy, the NVT ensemble transfer free energy in the infinite dilution limit, and $P\Delta V$ are all equal. It appears this is not the case for models with solvent-solvent interaction terms.

## 5.2 Making the Form of Solvent-Solvent Interactions More General

While we could pursue the derivation in Section 5.1 in more detail, instead we will stop here to note that we can model the solvent-solvent interaction $\Phi'$ as a general function of $\phi$ without much increase in complication, so long as we continue to include the $\delta$-function form for $\Phi'$; that is, we can take the interaction free energy $G_{\text{int}} = \Phi[\phi]$ (see equation 4.2 of Chapter 4) to be

$$G_{\text{int}} = \int \int d^3r\, d^3r'\, \delta(\mathbf{r} - \mathbf{r}')\phi(\mathbf{r})f\left(\phi(\mathbf{r}')\right) \tag{5.12}$$

Performing the same procedure as above, in which we minimize the functional $G$ with respect to $\phi$ and use the fact that $\mu$ is independent of external potential, gives an implicit expression for $\phi$:

$$\phi(\mathbf{r}) = \rho e^{\beta(f(\rho)-f(\phi(\mathbf{r})))} e^{-\beta\mathcal{V}(\mathbf{r})} \tag{5.13}$$

where $\rho$ is the equilibrium density in the absence of external potential.

From equation 5.13 it is clear that $\phi(\mathbf{r})$ is local: it depends only on the form of $f(\phi)$, the value of the external potential at $\mathbf{r}$, and constants; $\phi(\mathbf{r})$ does not depend on the value of $\phi$ anywhere other than $\mathbf{r}$.

We set the restriction on equation 5.13 that $\frac{\partial \phi}{\partial \mathcal{V}} < 0$. This has both physical and computational motivations. Physically, we want to ensure that our equations are such that an increase in external potential never leads to an increase in solvent density. Computationally, we need to ensure that $\phi(\mathcal{V})$ is a one-to-one function, and thus equation 5.13 has a unique solution for $\phi$ given $\mathcal{V}$. Performing the derivative in equation 5.13 gives

$$\frac{\partial \phi}{\partial \mathcal{V}} = \frac{-\beta\phi}{1 + \beta\frac{\partial f}{\partial \phi}\phi} \tag{5.14}$$

Since $\phi$ and $\beta$ are always positive, this implies that we require

$$\beta\frac{\partial f}{\partial \phi}\phi > -1 \tag{5.15}$$

Again, since $\phi$ and $\beta$ are always positive, a sufficient (but not necessary) condition is that $\frac{\partial f}{\partial \phi} > 0$. Relaxing this inequality requires knowing something about $\phi$, which for all but the simplest $f(\phi)$ is non-analytic. At least in practice, after $\phi(\mathcal{V})$ has been found numerically, equation 5.15 can be used to confirm that $\phi(\mathcal{V})$ is one-to-one.

As an example, we will take

$$f(\phi) = \frac{k}{2\beta}(\phi - \rho)^2 \qquad (5.16)$$

For this model we plot $\phi/\rho$ as a function of $\mathcal{V}/k_B T$ in Figure 5.1. Even though for this interaction function it is clear that $\frac{\partial f}{\partial \phi} > 0$ does not hold, visually the function can be seen to satisfy $\frac{\partial \phi}{\partial \mathcal{V}} < 0$. And indeed, when we use the calculated values of $\phi$ to plot the quantity in equation 5.15 in Figure 5.1, we see that the inequality is satisfied. Finding the general conditions on $f(\phi)$ for equation 5.15 to be satisfied is a topic for future work.

## 5.3  Fitting the Model to Data

We would like to fit the model of section 5.2 to simulation data, so that it reproduces the correct forces that the solvent exerts on the protein in any given protein configuration. In principle, we could allow both $\mathcal{V}(r)$ and $f(\phi)$ to be defined by a spline function of general form. This gives a model with as many parameters are there are points in these two splines, though in practice these will be constrained by equation 5.15, and potentially by further constraints we may place on how fast these functions can change. Below we will make the problem more tractable by developing a procedure to fix $\mathcal{V}(r)$.

It is important to note that $\mathcal{V}(r)$ can be different for each atom type in the protein—and indeed should be different. On the other hand, $f(\phi)$ is a single function and should depend only on the solvent used. Thus it is not sufficient to simulate the transfer of a single atom to the desired solvent and then find the $\mathcal{V}(r)$ and $f(\phi)$ which reproduce the solvent density and transfer free energy; instead we must consider various combinations of atom

Figure 5.1: **The solvent density $\phi$ *vs.* external potential** for the interaction function 5.16 for various values of $k$.

types at various distances from each other. Of particular interest must be the regions in which the solvent feels the potential from more than one atom, as we expect the solvent density in those regions to be the most difficult to capture. Because $f(\phi)$ is a non-linear function whose parameters are limited only to the number of spline points we choose to take (in what follows we will use twenty), using the transfer free energy from *e.g.* a single atom would not uniquely determine $f(\phi)$.

The procedure we will use to determine $\mathcal{V}(r)$ for various atom types and $f(\phi)$ can be illustrated by example. We take two uncharged $C_\alpha$ atoms from the CHARMM27 forcefield. We then treat these atoms as bonded and use bond contraints to fix them at various distances from each other ($d = 0.12, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8$ nm). For each distance we calculate

Figure 5.2: **Testing the condition in equation** 5.15; indeed $\beta\frac{\partial f}{\partial \phi}\phi > -1$.

the transfer free energy from vacuum to TIP3P water with the Bennet Acceptance Ratio method in GROMACS. We also get the density of water molecules around each configuration. We then seek to find $\mathcal{V}(r)$ and $f(\phi)$ which reproduce both the transfer free energies and the water density profiles for all eight distances.

Figure 5.3 shows the simulation results. The solvation shells for both atoms are clearly visible. Also of interest is the desolvation barrier observed in panel B. It is worth noting that desolvation barriers cannot be captured by SASA models as the surface area of the two atoms monotonically decreases as they are brought together.

Rather than allowing both $\mathcal{V}(r)$ and $f(\phi)$ to vary independently, we fixed $\mathcal{V}(r)$ for a given $f(\phi)$ by using $f(\phi)$ to calculate $\phi(\mathcal{V})$, then finding the

103

A



$d = 0.12$ nm
$\Delta G = 13.33 \pm 0.07$ kJ/mol

$d = 0.20$ nm
$\Delta G = 17.7 \pm 0.1$ kJ/mol

$d = 0.30$ nm
$\Delta G = 22.19 \pm 0.05$ kJ/mol

$d = 0.40$ nm
$\Delta G = 25.69 \pm 0.07$ kJ/mol

$d = 0.50$ nm
$\Delta G = 27.65 \pm 0.06$ kJ/mol

$d = 0.60$ nm
$\Delta G = 27.5 \pm 0.1$ kJ/mol

$d = 0.70$ nm
$\Delta G = 27.00 \pm 0.04$ kJ/mol

$d = 0.80$ nm
$\Delta G = 27.23 \pm 0.07$ kJ/mol

B



Figure 5.3: **Transfer of a pair of bonded carbon atoms from vacuum to water.** Transfers were performed using the CHARMM27 forcefield with TIP3P water. A) Density profiles for each distance; darker areas correspond to higher water densities. The excluded volume (lighter spheres around each carbon) and solvation shells (darker rings around the lighter spheres) are clearly visible here. B) Transfer free energy as a function of $d$. Notable is the desolvation barrier observable around $d = 0.5$nm. The local maximum in $\Delta G$ arises because the finite size of water molecules induces preferred and avoided separation distances of the carbons. Uncertainties on the transfer free energy are smaller than the symbol size.

inverse of this function $\phi^{-1}(\phi|f)$ and applying it to the observed simulation density profile $\phi_{\mathrm{sim}}(r)$ for a single carbon atom. So

$$\mathcal{V}(r) = \phi^{-1}(\phi_{\mathrm{sim}}(r)|f) \tag{5.17}$$

We then vary $f(\phi)$ (and hence $\mathcal{V}(r)$) to minimize the difference between the DFT calculated free energies and the explicit solvent free energies.

We parameterized $f(\phi)$ as a cubic spline with 20 spline points and used a trust-region minimization approach[199] to find those points. After performing several stages of minimization, we arrive at a function $f(\phi)$ which captures the results of Figure 5.3 reasonably well. This function is plotted in Figure 5.6, and the resulting transfer free energies are shown in Figure 5.4. The cDFT model parameterized here captures the shape of the $\Delta G$ *vs d* curve, and in particular captures the desolvation effect observed when $d = 0.25$ to $0.3$ nm. As mentioned above a solvent accessible surface area approach cannot capture this–shown is the best fit SASA curve for these free energies, even if the fit is of similar quality with respect to the residual.

There are several other comparisons we can make. The potential function $\mathcal{V}$ arrived at from equation 5.17 (which implicitly uses $f(\phi)$) can be compared to the actual potential the solvent molecules experience. In the case of TIP3P water and uncharged solute atoms, this comparison is straightforward as the water hydrogen atoms have no van der Waals radius. The comparison between the two potentials is shown in Figure 5.5. Relative to the true potential, the effective potential has extra local maxima and minima, which correspond to the solvation shells of water. These features must be present because we recast the many-body behaviour of water into the effects due to a local potential. Additionally, the well depth is somewhat greater for the effective potential, and the hard wall does not rise quite as fast. Still the overall form is roughly consistent with what we might expect *a priori*.

We can also compare the effective solvent-solvent interaction function $f(\phi)$ with a function derived directly from simulations. We start by recasting

Figure 5.4: **Transfer free energies with DFT and surface area**. The free energies from Figure 5.3 are compared with the best fit free energies from the cDFT model developed in the text and the best fit free energies from a surface area model. While the residuals of the fits are comparable, the surface area model fails to capture the desolvation barrier observed from 0.25 to 0.3 nm. The SASA curve was generated by taking the analytic form of the surface area for two spheres and allowing the effective solvent radius and the proportionality constant to be varied.

106

Figure 5.5: **Effective and actual potentials**. The red curve gives the potential arrived at by finding the $f(\phi)$ that best reproduces the two-carbon transfer energy data and requiring $\mathcal{V}(\mathbf{r})$ to reproduce the observed radial distribution function. The green curve is the actual potential (*i.e.* a van der Waals potential).

equation 5.13 as an equation for $f(\phi) - f(\rho)$:

$$f(\phi) - f(\rho) = -\mathcal{V}(r) - k_B T \ln\left(\frac{\phi}{\rho}\right) \tag{5.18}$$

This is in fact exactly how we extracted $\Phi'$ in Figure 4.2; here $f(\phi) - f(\rho)$ has taken the place of $\Phi'$. From a simulation of an isolated carbon atom we determine $\phi$, and then $\mathcal{V}$ in equation 5.18 is taken as the actual potential. Then we can find for each $r$ a $\mathcal{V}$ and a $f(\phi) - f(\rho)$ (or $\Phi'_f - \Phi'_i$ in our

Figure 5.6: **Solvent-solvent interaction potentials.** The red curve is the best fit $f(\phi)$ arrived at from the DFT approach, while the gree curve is the observed $\Phi'$ from the radial distribution function.

more general notation) and plot these against each other. This we do in Figure 5.6. Two interesting features emerge: one is that the overall shape of the functions is similar, and the other is the jump in $f(\phi)$ around $\phi = 80$ nm$^3$. This jump corresponds to a "doubling" back of the "true" solvent-solvent interaction function, which is no longer a single-valued function at these densities. It is encouraging that this feature emerged from a routine which simply varied $f(\phi)$, and minimized the difference between observed and calculated free energies. That the functions do not overlap entirely is to be expected, as the effective interaction potential involves projecting essentially non-local interactions onto purely local ones.

Finally, we can compare the transfer free energy of atoms not used in the fitting procedure to observed explicit solvent transfer free energies. To do this we calculated the transfer free energy of several uncharged amino acids in GROMACS with the CHARMM27 forcefield using the Bennet acceptance ratio method. These amino acids are composed of various atoms; in addition to several types of carbon (carbon atoms with slightly different force-field parameters are used in various positions in the amino acid), nitrogen, oxygen, sulfer, and hydrogen are also present. For each of these the effective $\mathcal{V}(r)$ was determined from observed radial distribution functions and $f(\phi)$ through equation 5.17 and the total transfer free energy calculated in the DFT approach. The result is shown in Figure 5.7. The results show general agreement, with a correlation coefficient of 0.7, and absolute magnitudes that are comparable. By comparison we can consider the best fit SASA prediction of the transfer free energy *vs* the explicit simulation free energy, which has a correlation coefficient of 0.65. To compare these another way, using the Kendall $\tau$ test, we arrive at a significance of $p = 0.001$ for the relationship between the DFT transfer free energy and the explicit solvent transfer free energy, and a significance of $p = 0.003$ for the relationship between SASA and the explicit solvent transfer free energy. We note here that the correlation between SASA and explicit transfer free energy is the best possible correlation, while the correlation between the DFT predictions and the explicit transfer free energy is a starting point which we expect would be improved by fitting $f(\phi)$ (and potentially the effective potentials $\mathcal{V}(r)$ ) to a more representative set of candidate structures rather than two carbon atoms.

## 5.4 Discussion

We have shown that the protein-water interaction can be captured at least as well as existing implicit solvent approaches by a classical density functional theory with purely local solvent-solvent interactions. The strengths of this method are that it captures important effects lacking in traditional solvent accessible surface area models, while still being fast to implement. The
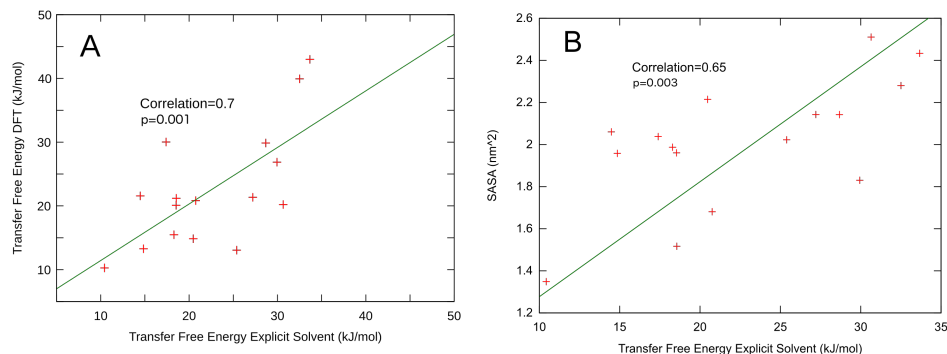
Figure 5.7: **A) Comparison of DFT predicted free energies with results from explicit solvent simulations**. The explicit solvent results were obtained with the Bennet acceptance ratio method implemented in GROMACS. The DFT predicted free energies were obtained as described in the text. B) Comparison of SASA predicted free energies with results from explicit solvent simulations.

reason for this speed is the approximation that allows the density to be a single valued function of the potential. Thus a table of $\phi(\mathcal{V})$ values can be computed once for a given $f(\phi)$ (which does not change from simulation to simulation) and then $\phi$ can be determined at any point from the atom-dependent $\mathcal{V}(r)$ functions. Given that each of these functions can be stored as a cubic spline, the total transfer free energy can be computed in a single integral over the space surrounding the protein, a procedure of comparable number of operations to the evaluation of surface area[200].

To make this last point more quantitative, we examine the calculation of solvent accessible surface area implemented in GROMACS, described by Eisenhaber *et al*[201]. Here, to calculate surface area, each atom in the molecule is assigned a set of points around it. The distance from each of these points to neighbouring atoms is calculated, and then each point is either declared buried or exposed. From this the exposed surface of the atom is determined. Such a method could also be applied to the DFT approach described in this chapter, with two differences: 1) the density must be evaluated at each point, and 2) the number of points needed to accurately

Figure 5.8: Testing convergence of a DFT free energy calculation. The free energy is plotted against the number of points for which the effective external potential was neither zero nor infinite, divided by the number of atoms. Also plotted is the deviation from true SASA vs the number of points per atoms from [201]. The two converge in a similar scale.

determine the transfer free energy may be different than the number needed to determine the solvent accessible surface area, particularly as the points would need to be distributed in a spherical shell around each atom rather than on a surface.

Based on the benchmarks described in [201], it requires on the order of $10^6$ floating point operations to compute the distance of neighbouring atoms to approximately $3 \times 10^5$ points. Evaluating a cubic spline at each of those points and summing the result can be expected to add on the order of 10 operations, an increase of approximately 30 %.

The question of how many points would be required is more difficult to answer without fully implementing the implicit solvent method and testing it on various proteins. We can make an estimate, though, from our initial calculations. The DFT values for $\Delta G$ in Figure 5.7 were computed with a three-dimensional integration grid which covered all the space surrounding

111

the protein. To make a better comparison to the method in [201], we can imagine a method in which only the points around each atom for which the effective external potential is neither zero nor infinite are considered. In Figure 5.8 we look at the free energy of transfer of a Leucine amino acid (which contains 22 atoms) calculated by DFT, *vs* the number of non-zero-non-infinite points. The value converges around 1000 points per atom. In [201] the number of points per atom considered ranged from 300 to 1500, depending on the accuracy wished. All together we estimate that the DFT calculation can be performed with approximately twice the operations a surface area calculation can be.

The procedure outlined here calculates the non-polar component of the transfer free energy. While density functional theory could be extended to consider the electrostatic aspects of the transfer free energy, the more limited approach described here allows this procedure to be integrated with existing methods, particularly GB/SA. Despite its flaws, GB/SA is undoubtedly the most widely used implicit solvent model in molecular dynamics simulations, and thus the ability for new models to integrate and improve upon it must be counted as an advantage.

There remains much work to do with this model, which we will consider in the next chapter, on future research directions.

# Chapter 6

# Conclusions and Future Directions

## 6.1  Summary

In this thesis we studied protein-solvent interactions and developed a classical density functional theory for those interactions. The interactions between the protein and its environment are inseparable from the intra-protein interactions. Proteins do not evolve in a vacuum, and the way in which the intra-protein interactions come together to create the folded, functioning protein depends intrinsically on the protein-solvent interactions.

In chapter 2 we examined some of the experimental literature on protein-solvent interactions. Specifically we looked at how the introduction of a cosolute changed the enthalpy and entropy of unfolding. We found that the change in unfolding entropy and enthalpy upon the introduction of a cosolute was an example of a broad class of phenomena known as entropy-enthalpy compensation, in which the change in free energy is small relative to the change in both entropy and enthalpy. There has been a long-standing debate concerning the significance of entropy-enthalpy compensation, and it has been suggested that the effect is due to highly correlated experimental uncertainties in entropy and enthalpy. We showed that there are indeed highly correlated uncertainties, but that the effect is still significant in light

of them, at least for the case of the transfer of various proteins from water to water with various cosolutes. This is not a general statement about all systems which seem to display entropy-enthalpy compensation, though, and wider application of procedures similar to the one we developed here to measure the uncertainty in entropy and enthalpy could establish which systems genuinely display the effect and which have uncertainties so large as to be unable to determine it.

In chapter 3 we continued to examine cosolutes (or osmolytes as we also refer to them), this time from the standpoint of developing a general theory. We develop a classical density functional theory for protein-cosolute interactions, which assumes that the density of water and cosolute are perfectly correlated, in the sense that the total density at any point is fixed, and further assumes that the cosolute-cosolute interaction energy doesn't change upon transfer of the protein. We show that the theory gives good results for a monatomic gas, even in the regime in which the gas is significantly non-ideal. This theory can reproduce the observed transfer free energies of moving side chains from water to water plus urea, and captures the temperature dependence of the transfer free energy. We showed that the cDFT theory developed in chapter 3 reduces to a SASA approach with suitable parameters. Our approach reproduces the general phenomenon of a desolvation barrier in a natural way as well. Finally we examine how a simple cDFT model can be related to entropy-enthalpy compensation and use the data from chapter 2 to restrict the parameter space of this model.

In chapter 4 we take a step back to examine the transfer problem in cDFT in more general terms. We show that the free energy of a particle bath at fixed temperature and pressure is independent of the external potential provided the external potential is finite in range–this is the dilute solute limit in the context of transfer. We discuss the implications of this finding–in particular the fact that all transfer free energy comes from the solute-solvent interactions–and note its usefulness in working with cDFT.

In chapter 5 we develop a classical density functional theory of water to implement as an implicit solvent. To make the mode fast to execute we assume solvent-solvent interaction can be modelled by a completely local

way. In exchange we allow the local function of $\phi$ to take on any form, subject to conditions to keep the solution physical. We fit the resulting theory to explicit water simulations. The model shows good agreement with explicit simulations, while capturing the desolvation barrier and temperature dependence, and while maintaining speed by requiring that $\phi$ depend only on the local potential.

Each of these chapters involved different aspects of protein-solvent interactions and their theoretical description, particularly within the framework of classical density functional theory.

## 6.2 Theoretical Work in cDFT

While classical density functional theory has a long history, there is still work to do on the theoretical end. The functional often described as the starting point for cDFT looks like this:

$$G = \int d^3r \, k_B T[\phi \ln \phi - \phi] + \mathcal{V}(\mathbf{r}) + \int d^3r' U(\mathbf{r} - \mathbf{r}')g(\mathbf{r} - \mathbf{r}')\phi(\mathbf{r})\phi(\mathbf{r}') \quad (6.1)$$

But, as we note in the introduction, this cannot be correct; for a hard-sphere model in the absence of an external potential the energy terms are zero, which implies an ideal gas solution no matter what density we take. Thus there must be a multi-particle entropy functional.

Entropy functionals, however, have seen remarkably little development. Early efforts by Nettleton and Green[202] involved taking the known form of the total density functional from radial distribution functions and reverse-engineering a numeric entropy functional. More recent approaches have looked at entropy functionals of 3-body and higher correlation terms[203], which limits their applicability to the use of cDFT in practice.

One alternative approach which seems promising is to express the entropy as an integral series in correlation functions; that is, the total entropy would be

$$S = S_{\text{ideal}} + \int d^3r \, [p(\mathbf{r})S_{\text{present}} + (1 - p(\mathbf{r}))S_{\text{absent}}] \quad (6.2)$$

where $S_{\text{ideal}}$ is the ideal gas entropy, $p(\mathbf{r})$ is the probability of finding a particle at position $\mathbf{r}$, $S_{\text{present}}$ is the entropy of the system of $N-1$ particles given that a particle is at $\mathbf{r}$, and $S_{\text{absent}}$ is the entropy of the system given that a particle is not at $\mathbf{r}$. This type of approach gives rise to an integral series reminiscent of the Orstein-Zernike equation.

We have not had time to pursue this line of analysis in this thesis, but it is a direction for future research. And whether or not this particular equation bears fruit, the search for analytic expressions for the entropy functional should not be abandoned.

## 6.3 Error Analysis in Thermodynamic Data

We developed a technique for determining the uncertainty in fitting heat capacity or fraction of unfolded data and extracting thermodynamic parameters. This technique is useful, but it does overlook one key source of uncertainty: the uncertainty in the underlying experiment. That the model uncertainty is important can be seen from the plots in figures 2.2. For many regions in those plots the difference between models is well outside the uncertainty calculated by the monte carlo technique. This implies the choice of model itself brings in a great deal of uncertainty. Given that most of the literature seems to use the temperature-independent $\Delta C_p$ model without checking if a linear or non-linear $\Delta C_p$ model would give the same results, this is a worrisome finding.

The uncertainty in the underlying experiment is more difficult to quantify. Looking at table 2.3 we see that, for example, RNase A, different experiments gave different values for $\Delta H_f^0$; that is, the enthalpy of unfolding in the absence of cosolute. The variance in $\Delta H_f^0$ values is somewhat higher than the variance in $T_f$ values, which raises the concern that this experimental uncertainty contributes to entropy-enthalpy compensation in a way we have not accounted for. Given the statistical significance of our result it is unlikely this effect will reverse our conclusion, but this is an important consideration for future work. Such work will likely require collaboration with experimentalists, as literature sources typically quote a single value for

116

each thermodynamic parameter, rather than a range of values indicative of the multiple samples they may or may not have tested.

## 6.4   Implementing the cDFT Implicit Solvent

We have developed an implicit solvent model based on classical density functional theory. Of course, the development of the model is the first step towards its adoption by researchers performing molecular dynamics. Implementing the model across force-fields and MD software packages is an enormous challenge well outside the scope of this thesis, or indeed even a single research group. Nonetheless this is the ultimate goal of the project.

One concern moving forward is what the target system should be for fitting the cDFT parameters. We have used explicit water simulations because we had easy access to them and they provide a high level of detail about energies and distribution functions. But the case could be made that experimental results would provide a better target system, and lead to a set of parameters which would not depend on any one forcefield. Experimental results have the disadvantage, though, of not cleanly separating the non-polar and polar contributions to the transfer free energy. One possibility would be to start with experimental transfer free energies and subtract off the predicted polar component (either in generalized Born or Poisson-Bolztman) to obtain quasi-experimental non-polar components of the transfer free energies.

We have assumed in this thesis that the cDFT model we developed would apply to the non-polar contributions to the transfer free energy. There is no reason *a priori* why we could not also apply cDFT to the polar contribution– that is, the electrostatic interactions. Doing so would require a significant modification of the work presented here. The water density in such an approach would no longer be a function of position only but also orientation; $\phi = \phi(\mathbf{r}, \Omega)$. The solvent-solvent interaction energy would then also be a function of both molecule separation and the orientation of each molecule, which greatly increases its computational cost. One way around this might be to break the density up into two fields: a position dentisy $\phi(\mathbf{r})$ and a po-

larization density $\mathbf{p}(\mathbf{r})$. Performing two three-dimensional integrals would be faster than one five-dimensional integral, but the accuracy of such an approximation needs to be assessed. It is not clear whether a purely local approach such as the one adopted in chapter 5 would continue to be effective. On the other hand, it is not obvious that it will not be effective, and extending the theory to polar interactions is a topic of future interest.

## 6.5   Parting Thoughts

All science relies on others to truly succeed. Findings left on their own do not advance human knowledge. But work on methods perhaps uniquely relies on the adoption of other scientists to bear fruit. In this thesis we have proposed a new method of looking at transfer free energies. The extent to which this method contributes to the field is, in one sense, out of our control. Regardless of its intrinsic merits, if it does not resonate with the community in a way that spurs others to work on implementing it in various forcefields and MD packages it will not be widely used and its impact will be small. On the other hand if it does find wide implementation it could have a greater impact than other models, regardless of their intrinsic merit. In working within the existing framework of decomposing the transfer free energy into polar and non-polar components, and in targeting the non-polar parts, we hope this model will be easily adaptable to existing MD approaches and therefore more likely to see wide implementation.

If we examine the literature on GB/SA, for example, after the initial paper by Qiu *et al*[48], following papers by various groups parameterized the model for AMBER[204], CHARMM[205], OPLS[206], and GROMOS[207] forcefields, suggested surface area calculation algorithms, and ported the model to various MD packages such as GROMACS[208] and NAMD[206]. This work underlaid the widespread adoption GB/SA enjoys today.

Moving forward then, the future work outlined in this chapter will be an important component of seeing the work in this thesis fulfill its promise, but equally important will be networking with other researchers to see these results implemented in a variety of existing software. We hope our shoulders

118

provide an ample platform for those that would follow.

# Bibliography

[1] G Sliwoski, S Kothiwale, J Meiler, and E W Lowe. Computational methods in drug discovery. *Pharmacological Reviews*, 66:334–395, 2014.

[2] William C Guest. *Template-Directed Protein Misfolding in Neurodegenerative Disease*. PhD thesis, The University of British Columbia, 2012.

[3] F Crick. Central dogma of molecular biology. *Nature*, 227:561, 1970.

[4] Andrea Ilari and Carmelinda Savino. Protein structure determination by x-ray crystallography. In Jonathan M Keith, editor, *Bioinformatics*, volume 452 of *Methods in Molecular Biology*, pages 63–87. Humana Press, 2008. ISBN 978-1-58829-707-5.

[5] M Billeter, G Wagner, and K W uthrich. Solution nmr structure determination of proteins revisisted. *Journal of Biomolecular NMR*, 42:155–158, 2008.

[6] Chengsong Liu, Dwayne Chu, Rhonda D. Wideman, R. Scott Houliston, Hannah J. Wong, and Elizabeth M. Meiering. Thermodynamics of denaturation of hisactophilin, a $\beta$-trefoil protein. *Biochemistry*, 40(13):3817–3827, 2001. doi: 10.1021/bi002609i. URL http://pubs.acs.org/doi/abs/10.1021/bi002609i. PMID: 11300762.

[7] Valerij Ya. Grinberg, Natalia V. Grinberg, Tatiana V. Burova, Michele Dalgalarrondo, and Thomas Haertlé. Ethanol-induced conformational transitions in holo--lactalbumin: Spectral and calorimetric studies. *Biopolymers*, 46(4):253–265, 1998. ISSN 1097-0282. doi: 10.1002/(SICI)1097-0282(19981005)46:4⟨253::AID-BIP7⟩3.0.CO;2-O.

URL
`http://dx.doi.org/10.1002/(SICI)1097-0282(19981005)46:`
`4<253::AID-BIP7>3.0.CO;2-O.`

[8] WF Claussen and MF Polglase. Solubilities and structures in aqueous aliphatic hydrocarbon solutions. *Journal of the American Chemical Society*, 74(19):4817–4819, 1952.

[9] R. Aveyard and A. S. C. Lawrence. Calorimetric studies on n-aliphatic alcohol + water and n-aliphatic alcohol + water detergent systems. *Trans. Faraday Soc.*, 60:2265–2278, 1964. doi: 10.1039/TF9646002265.

[10] E. M. Arnett and D. R. McKelvey. In J. F. Coetzee and C. D. Ritchie, editors, *Solute-Solvent Interactions*, chapter 6, pages 344–395. Dekker, New York, 1969.

[11] P. L. Privalov and S. J. Gill. Stability of protein-structure and hydrophobic interaction. *Adv. Prot. Chem.*, 39:191–234, 1988.

[12] Rufus Lumry and Shyamala Rajender. Enthalpyentropy compensation phenomena in water solutions of proteins and small molecules: A ubiquitous properly of water. *Biopolymers*, 9(10): 1125–1227, 1970. ISSN 1097-0282. doi: 10.1002/bip.1970.360091002. URL `http://dx.doi.org/10.1002/bip.1970.360091002`.

[13] Jack D. Dunitz. Win some, lose some: enthalpy-entropy compensation in weak intermolecular interactions. *Chemistry and Biology*, 2(11):709 – 712, 1995. ISSN 1074-5521. doi: http://dx.doi.org/10.1016/1074-5521(95)90097-7. URL `http://www.sciencedirect.com/science/article/pii/1074552195900977`.

[14] Hong Qian and J. J. Hopfield. Entropyenthalpy compensation: Perturbation and relaxation in thermodynamic systems. *The Journal of Chemical Physics*, 105(20):9292–9298, 1996. doi: http://dx.doi.org/10.1063/1.472728.

[15] Kim Sharp. Entropyenthalpy compensation: Fact or artifact? *Protein Science*, 10(3):661–667, 2001. ISSN 1469-896X. doi: 10.1110/ps.37801. URL `http://dx.doi.org/10.1110/ps.37801`.

[16] Andrew T Fenley, Hari S Muddana, and Michael K Gilson. Entropy–enthalpy transduction caused by conformational shifts can

obscure the forces driving protein–ligand binding. *Proceedings of the National Academy of Sciences*, 109(49):20006–20011, 2012.

[17] John D. Chodera and David L. Mobley. Entropy-enthalpy compensation: Role and ramifications in biomolecular ligand recognition and design. *Annual Review of Biophysics*, 42(1):121–142, 2013. doi: 10.1146/annurev-biophys-083012-130318. URL `http://www.annualreviews.org/doi/abs/10.1146/annurev-biophys-083012-130318`. PMID: 23654303.

[18] Emilio Gallicchio, Masahito Mogami Kubo, and Ronald M. Levy. Entropyenthalpy compensation in solvation and ligand binding revisited. *Journal of the American Chemical Society*, 120(18): 4526–4527, 1998. doi: 10.1021/ja974061h.

[19] A. Ben-Naim and Y. Marcus. Solvation thermodynamics of nonionic solutes. *J. Chem. Phys.*, 81(4):2016–2027, 1984. doi: http://dx.doi.org/10.1063/1.447824.

[20] Remo Perozzo, Gerd Folkers, and Leonardo Scapozza. Thermodynamics of proteinligand interactions: History, presence, and future aspects. *Journal of Receptors and Signal Transduction*, 24 (1-2):1–52, 2004. doi: 10.1081/RRS-120037896. URL `http://informahealthcare.com/doi/abs/10.1081/RRS-120037896`. PMID: 15344878.

[21] Virginie Lafont, Anthony A. Armstrong, Hiroyasu Ohtaka, Yoshiaki Kiso, L. Mario Amzel, and Ernesto Freire. Compensating enthalpic and entropic changes hinder binding affinity optimization. *Chemical Biology and Drug Design*, 69(6):413–422, 2007. ISSN 1747-0285. doi: 10.1111/j.1747-0285.2007.00519.x. URL `http://dx.doi.org/10.1111/j.1747-0285.2007.00519.x`.

[22] Vijay M. Krishnamurthy, Brooks R. Bohall, Vincent Semetey, and George M. Whitesides. The paradoxical thermodynamic basis for the interaction of ethylene glycol, glycine, and sarcosine chains with bovine carbonic anhydrase ii: an unexpected manifestation of enthalpy/entropy compensation. *Journal of the American Chemical Society*, 128(17):5802–5812, 2006. doi: 10.1021/ja060070r. URL `http://pubs.acs.org/doi/abs/10.1021/ja060070r`. PMID: 16637649.

[23] Derek R. Dee, Yasumi Horimoto, and Rickey Y. Yada. Conserved prosegment residues stabilize a late-stage folding transition state of pepsin independently of ground states. *PLoS ONE*, 9(7):e101339, 07 2014. doi: 10.1371/journal.pone.0101339.

[24] Andrew J. Baldwin, Tuomas P. J. Knowles, Gian Gaetano Tartaglia, Anthony W. Fitzpatrick, Glyn L. Devlin, Sarah Lucy Shammas, Christopher A. Waudby, Maria F. Mossuto, Sarah Meehan, Sally L. Gras, John Christodoulou, Spencer J. Anthony-Cahill, Paul D. Barker, Michele Vendruscolo, and Christopher M. Dobson. Metastability of native proteins and the phenomenon of amyloid formation. *Journal of the American Chemical Society*, 133(36): 14160–14163, 2011. doi: 10.1021/ja2017703. PMID: 21650202.

[25] S. S. Plotkin and J. N. Onuchic. Understanding protein folding with energy landscape theory i: Basic concepts. *Quart. Rev. Biophys.*, 35 (2):111–167, 2002.

[26] R. R. Krug, W. G. Hunter, and R. A. Grieger. Statistical interpretation of enthalpy-entropy compensation. *Nature*, 261: 566–567, 1976. doi: 10.1038/261566a0. URL `http://dx.doi.org/10.1038/261566a0`.

[27] E. Gallicchio, M. M. Kubo, and R. M. Levy. Enthalpy entropy and cavity decomposition of alkane hydration free energies: Numerical results and implications for theories of hydrophobic solvation. *The Journal of Physical Chemistry B*, 104(26):6271–6285, 2000. doi: 10.1021/jp0006274. URL `http://dx.doi.org/10.1021/jp0006274`.

[28] Erik Persson and Bertil Halle. Cell water dynamics on multiple time scales. *Proc. Natl. Acad. Sci. U. S. A.*, 105(17):6266–6271, 2008. doi: 10.1073/pnas.0709585105.

[29] PH Yancey, ME Clark, SC Hand, RD Bowlus, and GN Somero. Living with water stress: evolution of osmolyte systems. *Science*, 217 (4566):1214–1222, 1982. doi: 10.1126/science.7112124. URL `http://www.sciencemag.org/content/217/4566/1214.abstract`.

[30] R. John Ellis and Allen P. Minton. Protein aggregation in crowded environments. *Biol. Chem.*, 387:485–497, 2006.

[31] Huan-Xiang Zhou, Germán Rivas, and Allen P. Minton. Macromolecular crowding and confinement: Biochemical,

biophysical, and potential physiological consequences. *Annu. Rev. Biophys.*, 37(1):375–397, 2008.

[32] Allen P. Minton. The influence of macromolecular crowding and macromolecular confinement on biochemical reactions in physiological media. *J. Biol. Chem.*, 276:10577–10580, 2001.

[33] L. Ellgaard and A. Helenius. Quality control in the endoplasmic reticulum. *Nat. Rev. Mol. Cell. Biol.*, 4:181–191, 2003.

[34] Deepak R. Canchi, Dietmar Paschek, and Angel E. García. Equilibrium study of protein denaturation by urea. *J. Am. Chem. Soc.*, 132(7):2338–2344, 2010. doi: 10.1021/ja909348c. URL http://pubs.acs.org/doi/abs/10.1021/ja909348c.

[35] A. Linhananta, S. Hadizadeh, and S. S. Plotkin. An effective solvent theory connecting the underlying mechanisms of osmolytes and denaturants for protein stability. *Biophys. J.*, 100:459–468, 2011.

[36] Miriam Friedel, Daniel J. Sheeler, and Joan-Emma Shea. Effects of confinement and crowding on the thermodynamics and kinetics of folding of a minimalist beta-barrel protein. *J. Chem. Phys.*, 118(17): 8106–8113, 2003. doi: 10.1063/1.1564048.

[37] M. S. Cheung, D. Klimov, and D. Thirumalai. Molecular crowding enhances native state stability and refolding rates of globular proteins. *Proc. Natl. Acad. Sci. U. S. A.*, 102:4753 – 4758, 2005.

[38] C. Ionescu-Tirgoviste and F. Despa. Biophysical alteration of the secretory track in $\beta$-cells due to molecular overcrowding: the relevance for diabetes. *Integr. Biol.*, 3:173–179, 2011.

[39] A. P. Minton. Implications of macromolecular crowding for protein assembly. *Curr. Opin. Struct. Biol.*, 10:34–39, 2000.

[40] M. Wojciechowski and M. Cieplak. Effects of confinement and crowding on folding of model proteins. *BioSystems*, 94:248–252, 2008.

[41] William L. Jorgensen, Jayaraman Chandrasekhar, Jeffry D. Madura, Roger W. Impey, and Michael L. Klein. Comparison of simple potential functions for simulating liquid water. *The Journal of Chemical Physics*, 79(2):926–935, 1983. doi: 10.1063/1.445869.

[42] K. Vanommeslaeghe, E. Hatcher, C. Acharya, S. Kundu, S. Zhong, J. Shim, E. Darian, O. Guvench, P. Lopes, I. Vorobyov, and A. D. Mackerell. Charmm general force field: A force field for drug-like molecules compatible with the charmm all-atom additive biological force fields. *Journal of Computational Chemistry*, 31(4):671–690, 2010. ISSN 1096-987X. doi: 10.1002/jcc.21367. URL `http://dx.doi.org/10.1002/jcc.21367`.

[43] Christopher J. Cramer and Donald G. Truhlar. Implicit solvation models: equilibria, structure, spectra, and dynamics. *Chemical Reviews*, 99(8):2161–2200, 1999. doi: 10.1021/cr960149m. URL `http://pubs.acs.org/doi/abs/10.1021/cr960149m`.

[44] Jianhan Chen, Charles L Brooks III, and Jana Khandogin. Recent advances in implicit solvent-based methods for biomolecular simulations. *Current Opinion in Structural Biology*, 18(2):140 – 148, 2008. ISSN 0959-440X. doi: http://dx.doi.org/10.1016/j.sbi.2008.01.003. URL `http://www.sciencedirect.com/science/article/pii/S0959440X08000079`. Theory and simulation / Macromolecular assemblages.

[45] Christopher J Cramer and Donald G Truhlar. Implicit solvation models: equilibria, structure, spectra, and dynamics. *Chem. Rev.-Columbus*, 99(8):2161, 1999.

[46] Nathan A. Baker, David Sept, Simpson Joseph, Michael J. Holst, and J. Andrew McCammon. Electrostatics of nanosystems: Application to microtubules and the ribosome. *Proc. Natl. Acad. Sci. USA*, 98(18):10037–10041, 2001. doi: 10.1073/pnas.181342398.

[47] Wonpil Im, Dmitrii Beglov, and Benot Roux. Continuum solvation model: Computation of electrostatic forces from numerical solutions to the poisson-boltzmann equation. *Computer Physics Communications*, 111(13):59 – 75, 1998. ISSN 0010-4655. doi: http://dx.doi.org/10.1016/S0010-4655(98)00016-2. URL `http://www.sciencedirect.com/science/article/pii/S0010465598000162`.

[48] Di Qiu, Peter S. Shenkin, Frank P. Hollinger, and W. Clark Still. The gb/sa continuum model for solvation. a fast analytical method for the calculation of approximate born radii. *The Journal of Physical*

*Chemistry A*, 101(16):3005–3014, 1997. doi: 10.1021/jp961992r. URL http://pubs.acs.org/doi/abs/10.1021/jp961992r.

[49] Jiang Zhu, Yunyu Shi, and Haiyan Liu. Parametrization of a generalized born/solvent-accessible surface area model and applications to the simulation of protein dynamics. *The Journal of Physical Chemistry B*, 106(18):4844–4853, 2002. doi: 10.1021/jp020058v. URL http://pubs.acs.org/doi/abs/10.1021/jp020058v.

[50] Jennifer L. Knight and Charles L. Brooks. Surveying implicit solvent models for estimating small molecule absolute hydration free energies. *Journal of Computational Chemistry*, 32(13):2909–2923, 2011. ISSN 1096-987X. doi: 10.1002/jcc.21876. URL http://dx.doi.org/10.1002/jcc.21876.

[51] D Hawkins, C Cramer, and D Truhlar. Parametrized models of aqueous free energies of solvation based on pairwise descreening of solute atomic charges from a dielectric medium. *Journal of Physical Chemistry*, 100:19824–19839, 1996.

[52] A Onufriev, D Bashford, and D Case. Exploring protein native states and large-scale conformational changes with a modified generalized born model. *PROTEINS: Structure, Function, and Genetics*, 55: 383–394, 2004.

[53] B. Lee and F.M. Richards. The interpretation of protein structures: estimation of static accessibility. *J. Mol. Biol.*, 55:379–400, 1971.

[54] D. Eisenberg and A. D. McLachlan. Solvation energy in protein folding and binding. *Nature*, 319:199–203, 1986.

[55] Chunhu Tan, Yu-Hong Tan, and Ray Luo. Implicit nonpolar solvent models. *The Journal of Physical Chemistry B*, 111(42):12263–12274, 2007. doi: 10.1021/jp073399n. URL http://pubs.acs.org/doi/abs/10.1021/jp073399n.

[56] Jianhan Chen and Charles L. Brooks. Critical importance of length-scale dependence in implicit modeling of hydrophobic interactions. *Journal of the American Chemical Society*, 129(9): 2444–2445, 2007. doi: 10.1021/ja068383+. URL http://pubs.acs.org/doi/abs/10.1021/ja068383.

[57] Allison Ferguson, Zhirong Liu, and Hue Sun Chan. Desolvation barrier effects are a likely contributor to the remarkable diversity in the folding rates of small proteins. *Journal of Molecular Biology*, 389 (3):619 – 636, 2009.

[58] Urs Haberthür and Amedeo Caflisch. Facts: Fast analytical continuum treatment of solvation. *Journal of Computational Chemistry*, 29(5):701–715, 2008. ISSN 1096-987X. doi: 10.1002/jcc.20832. URL http://dx.doi.org/10.1002/jcc.20832.

[59] Jane R. Allison, Katharina Boguslawski, Franca Fraternali, and Wilfred F. van Gunsteren. A refined, efficient mean solvation force model that includes the interior volume contribution. *The Journal of Physical Chemistry B*, 115(15):4547–4557, 2011.

[60] Kunitsugu Soda. Solvent exclusion effect predicted by the scaled particle theory as an important factor of the hydrophobic effect. *J. Phys. Soc. Jpn.*, 62(5):1782–1793, 1993.

[61] A.J. Saunders, P. R. Davis-Searles, D. L. Allen, G. J. Pielak, , and D. A. Erie. Osmolyte-induced changes in protein conformation equilibria. *Biopolymers*, 53:293–307, 2000.

[62] John A. Schellman. Protein stability in mixed solvents: A balance of contact interaction and excluded volume. *Biophys. J.*, 85(1):108–125, 2003.

[63] Takashi Imai, Yuichi Harano, Masahiro Kinoshita, Andriy Kovalenko, and Fumio Hirata. Theoretical analysis on changes in thermodynamic quantities upon protein folding: Essential role of hydration. *J. Chem. Phys.*, 126:225102, 2007.

[64] Jason A. Wagoner and Nathan A. Baker. Assessing implicit models for nonpolar mean solvation forces: The importance of dispersion and volume terms. *Proc. Natl. Acad. Sci. U. S. A.*, 103(22): 8331–8336, 2006. doi: 10.1073/pnas.0600118103.

[65] Tyler Luchko, Sergey Gusarov, Daniel R. Roe, Carlos Simmerling, David A. Case, Jack Tuszynski, and Andriy Kovalenko. Three-dimensional molecular theory of solvation coupled with molecular dynamics in amber. *Journal of Chemical Theory and Computation*, 6(3):607–624, 2010. doi: 10.1021/ct900460m. URL http://pubs.acs.org/doi/abs/10.1021/ct900460m.

[66] David Chandler and Hans C. Andersen. Optimized cluster expansions for classical fluids. ii. theory of molecular liquids. *The Journal of Chemical Physics*, 57(5):1930–1937, 1972. doi: http://dx.doi.org/10.1063/1.1678513. URL `http://scitation.aip.org/content/aip/journal/jcp/57/5/10.1063/1.1678513`.

[67] Andriy Kovalenko and Fumio Hirata. Self-consistent description of a metalwater interface by the kohnsham density functional theory and the three-dimensional reference interaction site model. *The Journal of Chemical Physics*, 110(20):10095–10112, 1999. doi: http://dx.doi.org/10.1063/1.478883. URL `http://scitation.aip.org/content/aip/journal/jcp/110/20/10.1063/1.478883`.

[68] Andreas Vitalis and Rohit V. Pappu. Absinth: A new continuum solvation model for simulations of polypeptides in aqueous solutions. *Journal of Computational Chemistry*, 30(5):673–699, 2009. ISSN 1096-987X. doi: 10.1002/jcc.21005. URL `http://dx.doi.org/10.1002/jcc.21005`.

[69] W. Kohn and L. J. Sham. Self-consistent equations including exchange and correlation effects. *Phys. Rev.*, 140:A1133–A1138, 1965.

[70] T Van Mourik, M B uhl, and M-P Gaigeot. Density functional theory across chemistry, physics and biology. *Philosophical Transactions Series A, Mathematical, Physical, and Engineering Sciences*, 372, 2014.

[71] Kieron Burke, Jan Werschnik, and E. K. U. Gross. Time-dependent density functional theory: Past, present, and future. *The Journal of Chemical Physics*, 123(6):062206, 2005.

[72] N David Mermin. Thermal properties of the inhomogeneous electron gas. *Phys. Rev.*, 137(5A):A1441, 1965.

[73] C. Ebner, W. F. Saam, and D. Stroud. Density-functional theory of simple classical fluids. i. surfaces. *Phys. Rev. A*, 14:2264–2273, Dec 1976. doi: 10.1103/PhysRevA.14.2264.

[74] R Evans. The nature of the liquid-vapour interface and other topics in the statistical mechanics of non-uniform, classical fluids. *Adv. Phys.*, 28(2):143–200, 1979.

[75] David Chandler, John D. McCoy, and Sherwin J. Singer. Density functional theory of nonuniform polyatomic systems. i. general formulation. *J. Chem. Phys.*, 85(10):5971–5976, 1986. doi: 10.1063/1.451510. URL http://link.aip.org/link/?JCP/85/5971/1.

[76] R. Evans. Density functionals in the theory of inhomogeneous fluids. In D. Henderson, editor, *Fundamentals of inhomogeneous fluids*. Dekker, New York, 1992.

[77] T. Biben, J. P. Hansen, and Y. Rosenfeld. Generic density functional for electric double layers in a molecular solvent. *Phys. Rev. E*, 57: R3727–R3730, Apr 1998. doi: 10.1103/PhysRevE.57.R3727.

[78] John M. Richardson. Variational theory of the radial distribution function. *The Journal of Chemical Physics*, 23(12):2304–2308, 1955. doi: http://dx.doi.org/10.1063/1.1740743. URL http://scitation.aip.org/content/aip/journal/jcp/23/12/10.1063/1.1740743.

[79] Jianzhong Wu. Density functional theory for chemical engineering: From capillarity to soft materials. *American Institute of Chemical Engineers Journal*, 52:1169, 2006.

[80] Christopher P. Emborsky, Zhengzheng Feng, Kenneth R. Cox, and Walter G. Chapman. Recent advances in classical density functional theory for associating and polyatomic molecules. *Fluid Phase Equilibria*, 306(1):15 – 30, 2011. ISSN 0378-3812. doi: http://dx.doi.org/10.1016/j.fluid.2011.02.007. URL http://www.sciencedirect.com/science/article/pii/S0378381211000653.

[81] Akira R. Kinjo and Shoji Takada. Effects of macromolecular crowding on protein folding and aggregation studied by density functional theory: statics. *Phys. Rev. E*, 66:031911, Sep 2002. doi: 10.1103/PhysRevE.66.031911. URL http://link.aps.org/doi/10.1103/PhysRevE.66.031911.

[82] Daniel Borgis, Lionel Gendre, and Rosa Ramirez. Molecular density functional theory: Application to solvation and electron-transfer thermodynamics in polar solvents. *The Journal of Physical Chemistry B*, 116(8):2504–2512, 2012. doi: 10.1021/jp210817s.

[83] S. S. Plotkin and J. N. Onuchic. Investigation of routes and funnels in protein folding by free energy functional methods. *Proc. Natl. Acad. Sci. USA*, 97:6509–6514, 2000.

[84] S. S. Plotkin and J. N. Onuchic. Structural and energetic heterogeneity in protein folding i: Theory. *J. Chem. Phys.*, 116(12): 5263–5283, 2002.

[85] Y. Singh, J. P. Stoessel, and P. G. Wolynes. Hard-sphere glass and the density-functional theory of aperiodic crystals. *Phys. Rev. Lett.*, 54:1059–1062, Mar 1985. doi: 10.1103/PhysRevLett.54.1059. URL `http://link.aps.org/doi/10.1103/PhysRevLett.54.1059`.

[86] T. R. Kirkpatrick and P. G. Wolynes. Connections between some kinetic and equilibrium theories of the glass transition. *Phys. Rev. A*, 35(7):3072–3080, 1987.

[87] Randall W. Hall and Peter G. Wolynes. The aperiodic crystal picture and free energy barriers in glasses. *J. Chem. Phys.*, 86(5): 2943–2948, 1987. doi: 10.1063/1.452045. URL `http://link.aip.org/link/?JCP/86/2943/1`.

[88] X. Y. Xia and P. G. Wolynes. Fragilities of liquids predicted from the random first order transition theory of glasses. *Proc. Natl. Acad. Sci. U. S. A.*, 97:2990–2994, 2000.

[89] Xiaoyu Xia and Peter G. Wolynes. Microscopic theory of heterogeneity and nonexponential relaxations in supercooled liquids. *Phys. Rev. Lett.*, 86:5526–5529, Jun 2001. doi: 10.1103/PhysRevLett.86.5526. URL `http://link.aps.org/doi/10.1103/PhysRevLett.86.5526`.

[90] Jacob D Stevenson, Jörg Schmalian, and Peter G Wolynes. The shapes of cooperatively rearranging regions in glass-forming liquids. *Nat. Phys.*, 2(4):268–274, 2006.

[91] Vassiliy Lubchenko and Peter G. Wolynes. Theory of structural glasses and supercooled liquids. *Annu. Rev. Phys. Chem.*, 58(1): 235–266, 2007. doi: 10.1146/annurev.physchem.58.032806.104653.

[92] B. A. Shoemaker, J. Wang, and P. G. Wolynes. Structural correlations in protein folding funnels. *Proc. Nat. Acad. Sci. U. S. A.*, 94:777–782, 1997.

[93] B. A. Shoemaker, J. Wang, and P. G. Wolynes. Exploring structures in protein folding funnels with free energy functionals: the transition state ensemble. *J. Mol. Biol.*, 287:675–694, 1999.

[94] J. J. Portman, S. Takada, and P. G. Wolynes. Microscopic theory of protein folding rates. I. Fine structure of the free energy profile and folding routes from a variational approach. *J. Chem. Phys.*, 114: 5069–5081, 2001.

[95] T. K. Vanderlick, L. E. Scriven, and H. T. Davis. Molecular theories of confined fluids. *The Journal of Chemical Physics*, 90(4): 2422–2436, 1989. doi: http://dx.doi.org/10.1063/1.455985. URL `http://scitation.aip.org/content/aip/journal/jcp/90/4/10.1063/1.455985`.

[96] Yu Chen Shen and David W Oxtoby. Density functional theory of crystal growth: Lennard-jones fluids. *Journal of Chemical Physics*, 104:4233, 1996.

[97] Guillaume Jeanmairet, Maximilien Levesque, Rodolphe Vuilleumier, and Daniel Borgis. Molecular density functional theory of water. *Journal of Physical Chemistry Letters*, 4:619–624, 2013.

[98] Guillaume Jeanmairet, Maximilien Levesque, and Daniel Borgis. Molecular density functional theory of water describing hydrophobicity at short and long length scales. *Journal of Chemical Physics*, 139:154101, 2013.

[99] Sander Pronk, Szilrd Pll, Roland Schulz, Per Larsson, Pr Bjelkmar, Rossen Apostolov, Michael R. Shirts, Jeremy C. Smith, Peter M. Kasson, David van der Spoel, Berk Hess, and Erik Lindahl. Gromacs 4.5: a high-throughput and highly parallel open source molecular simulation toolkit. *Bioinformatics*, 29(7):845–854, 2013. doi: 10.1093/bioinformatics/btt055. URL `http://bioinformatics.oxfordjournals.org/content/29/7/845.abstract`.

[100] Dmitry N. Ivankov and Alexei V. Finkelstein. Prediction of protein folding rates from the amino acid sequence-predicted secondary structure. *Proceedings of the National Academy of Sciences of the United States of America*, 101(24):8942–8944, 2004. doi: 10.1073/pnas.0402659101. URL `http://www.pnas.org/content/101/24/8942.abstract`.

[101] David A. Sivak, John D. Chodera, and Gavin E. Crooks. Using nonequilibrium fluctuation theorems to understand and correct errors in equilibrium and nonequilibrium simulations of discrete langevin dynamics. *Phys. Rev. X*, 3:011007, Jan 2013. doi: 10.1103/PhysRevX.3.011007. URL `http://link.aps.org/doi/10.1103/PhysRevX.3.011007`.

[102] Yaoqi Zhou, Martin Karplus, Keith D. Ball, and R. Stephen Berry. The distance fluctuation criterion for melting: Comparison of square-well and morse potential models for clusters and homopolymers. *The Journal of Chemical Physics*, 116(5), 2002.

[103] R. A. Buckingham. The classical equation of state of gaseous helium, neon and argon. *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 168(933):264–283, 1938. ISSN 0080-4630. doi: 10.1098/rspa.1938.0173.

[104] B. R. Brooks, C. L. Brooks, III, A. D. Mackerell, Jr., L. Nilsson, R. J. Petrella, B. Roux, Y. Won, G. Archontis, C. Bartels, S. Boresch, A. Caflisch, L. Caves, Q. Cui, A. R. Dinner, M. Feig, S. Fischer, J. Gao, M. Hodoscek, W. Im, K. Kuczera, T. Lazaridis, J. Ma, V. Ovchinnikov, E. Paci, R. W. Pastor, C. B. Post, J. Z. Pu, M. Schaefer, B. Tidor, R. M. Venable, H. L. Woodcock, X. Wu, W. Yang, D. M. York, and M. Karplus. CHARMM: The Biomolecular Simulation Program. *JOURNAL OF COMPUTATIONAL CHEMISTRY*, 30(10, Sp. Iss. SI):1545–1614, JUL 30 2009.

[105] M. J. Frisch, G. W. Trucks, H. B. Schlegel, G. E. Scuseria, M. A. Robb, J. R. Cheeseman, J. A. Montgomery, Jr., T. Vreven, K. N. Kudin, J. C. Burant, J. M. Millam, S. S. Iyengar, J. Tomasi, V. Barone, B. Mennucci, M. Cossi, G. Scalmani, N. Rega, G. A. Petersson, H. Nakatsuji, M. Hada, M. Ehara, K. Toyota, R. Fukuda, J. Hasegawa, M. Ishida, T. Nakajima, Y. Honda, O. Kitao, H. Nakai, M. Klene, X. Li, J. E. Knox, H. P. Hratchian, J. B. Cross, V. Bakken, C. Adamo, J. Jaramillo, R. Gomperts, R. E. Stratmann, O. Yazyev, A. J. Austin, R. Cammi, C. Pomelli, J. W. Ochterski, P. Y. Ayala, K. Morokuma, G. A. Voth, P. Salvador, J. J. Dannenberg, V. G. Zakrzewski, S. Dapprich, A. D. Daniels, M. C. Strain, O. Farkas, D. K. Malick, A. D. Rabuck, K. Raghavachari, J. B. Foresman, J. V. Ortiz, Q. Cui, A. G. Baboul, S. Clifford,

J. Cioslowski, B. B. Stefanov, G. Liu, A. Liashenko, P. Piskorz, I. Komaromi, R. L. Martin, D. J. Fox, T. Keith, M. A. Al-Laham, C. Y. Peng, A. Nanayakkara, M. Challacombe, P. M. W. Gill, B. Johnson, W. Chen, M. W. Wong, C. Gonzalez, and J. A. Pople. Gaussian 03, Revision C.02. Gaussian, Inc., Wallingford, CT, 2004.

[106] William L. Jorgensen, David S. Maxwell, and Julian Tirado-Rives. Development and testing of the opls all-atom force field on conformational energetics and properties of organic liquids. *Journal of the American Chemical Society*, 118(45):11225–11236, 1996.

[107] Wendy D. Cornell, Piotr Cieplak, Christopher I. Bayly, Ian R. Gould, Kenneth M. Merz, David M. Ferguson, David C. Spellmeyer, Thomas Fox, James W. Caldwell, and Peter A. Kollman. A second generation force field for the simulation of proteins, nucleic acids, and organic molecules. *Journal of the American Chemical Society*, 117(19):5179–5197, 1995. doi: 10.1021/ja00124a002.

[108] Nathan Schmid, AndreasP. Eichenberger, Alexandra Choutko, Sereina Riniker, Moritz Winger, AlanE. Mark, and WilfredF. van Gunsteren. Definition and testing of the gromos force-field versions 54a7 and 54b7. *European Biophysics Journal*, 40(7):843–856, 2011. ISSN 0175-7571. doi: 10.1007/s00249-011-0700-9.

[109] Valentina Tozzini. Coarse-grained models for proteins. *Current Opinion in Structural Biology*, 15(2):144 – 150, 2005. ISSN 0959-440X. doi: http://dx.doi.org/10.1016/j.sbi.2005.02.005. URL `http://www.sciencedirect.com/science/article/pii/S0959440X05000515`. Theory and simulation/Macromolecular assemblages.

[110] N Gō. Theoretical studies of protein folding. *Annual Review of Biophysics and Bioengineering*, 12(1):183–210, 1983. doi: 10.1146/annurev.bb.12.060183.001151. URL `http://dx.doi.org/10.1146/annurev.bb.12.060183.001151`. PMID: 6347038.

[111] Aram Davtyan, Nicholas P. Schafer, Weihua Zheng, Cecilia Clementi, Peter G. Wolynes, and Garegin A. Papoian. Awsem-md: Protein structure prediction using coarse-grained physical potentials and bioinformatically based local structure biasing. *The Journal of Physical Chemistry B*, 116(29):8494–8503, 2012. doi:

10.1021/jp212541y. URL
`http://pubs.acs.org/doi/abs/10.1021/jp212541y`.

[112] Daniel L. Ensign, Peter M. Kasson, and Vijay S. Pande.
Heterogeneity even at the speed limit of folding: Large-scale
molecular dynamics study of a fast-folding variant of the villin
headpiece. *Journal of Molecular Biology*, 374(3):806 – 816, 2007.
ISSN 0022-2836. doi: http://dx.doi.org/10.1016/j.jmb.2007.09.069.
URL `http://www.sciencedirect.com/science/article/pii/`
`S0022283607012685`.

[113] D.E. Shaw, R.O. Dror, J.K. Salmon, J.P. Grossman, K.M.
Mackenzie, J.A. Bank, C. Young, M.M. Deneroff, B. Batson, K.J.
Bowers, E. Chow, M.P. Eastwood, D.J. Ierardi, J.L. Klepeis, J.S.
Kuskin, R.H. Larson, K. Lindorff-Larsen, P. Maragakis, M.A.
Moraes, S. Piana, Y. Shan, and B. Towles. Millisecond-scale
molecular dynamics simulations on anton. In *High Performance
Computing Networking, Storage and Analysis, Proceedings of the
Conference on*, pages 1–11, Nov 2009. doi: 10.1145/1654059.1654126.

[114] John G. Kirkwood. Statistical mechanics of fluid mixtures. *The
Journal of Chemical Physics*, 3(5):300–313, 1935. doi:
10.1063/1.1749657. URL
`http://link.aip.org/link/?JCP/3/300/1`.

[115] Miguel Jorge, Nuno M. Garrido, Antonio J. Queimada, Ioannis G.
Economou, and Eugenia A. Macedo. Effect of the integration
method on the accuracy and computational efficiency of free energy
calculations using thermodynamic integration. *Journal of Chemical
Theory and Computation*, 6(4):1018–1027, 2010. doi:
10.1021/ct900661c.

[116] Charles H. Bennett. Efficient estimation of free energy differences
from Monte Carlo data. *J Comput Phys*, 22:245–268, 1976.

[117] C. Jarzynski. Nonequilibrium equality for free energy differences.
*Phys. Rev. Lett.*, 78:2690–2693, Apr 1997.

[118] Jeff Gore, Felix Ritort, and Carlos Bustamante. Bias and error in
estimates of equilibrium free-energy differences from nonequilibrium
measurements. *Proceedings of the National Academy of Sciences*, 100
(22):12564–12569, 2003.

134

[119] Marc Souaille and Benoit Roux. Extension to the weighted histogram analysis method: combining umbrella sampling with free energy calculations. *Computer Physics Communications*, 135(1):40 – 57, 2001. ISSN 0010-4655. doi: http://dx.doi.org/10.1016/S0010-4655(00)00215-0. URL `http://www.sciencedirect.com/science/article/pii/S0010465500002150`.

[120] C. B. Anfinsen and H. A. Scheraga. Experimental and theoretical aspects of protein folding. In *Advances in Protein Chemistry*, volume 29, pages 205–301, New York, 1975. Academic Press.

[121] Sophie E. Jackson and Alan R. Fersht. Folding of chymotrypsin inhibitor 2. 1. Evidence for a two state transition. *Biochemistry*, 30: 10428–10435, 1991.

[122] G. I. Makhatadze and P. L. Privalov. Energetics of protein structure. *Adv. Protein Chem.*, 47:307–425, 1995.

[123] P.L. Privalov, E.I. Tiktopulo, S.Yu. Venyaminov, Yu.V. Griko, G.I. Makhatadze, and N.N. Khechinashvili. Heat capacity and conformation of proteins in the denatured state. *Journal of Molecular Biology*, 205(4):737 – 750, 1989. doi: http://dx.doi.org/10.1016/0022-2836(89)90318-5.

[124] K. A. Dill, S. Bromberg, K. Yue, K. M. Fiebig, D. P. Yee, P. D. Thomas, and H. S. Chan. Principles of protein folding—A perspective from simple exact models. *Protein Science*, 4:561–602, 1995.

[125] S. S. Plotkin and J. N. Onuchic. Understanding protein folding with energy landscape theory ii: Quantitative aspects. *Quart. Rev. Biophys.*, 35(3):205–286, 2002.

[126] Christopher D. Snow, Houbi Nguyen, Vijay S. Pande, and Martin Gruebele. Absolute comparison of simulated and experimental protein-folding dynamics. *Nature*, 420(6):102–106, 2002.

[127] Thomas R Weikl, Matteo Palassini, and Ken A Dill. Cooperativity in two-state protein folding kinetics. *Protein Science*, 13(3):822–829, 2004.

[128] M. R. Ejtehadi, S. P. Avall, and S. S. Plotkin. Three-body interactions improve the prediction of rate and mechanism in protein folding models. *Proc. Natl. Acad. Sci. USA*, 101(42):15088–15093, 2004.

[129] Elizabeth Rhoades, Mati Cohen, Benjamin Schuler, and Gilad Haran. Two-state folding observed in individual protein molecules. *Journal of the American Chemical Society*, 126(45):14686–14687, 2004.

[130] Beatriz Ibarra-Molero and Jose Manuel Sanchez-Ruiz. Statistical differential scanning calorimetry: Probing protein folding-unfolding ensembles. In Victor Muñoz, editor, *Protein folding, misfolding and aggregation: Classical themes and novel approaches*, RSC Biomolecular Sciences, pages 85–103. Royal Society of Chemistry, Cambridge, UK, 2008.

[131] Douglas Poland. Empirical protein partition functions. *The Journal of Physical Chemistry B*, 116(23):6683–6693, 2012. doi: 10.1021/jp211794u.

[132] Peter L. Privalov. *Microcalorimetry of Macromolecules: The Physical Basis of Biological Structures*. John Wiley & Sons, Inc., New York, 2012.

[133] Lucas N. R. Wafer, Werner W. Streicher, and George I. Makhatadze. Thermodynamics of the trp-cage miniprotein unfolding in urea. *Proteins: Structure, Function, and Bioinformatics*, 78(6):1376–1381, 2010. ISSN 1097-0134. doi: 10.1002/prot.22681.

[134] Andrew T. Fenley, Hari S. Muddana, and Michael K. Gilson. Entropyenthalpy transduction caused by conformational shifts can obscure the forces driving proteinligand binding. *Proceedings of the National Academy of Sciences*, 109(49):20006–20011, 2012. doi: 10.1073/pnas.1213180109. URL http://www.pnas.org/content/109/49/20006.abstract.

[135] Wayne J. Becktel and John A. Schellman. Protein stability curves. *Biopolymers*, 26:1859–1877, 1987.

[136] Patrick L Wintrode, George I Makhatadze, and Peter L Privalov. Thermodynamics of ubiquitin unfolding. *Proteins: Structure, Function, and Bioinformatics*, 18(3):246–253, 1994.

[137] GI Makhatadze and PL Privalov. Heat capacity of proteins: I. partial molar heat capacity of individual amino acid residues in aqueous solution: hydration effect. *J. Mol. Biol.*, 213(2):375–384, 1990.

[138] Mark E. Zweifel and Doug Barrick. Relationships between the temperature dependence of solvent denaturation and the denaturant dependence of protein stability curves. *Biophysical Chemistry*, 101-102(0):221 – 237, 2002. doi: http://dx.doi.org/10.1016/S0301-4622(02)00181-3. URL `http://www.sciencedirect.com/science/article/pii/S0301462202001813`.

[139] Eric M. Nicholson and J. Martin Scholtz. Conformational stability of the escherichia coli hpr protein: test of the linear extrapolation method and a thermodynamic characterization of cold denaturation. *Biochemistry*, 35(35):11369–11378, 1996. doi: 10.1021/bi960863y.

[140] James U. Bowie and Robert T. Sauer. Equilibrium dissociation and unfolding of the arc repressor dimer. *Biochemistry*, 28(18): 7139–7143, 1989. doi: 10.1021/bi00444a001. URL `http://pubs.acs.org/doi/abs/10.1021/bi00444a001`.

[141] Fan-Guo Meng, Yuan-Kai Hong, Hua-Wei He, Arkadii E. Lyubarev, Boris I. Kurganov, Yong-Bin Yan, and Hai-Meng Zhou. Osmophobic effect of glycerol on irreversible thermal denaturation of rabbit creatine kinase. *Biophysical Journal*, 87(4):2247 – 2254, 2004. ISSN 0006-3495. doi: http://dx.doi.org/10.1529/biophysj.104.044784. URL `http://www.sciencedirect.com/science/article/pii/S0006349504737029`.

[142] Nataša Poklar, Nina Petrovčič, Miha Oblak, and Gorazd Vesnaver. Thermodynamic stability of ribonuclease a in alkylurea solutions and preferential solvation changes accompanying its thermal denaturation: A calorimetric and spectroscopic study. *Protein Science*, 8(4):832–840, 1999. ISSN 1469-896X. doi: 10.1110/ps.8.4.832. URL `http://dx.doi.org/10.1110/ps.8.4.832`.

[143] Alan Fersht. *Structure and Mechanism in Protein Science: A Guide to Enzyme Catalysis and Protein Folding*. W H Freeman & Co, 3rd edition, 1998.

[144] Edward P. O'Brien, Guy Ziv, Gilad Haran, Bernard R. Brooks, and D. Thirumalai. Effects of denaturants and osmolytes on proteins are

accurately predicted by the molecular transfer model. *Proc. Natl. Acad. Sci. U. S. A.*, 105(36):13403–13408, 2008. doi: 10.1073/pnas.0802113105.

[145] M. Auton and D. W. Bolen. Predicting the energetics of osmolyte-induced protein folding/unfolding. *Proc. Natl. Acad. Sci. U. S. A.*, 102(42):15065–15068, 2005.

[146] Jai K. Kaushik and Rajiv Bhat. Why is trehalose an exceptional protein stabilizer?: An analysis of the thermal stability of proteins in the presence of the compatible osmolyte trehalose. *Journal of Biological Chemistry*, 278(29):26458–26465, 2003. doi: 10.1074/jbc.M300815200. URL `http://www.jbc.org/content/278/29/26458.abstract`.

[147] Rajendrakumar Singh, Inamul Haque, and Faizan Ahmad. Counteracting osmolyte trimethylamine n-oxide destabilizes proteins at ph below its pka: Measurements of thermodynamic parameters of proteins in the presence and absence of trimethylamine n-oxide. *Journal of Biological Chemistry*, 280(12):11035–11042, 2005. doi: 10.1074/jbc.M410716200. URL `http://www.jbc.org/content/280/12/11035.abstract`.

[148] Fabrizio Chiti, Nico A. J. van Nuland, Niccoló Taddei, Francesca Magherini, Massimo Stefani, Giampietro Ramponi, and Christopher M. Dobson. Conformational stability of muscle acylphosphatase: the role of temperature, denaturant concentration, and ph. *Biochemistry*, 37(5):1447–1455, 1998. doi: 10.1021/bi971692f. URL `http://pubs.acs.org/doi/abs/10.1021/bi971692f`.

[149] Vishwas R. Agashe and Jayant B. Udgaonkar. Thermodynamics of denaturation of barstar: Evidence for cold denaturation and evaluation of the interaction with guanidine hydrochloride. *Biochemistry*, 34(10):3286–3299, 1995. doi: 10.1021/bi00010a019. URL `http://pubs.acs.org/doi/abs/10.1021/bi00010a019`.

[150] S. Padmanabhan, D. V. Laurents, A. M. Fernández, M. Elias-Arnanz, J. Ruiz-Sanz, P. L. Mateo, M. Rico, and V. V. Filimonov. Thermodynamic analysis of the structural stability of phage 434 cro protein. *Biochemistry*, 38(47):15536–15547, 1999. doi:

10.1021/bi991757. URL
http://pubs.acs.org/doi/abs/10.1021/bi991757.

[151] T. Kaminyama, Y. Sadahide, Y. Nogusa, and K. Gekko. *Biochim. Biophys. Acta*, 1434:44–57, 1999.

[152] James W. Bryson, John R. Desjarlais, Tracy M. Handel, and William F. Degrado. From coiled coils to small globular proteins: Design of a native-like three-helix bundle. *Protein Science*, 7(6): 1404–1414, 1998. ISSN 1469-896X. doi: 10.1002/pro.5560070617. URL http://dx.doi.org/10.1002/pro.5560070617.

[153] F. Catanzano, A. Gambuti, G. Graziano, and G. Barone. Interaction with d-glucose and thermal denaturation of yeast hexokinase b: A dsc study. *Journal of Biochemistry*, 121(3):568–577, 1997. URL http://jb.oxfordjournals.org/content/121/3/568.abstract.

[154] Nisar Ahmad, V. R. Srinivas, G. Bhanuprakash Reddy, and Avadhesha Surolia. Thermodynamic characterization of the conformational stability of the homodimeric protein, pea lectin. *Biochemistry*, 37(47):16765–16772, 1998. doi: 10.1021/bi9811720. URL http://pubs.acs.org/doi/abs/10.1021/bi9811720.

[155] E. L. Kovrigin and S. A. Potekhin. *Biofizika*, 41:1201–1206, 1996.

[156] S. Knapp, Rudolf Ladenstein, and Erwin A. Galinski. Extrinsic protein stabilization by the naturally occurring osmolytes $\beta$-hydroxyectoine and betaine. *Extremophiles*, 3(3):191–198, 1999. ISSN 1431-0651. doi: 10.1007/s007920050116. URL http://dx.doi.org/10.1007/s007920050116.

[157] Marcelo M. Santoro, Yufeng Liu, Saber M. A. Khan, Li Xiang Hou, and D. W. Bolen. Increased thermal stability of proteins in the presence of naturally occurring osmolytes. *Biochemistry*, 31(23): 5278–5283, 1992. doi: 10.1021/bi00138a006. URL http://pubs.acs.org/doi/abs/10.1021/bi00138a006.

[158] Su Xu, Sanbo Qin, and Xian-Ming Pan. Thermal and conformational stability of ssh10b protein from archaeon sulfolobus shibattae. *Biochem J*, 382(Pt 2):433–440, 2004.

[159] Yaqiang Wang, Mohona Sarkar, Austin E. Smith, Alexander S. Krois, and Gary J. Pielak. Macromolecular crowding and protein

stability. *Journal of the American Chemical Society*, 134(40):
16614–16618, 2012. doi: 10.1021/ja305300m. URL
`http://pubs.acs.org/doi/abs/10.1021/ja305300m`.

[160] George I. Makhatadze, Marin M. Lopez, John M. Richardson III,
and Susan T. Thomas. Anion binding to the ubiquitin molecule.
*Protein Science*, 7(3):689–697, 1998. ISSN 1469-896X. doi:
10.1002/pro.5560070318. URL
`http://dx.doi.org/10.1002/pro.5560070318`.

[161] Michael Senske, Lisa Trk, Benjamin Born, Martina Havenith,
Christian Herrmann, and Simon Ebbinghaus. Protein stabilization
by macromolecular crowding through enthalpy rather than entropy.
*Journal of the American Chemical Society*, 136(25):9036–9041, 2014.

[162] A Ben-Naim. Hydrophobic interaction and structural changes in the
solvent. *Biopolymers*, 14(7):1337–1355, 1975.

[163] Ernest Grunwald. Thermodynamic properties, propensity laws, and
solvent models in solutions in self-associating solvents. application to
aqueous alcohol solutions. *Journal of the American Chemical
Society*, 106(19):5414–5420, 1984. doi: 10.1021/ja00331a006.

[164] HsiangAi Yu and Martin Karplus. A thermodynamic analysis of
solvation. *The Journal of Chemical Physics*, 89(4):2366–2379, 1988.
doi: http://dx.doi.org/10.1063/1.455080.

[165] B Lee. Enthalpy-entropy compensation in the thermodynamics of
hydrophobicity. *Biophysical chemistry*, 51(2):271–278, 1994.

[166] Eric A. Mills and Steven S. Plotkin. Density functional theory for
protein transfer free energy. *The Journal of Physical Chemistry B*,
117:13278–13290, 2013. doi: 10.1021/jp403600q.

[167] Daniel Borgis, Lionel Gendre, and Rosa Ramirez. Molecular density
functional theory: Application to solvation and electron-transfer
thermodynamics in polar solvents. *J. Phys. Chem. B*, 116(8):
2504–2512, 2012. doi: 10.1021/jp210817s.

[168] Charles. Tanford. Isothermal unfolding of globular proteins in
aqueous urea solutions. *J. Am. Chem. Soc.*, 86(10):2050–2059, 1964.
doi: 10.1021/ja01064a028.

[169] Jonathan E. Kohn, Ian S. Millett, Jaby Jacob, Bojan Zagrovic, Thomas M. Dillon, Nikolina Cingel, Robin S. Dothager, Soenke Seifert, P. Thiyagarajan, Tobin R. Sosnick, M. Zahid Hasan, Vijay S. Pande, Ingo Ruczinski, Sebastian Doniach, and Kevin W. Plaxco. Random-coil behavior and the dimensions of chemically unfolded proteins. *Proc. Natl. Acad. Sci. U. S. A.*, 101(34):12491–12496, 2004.

[170] Jeffrey K. Noel, Paul C. Whitford, and José N. Onuchic. The shadow map: A general contact definition for capturing the dynamics of biomolecular folding and function. *J. Phys. Chem. B*, 116(29): 8692–8702, 2012. doi: 10.1021/jp300852d. URL `http://pubs.acs.org/doi/abs/10.1021/jp300852d`.

[171] Jeffry K. Noel, Paul C. Whitford, Karissa Y. Sanbonmatsu, and José N. ONuchi. Smog@ctbp: simplified deployment of structure based models in gromacs. *Nucleic Acids Res.*, 38:W657–661, 2010.

[172] C. Clementi, H. Nymeyer, and J. N. Onuchic. Topological and energetic factors: what determines the structural details of the transition state ensemble and "en-route" intermediates for protein folding? An investigation for small globular proteins. *J. Mol. Biol.*, 298:937–953, 2000.

[173] G. J. Kleywegt and T. A. Jones. Detection, delineation, measurement and display of cavities in macromolecular structures. *Acta Crystallogr., Sect. D: Biol. Crystallogr.*, 50(2):178–185, Mar 1994. doi: 10.1107/S0907444993011333. URL `http://dx.doi.org/10.1107/S0907444993011333`.

[174] Kristine M. Kast, Jürgen Brickmann, Stefan M. Kast, and R. Stephen Berry. Binary phases of aliphatic N-oxides and water: Force field development and molecular dynamics simulation. *J. Phys. Chem. A*, 107(27):5342–5351, 2003. doi: 10.1021/jp027336a. URL `http://pubs.acs.org/doi/abs/10.1021/jp027336a`.

[175] John A. Schellman. The thermodynamics of solvent exchange. *Biopolymers*, 34(8):1015–1026, 1994. doi: 10.1002/bip.360340805.

[176] J. K. Myers, C. N. Pace, and J. M. Scholtz. Denaturant $m$ values and heat capacity changes: Relation to changes in accessible surface areas of protein unfolding. *Protein Sci.*, 4:2138–2148, 1995.

[177] Darwin O. V. Alonso and Ken A. Dill. Solvent denaturation and stabilization of globular proteins. *Biochemistry*, 30:5974–5985, 1991.

[178] T. L. Hill. *An Introduction to Statistical Thermodynamics*. Courier Dover Publications, New York, 1960.

[179] Norman F Carnahan and Kenneth E Starling. Equation of state for nonattracting rigid spheres. *J. Chem. Phys.*, 51:635, 1969.

[180] Michael Plischke and Birger Bergersen. *Statistical Physics*. World Scientific, Singapore, 3rd edition, 2006.

[181] R. O. Jones and O. Gunnarsson. The density functional formalism, its applications and prospects. *Rev. Mod. Phys.*, 61:689–746, Jul 1989. doi: 10.1103/RevModPhys.61.689. URL http://link.aps.org/doi/10.1103/RevModPhys.61.689.

[182] Joseph T. Slusher. Accurate estimates of infinite-dilution chemical potentials of small hydrocarbons in water via molecular dynamics simulation. *J. Phys. Chem. B*, 103(29):6075–6079, 1999. doi: 10.1021/jp990709w.

[183] J. T. Wescott, L. R. Fisher, and S. Hanna. Use of thermodynamic integration to calculate the hydration free energies of n-alkanes. *J. Chem. Phys.*, 116(6):2361–2369, 2002. doi: 10.1063/1.1431588.

[184] Michael R. Shirts, Jed W. Pitera, William C. Swope, and Vijay S. Pande. Extremely precise free energy calculations of amino acid side chain analogs: Comparison of common molecular mechanics force fields for proteins. *J. Chem. Phys.*, 119(11):5740–5761, 2003. doi: 10.1063/1.1587119.

[185] A A Zamyatnin. Amino acid, peptide, and protein volume in solution. *Annu. Rev. Biophys. Bioeng.*, 13(1):145–165, 1984. doi: 10.1146/annurev.bb.13.060184.001045.

[186] Matthew Auton and D. Wayne Bolen. Additive transfer free energies of the peptide backbone unit that are independent of the model compound and the choice of concentration scale. *Biochemistry*, 43 (5):1329–1342, 2004. doi: 10.1021/bi035908r.

[187] Rosa Ramirez and Daniel Borgis. Density functional theory of solvation and its relation to implicit solvent models. *J. Phys. Chem. B*, 109(14):6754–6763, 2005. doi: 10.1021/jp045453v.

[188] Srabani Roy and Biman Bagchi. Solvation dynamics in liquid water. a novel interplay between librational and diffusive modes. *J. Chem. Phys.*, 99(12):9938–9943, 1993. doi: 10.1063/1.465392.

[189] A. Chandra and B. Bagchi. Molecular theory of solvation and solvation dynamics in a binary dipolar liquid. *J. Chem. Phys.*, 94 (12):8367–8377, 1991. doi: 10.1063/1.460068.

[190] Akira Yoshimori, Tyler J. F. Day, and G. N. Patey. An investigation of dynamical density functional theory for solvation in simple mixtures. *J. Chem. Phys.*, 108(15):6378–6386, 1998. doi: 10.1063/1.476044.

[191] B. Götzelmann, R. Evans, and S. Dietrich. Depletion forces in fluids. *Phys. Rev. E*, 57:6785–6800, Jun 1998. doi: 10.1103/PhysRevE.57.6785. URL `http://link.aps.org/doi/10.1103/PhysRevE.57.6785`.

[192] Fumio Oosawa and Sho Asakura. Surface tension of high-polymer solutions. *J. Chem. Phys.*, 22(7):1255–1255, 1954. doi: 10.1063/1.1740346. URL `http://link.aip.org/link/?JCP/22/1255/1`.

[193] Phil Attard. Spherically inhomogeneous fluids. ii. hard-sphere solute in a hard-sphere solvent. *J. Chem. Phys.*, 91(5):3083–3089, 1989. doi: 10.1063/1.456931. URL `http://link.aip.org/link/?JCP/91/3083/1`.

[194] Phil Attard and G. N. Patey. Hypernetted-chain closure with bridge diagrams. asymmetric hard sphere mixtures. *J. Chem. Phys.*, 92(8): 4970–4982, 1990. doi: 10.1063/1.458556. URL `http://link.aip.org/link/?JCP/92/4970/1`.

[195] H. N. W. Lekkerkerker, W. C. K. Poon, P. N. Pusey, A. Stroobants, and P. B. Warren. Phase-behavior of colloid plus polymer mixtures. *Europhys. Lett.*, 20(6):559–564, 1992.

[196] Ronald Dickman, Phil Attard, and Veronika Simonian. Entropic forces in binary hard sphere mixtures: Theory and simulation. *J. Chem. Phys.*, 107(1):205–213, 1997. doi: 10.1063/1.474367. URL `http://link.aip.org/link/?JCP/107/205/1`.

[197] Jian-Min Yuan, Chia-Lin Chyan, Huan-Xiang Zhou, Tse-Yu Chung, Haibo Peng, Guanghui Ping, and Guoliang Yang. The effects of macromolecular crowding on the mechanical stability of protein molecules. *Protein Sci.*, 17(12):2156–2166, 2008.

[198] Rosa Ramirez, Michel Mareschal, and Daniel Borgis. Direct correlation functions and the density functional theory of polar solvents. *Chemical Physics*, 319:261–272, 2005.

[199] Thomas F. Coleman and Yuying Li. An interior trust region approach for nonlinear minimization subject to bounds. *SIAM Journal on Optimization*, 6(2):418–445, 1996.

[200] B. Lee and F.M. Richards. The interpretation of protein structures: Estimation of static accessibility. *Journal of Molecular Biology*, 55 (3):379 – IN4, 1971. ISSN 0022-2836.

[201] F Eisenhaber, P Lijnzaad, P Argos, C Sander, and M Scharf. The double cubic lattice method: Efficient approaches to numerical integration of surface area and volume and to dot surface contouring of molecular assemblies. *Journal of Computational Chemistry*, 16: 273–284, 1995.

[202] R E Nettleton and M S Green. Expression in terms of molecular distribution functions for the entropy density in an infinite system. *Journal of Chemical Physics*, 29:1365, 1958.

[203] A Singer. Maximum entropy formulation of the kirkwood superposition approximation. *Journal of Chemical Physics*, 121:3657, 2004.

[204] Assessing the performance of mm/pbsa and mm/gbsa methods. 3. the impact of force fields and ligand charge models. *Journal of Physical Chemistry B*, 117:8408–8421, 2013.

[205] V Z Spassov, L Yan, and S Szalma. Introducing an implicit membrane in generalized born/solvent accessibility continuum solvent models. *Journal of Physical Chemistry B*, 106:8726–8738, 2002.

[206] Kyle A. Beauchamp, Yu-Shan Lin, Rhiju Das, and Vijay S. Pande. Are protein force fields getting better? a systematic benchmark on 524 diverse nmr measurements. *Journal of Chemical Theory and Computation*, 8(4):1409–1414, 2012.

[207] P Larsson and E Lindahl. A high-performance parallel-generalized born implementation enabled by tabulated interaction rescaling. *Journal of Computational Chemistry*, 31:2593–2600, 2010.

[208] B Hess, C Kutzner, D van der Spoel, and E Lindahl. Gromacs 4: Algorithms for highly efficient, load-balanced, and scalable molecular simulation. *Journal of Chemical Theory and Computation*, 4: 435–447, 2008.

# Appendix A

# Proof of Theorems in DFT

In this appendix we prove the Hohenberg-Kohn theorem. To start we define

$$\text{Tr}... = \sum_{N=0}^{\infty} \frac{1}{h^{3N} N!} \int \int ...dr^{3N} dp^{3N} \tag{A.1}$$

and

$$\Xi = \text{Tr} e^{-\beta(\mathcal{H} - N\mu)} \tag{A.2}$$

With these defined we can begin with the following:

$$\Omega[f] = \text{Tr} f(\mathcal{H} - N\mu + k_B T \ln f) \tag{A.3}$$

**Lemma A.0.1** *Let $f_0$ be the equilibrium density. Then for any other density $f$, $\Omega[f] > \Omega[f_0]$.*

PROOF: The equilibrium density can be expressed as

$$f_o = \frac{e^{-\beta(\mathcal{H} - N\mu)}}{\Xi} \tag{A.4}$$

So then

$$\Omega[f_0] = \text{Tr}(\mathcal{H} - N\mu - k_B T \ln \Xi - \mathcal{H} + N\mu) \tag{A.5}$$

$$= -k_B T \ln \Xi \tag{A.6}$$

We can then write

$$\Omega[f] - \Omega[f_0] = \text{Tr} f(\mathcal{H} - N\mu + k_B T \ln f + k_B T \ln \Xi) \tag{A.7}$$

Since $\mathcal{H} - N\mu + k_B T \ln \Xi = -k_B T \ln f_0$,

$$\Omega[f] - \Omega[f_0] = \text{Tr} k_B T (f \ln f - f \ln f_0) \tag{A.8}$$

Since $\text{Tr} f = 1$, we can write

$$\Omega[f] - \Omega[f_0] = \text{Tr} k_B T (f \ln f - f \ln f_0 + f_0 - f) = \text{Tr} f_0 k_B T \left(\frac{f}{f_0} \ln \frac{f}{f_0} + 1 - \frac{f}{f_0}\right) \tag{A.9}$$

Since $x \ln x \geq x - 1$ for all $x$, where the equality is only true when $x = 1$, $\Omega[f] - \Omega[f_0] > 0$, and hence Lemma A.0.1 is proved.

We can now show the following theorem:

**Theorem A.0.2** *The functional* $\mathcal{F}[\phi] = \text{Tr} f_0 (\mathcal{H} - N\mu + k_B T \ln f_0)$ *is a unique functional of* $\phi$, *the single particle density.*

To prove this we assume the opposite: that there are two densities $f_0$ and $f_0'$, both of which are equilibrium densities for the given hamiltonian. Then, by Lemma A.0.1, we have

$$\mathcal{F}[\phi'] > \mathcal{F}[\phi] \tag{A.10}$$

But, the choice of $f_0'$ and $f_0$ was arbitrary; the prime could have been on either function. So equation A.10 cannot be true. Thus the density is a unique functional of the hamiltonian and hence of the external potential, proving theorem A.

Finally, we have

**Theorem A.0.3** *The minimum value of $\mathcal{F}[\phi]$ is the free energy of the system and occurs when $\phi = \phi_0$, the equilibrium density*

This follows because if $\mathcal{F}[\phi]$ at the equilibrium density is $\Omega[f_0]$ in equation A.6. Theorem shows that there cannot be another $\phi$ that satisfies this condition, and lemma A.0.1 shows that $\mathcal{F}$ is larger for any other $\phi$. Thus theorem A.0.3 is shown.

# Appendix B

# Simulation Parameters

In this appendix we list the van der Waals parameters used in simulations throughout this thesis. The vdW interaction between two atoms $i$ and $j$ is given by

$$V_{ij} = 4\sqrt{\epsilon_i \epsilon_j} \left( \left( \frac{\sigma_i + \sigma_j}{2r_{ij}} \right)^1 2 - \left( \frac{\sigma_i + \sigma_j}{2r_{ij}} \right)^6 \right) \qquad \text{(B.1)}$$

were $\sigma_i$ and $\epsilon_i$ are the van der Waals parameters for atom $i$ (and likewise for atom $j$) and $r_{ij}$ is the distance between atoms $i$ and $j$. The following atom types are from the CHARMM forcefield.

Table B.1: van der Waals parameters for atoms used in simulations in this thesis.

| Atom type | $\sigma$ (nm) | $\epsilon$ (kJ/mol) |
|:---:|:---:|:---:|
| CA | 0.355005321205 | 0.29288 |
| CC | 0.356359487256 | 0.29288 |
| CP1 | 0.405358916754 | 0.08368 |
| CP2 | 0.387540942391 | 0.23012 |
| CP3 | 0.387540942391 | 0.23012 |
| CT1 | 0.405358916754 | 0.08368 |
| CT2 | 0.387540942391 | 0.23012 |
| CT3 | 0.367050271874 | 0.33472 |
| H | 0.0400013524445 | 0.192464 |
| HA | 0.235197261589 | 0.092048 |
| HB | 0.235197261589 | 0.092048 |
| HP | 0.242003727796 | 0.12552 |
| HR1 | 0.160361769265 | 0.192464 |
| HR2 | 0.12472582054 | 0.192464 |
| HR3 | 0.261567863646 | 0.0326352 |
| NC2 | 0.329632525712 | 0.8368 |
| NH1 | 0.329632525712 | 0.8368 |
| NH2 | 0.329632525712 | 0.8368 |
| NH3 | 0.329632525712 | 0.8368 |
| NPH | 0.329632525712 | 0.8368 |
| NR1 | 0.329632525712 | 0.8368 |
| NR2 | 0.329632525712 | 0.8368 |
| NR3 | 0.329632525712 | 0.8368 |
| O | 0.302905564168 | 0.50208 |
| OC | 0.302905564168 | 0.50208 |
| S | 0.356359487256 | 1.8828 |
| OWT3 | 0.315058 | 0.636386 |
| HWT3 | 0.0 | 0.0 |