# Supporting information for:

# Constant-pH Molecular Dynamics Simulations for Large Biomolecular Systems

Brian K. Radak,[†] Christophe Chipot,[‡,¶] Donghyuk Suh,[§] Sunhwan Jo,[†,@] Wei Jiang,[†] James C. Phillips,[‖] Klaus Schulten,[‖,¶] and Benoît Roux[*,⊥,§,#]

†Leadership Computing Facility, Argonne National Laboratory, Argonne, IL 60439-8643, USA

‡Laboratoire International Associé Centre National de la Recherche Scientifique et University of Illinois at Urbana-Champaign, Unité Mixte de Recherche n° 7565, Université de Lorraine, Université de Lorraine, B.P. 70239, 54506 Vandœuvre-lès-Nancy cedex France

¶Department of Physics, University of Illinois at Urbana-Champaign, Urbana, IL, 61801-2325, USA

§Department of Chemistry, University of Chicago, Chicago, IL, 60637-1454, USA

‖Theoretical and Computational Biophysics Group, Beckman Institute for Advanced Science and Technology, University of Illinois at Urbana-Champaign, Urbana, IL, 61801-2325, USA

⊥Department of Biochemistry and Molecular Biology, University of Chicago, Chicago, IL, 60637-1454, USA

#Center for Nanoscale Materials, Argonne National Laboratory, Argonne, IL 60439-8643, USA

@Current Address: Silcs Bio, LLC, Baltimore, MD, 21201-1193, USA

E-mail: roux@uchicago.edu

# 1 Detailed Description of Force Field Modifications

## 1.1 Equivalence of Standard and Dummy Atom Force Fields

In the main text it was indicated that all noninteracting protonation sites remain in the system as "dummy" atoms and that this modification affects the dynamics but not the thermodynamics of the system. Here this assertion is verified in detail.

Consider a classical system with Hamiltonian $H(\boldsymbol{p}, \boldsymbol{r}) = K(\boldsymbol{p}) + U(\boldsymbol{r})$, where the kinetic and potential energies, $K$ and $U$, are strictly functions of the momenta, $\boldsymbol{p}$ and coordinates, $\boldsymbol{r}$. The system has the canonical distribution

$$\rho(\boldsymbol{p}, \boldsymbol{r}) = \frac{e^{-\beta H(\boldsymbol{p}, \boldsymbol{r})}}{Q}, \tag{1}$$

where $\beta \equiv 1/(k_{\mathrm{B}}T)$ is the inverse temperature and $Q$ is the partition function:

$$Q = \int d\boldsymbol{p} d\boldsymbol{r} \, e^{-\beta H(\boldsymbol{p}, \boldsymbol{r})}. \tag{2}$$

Additional, non-physical, "dummy" particles can be added to the system, with coordinates and momenta $\boldsymbol{\xi}$ and $\boldsymbol{p}_\xi$, respectively, which clearly modify the distribution. In general, the dummy particles will experience a potential $U_{\mathrm{D}}(\boldsymbol{r}, \boldsymbol{\xi})$, which depends on the coordinates of all of the particles in the system. The modified distribution is thus

$$\rho(\boldsymbol{p}, \boldsymbol{r}, \boldsymbol{p}_\xi, \boldsymbol{\xi}) = \frac{e^{-\beta H(\boldsymbol{p}, \boldsymbol{r})} e^{-\beta\left[K(\boldsymbol{p}_\xi) + U_{\mathrm{D}}(\boldsymbol{r}, \boldsymbol{\xi})\right]}}{Q_{\mathrm{D}}}, \tag{3}$$

where $Q_{\mathrm{D}}$ and $Q$ are similarly defined. This distribution is clearly *not the same* as $\rho(\boldsymbol{p}, \boldsymbol{r})$. Multiplying by $Q/Q$ it can be seen that:

$$\rho(\boldsymbol{p}, \boldsymbol{r}, \boldsymbol{p}_\xi, \boldsymbol{\xi}) = \frac{Q}{Q_{\mathrm{D}}} e^{-\beta\left[K(\boldsymbol{p}_\xi) + U_{\mathrm{D}}(\boldsymbol{r}, \boldsymbol{\xi})\right]} \rho(\boldsymbol{p}, \boldsymbol{r}). \tag{4}$$

Despite the inequivalence of the distributions, the averages of mechanical observables are the same for both distributions when the potential $U_{\mathrm{D}}(\boldsymbol{r}, \boldsymbol{\xi})$ is defined in a certain way. In the original

system the average of some mechanical observable is

$$\langle A \rangle = \int d\boldsymbol{p} d\boldsymbol{r} \, A(\boldsymbol{p}, \boldsymbol{r}) \rho(\boldsymbol{p}, \boldsymbol{r}) \tag{5}$$

and the modified average is

$$\langle A \rangle_{\mathrm{D}} = \frac{Q}{Q_{\mathrm{D}}} \int d\boldsymbol{p}_\xi \, e^{-\beta K(\boldsymbol{p}_\xi)} \int d\boldsymbol{p} d\boldsymbol{r} d\boldsymbol{\xi} \, A(\boldsymbol{p}, \boldsymbol{r}) \rho(\boldsymbol{p}, \boldsymbol{r}) e^{-\beta U_{\mathrm{D}}(\boldsymbol{r}, \boldsymbol{\xi})}. \tag{6}$$

Again, these equations are clearly not the same. However, it is now assumed that there is some transformation from $(\boldsymbol{r}, \boldsymbol{\xi})$ to some other dummy coordinate frame $(\boldsymbol{r}, \boldsymbol{\xi}')$ such that $U_{\mathrm{D}}(\boldsymbol{r}, \boldsymbol{\xi}) = U_{\mathrm{D}}(\boldsymbol{\xi}')$. This is true, for example, if the interactions are only governed by three non-dummy particles in terms of a unique set of internal coordinates. By construction, the Jacobian of the transformation $\boldsymbol{J}(\boldsymbol{r}, \boldsymbol{\xi}')$ is then block diagonal such that $|\boldsymbol{J}(\boldsymbol{r}, \boldsymbol{\xi}')| = |\boldsymbol{J}(\boldsymbol{r})||\boldsymbol{J}(\boldsymbol{\xi}')|$, where $|\cdot|$ indicates the determinant. It follows that

$$Q_{\mathrm{D}} = Q \int d\boldsymbol{p}_\xi d\boldsymbol{\xi}' \, |\boldsymbol{J}(\boldsymbol{\xi}')| e^{-\beta[K(\boldsymbol{p}_\xi) + U_{\mathrm{D}}(\boldsymbol{\xi}')]}, \tag{7}$$

and the expression for the modified average can be reduced such that

$$\langle A \rangle_{\mathrm{D}} = \frac{Q}{Q_{\mathrm{D}}} \int d\boldsymbol{p}_\xi d\boldsymbol{\xi}' \, |\boldsymbol{J}(\boldsymbol{\xi}')| e^{-\beta[K(\boldsymbol{p}_\xi) + U_{\mathrm{D}}(\boldsymbol{\xi}')]} \langle A \rangle = \langle A \rangle. \tag{8}$$

This equality also holds for holonomic constraints, provided that the coordinate transformation includes the degrees of freedom that are fixed. This is frequently the case for rigid molecular mechanics force fields constructed in bond-angle-torsion coordinates since the frozen degrees of freedom are usually bonds or angles.

## 2 Optimization of the Dummy Proton Force Field

The simplest choice for a dummy atom potential is to simply inherit one of the bond, angle, and torsion terms from the original force field. All other interactions are then "zeroed out" so as to meet the single bond-angle-torsion frame required to reproduce thermodynamics. However convenient this may be, it completely neglects all adjustments from Lennard-Jones and electrostatic

interactions and may lead to strange behavior. For example, retaining only these terms for a carboxylate proton would lead to a 1 : 1 ratio of *syn* and *anti* tautomer populations, when the non-dummy force field predicts a ratio closer to 13 : 1 (based on simulations of acetic acid, data not shown). In practice, it is relatively straightforward to pick a convenient coordinate frame and then reparameterize a single set of interactions to reproduce the potential of mean force (PMF) from the original force field. Since a simple functional form is being used, a complete PMF may not even be necessary, since the most desireable property for multiconformer systems is the relative populations and these can be computed with simple indicator functions.

## 2.1  A Simple Two State Model for Carboxylate Conformation

A carboxylate proton has two rotable states, *syn* and *anti*. The former is generally considerably more stable while the latter is likely only appreciably populated in nonhomogeneous solvated systems. It is convenient to define these states in terms of the C-C-O-H torsion $\phi$ such that *syn* corresponds to $\phi \approx \pi$ and *anti* to $\phi \approx 0$ (with $\phi$ in radians and periodic on the interval $[-\pi, \pi]$). A simple model for the potential of mean force $U(\phi)$ is given by a two term Fourier series,

$$
\begin{aligned}
U(\phi) &= U_1 \left(1 + \cos \phi\right) + U_2 \left[1 + \cos \left(2\phi - \pi\right)\right] \\
&= U_1 + U_2 + U_1 \cos \phi - U_2 \cos 2\phi,
\end{aligned}
\tag{9}
$$

where $U_1$ and $U_2$ are positive constants with the dimension of energy. For chemical applications, it is useful to also assume that $U_2 \geq U_1$, such that the *anti* state is always higher in energy. This potential has simple first and second derivatives:

$$
\frac{dU}{d\phi} = -U_1 \sin \phi + 2U_2 \sin 2\phi = 4U_2 \sin \phi \left(\cos \phi - \frac{U_1}{4U_2}\right)
\tag{10a}
$$

$$
\frac{d^2U}{d\phi^2} = -U_1 \cos \phi + 4U_2 \cos 2\phi = 8U_2 \cos \phi \left(\cos \phi - \frac{U_1}{8U_2}\right) - 4U_2,
\tag{10b}
$$

from which it can be seen that, in addition to minima at 0 and $\pm\pi$, there also exist maxima at $\pm\phi^{\ddagger}$, where $\phi^{\ddagger} = \arccos\left(U_1/4U_2\right)$. A more detailed defintion of the states is now possible such that *syn* corresponds to the the ranges $[-\pi, -\phi^{\ddagger})$ and $(\phi^{\ddagger}, \pi]$ and *anti* to the range $(-\phi^{\ddagger}, \phi^{\ddagger})$. The conformational "transition states," $\pm\phi^{\ddagger}$, can arguably be either excluded from or included in both

ranges.

We would like to know what the relative population of each state is. From the above, separate configuration integrals can be written for each state (for brevity, the subscript s is used for *syn* and a for *anti*, note also the use of symmetry):

$$Z_s = 2 \int_{\phi^\ddagger}^{\pi} d\phi \, e^{-\beta U(\phi)} \quad \text{and} \quad Z_a = 2 \int_0^{\phi^\ddagger} d\phi \, e^{-\beta U(\phi)}. \tag{11}$$

As usual $\beta \equiv 1/k_B T$ is the inverse temperature. The relative populations of the two states, $P_s$ and $P_a$, are determined by the relative free energy $\Delta F$ between states in the usual way:

$$\frac{P_a}{P_s} = \frac{Z_s}{Z_a} = e^{-\beta \Delta F}. \tag{12}$$

The first order rate constants also follow from basic transition state theory arguments:

$$k_{as} = \omega Z_a^{-1} e^{-\beta U(\phi^\ddagger)} \quad \text{and} \quad k_{sa} = \omega Z_s^{-1} e^{-\beta U(\phi^\ddagger)}. \tag{13}$$

The (real) transition frequency $\omega$ is taken under a harmonic approximation with the additional assumption that the dummy atom alone is experiencing the potential (i.e., $m = 1.008$ g/mol):

$$\omega \approx \sqrt{\frac{16 U_2^2 - U_1^2}{4 m U_2}}. \tag{14}$$

All of these quantities ($\Delta F$, $k_{as}$, and $k_{sa}$) are easily computed by numerical quadrature.

## 2.2 Optimizing Sampling Efficiency

The dynamics of a two state discrete Markov model can be described by a transition matrix,

$$\boldsymbol{T} = \begin{pmatrix} 1 - P_{sa} & P_{sa} \\ P_{as} & 1 - P_{as} \end{pmatrix}, \tag{15}$$

where

$$P_{as} = \frac{k_{as}}{\omega} \quad \text{and} \quad P_{sa} = \frac{k_{sa}}{\omega}. \tag{16}$$

The system has one non-stationary eigenvector with eigenvalue $\lambda = 1 - P_{\mathrm{sa}} - P_{\mathrm{as}}$ and a characteristic time $\tau = -1/\ln\lambda$, which can be taken as a physical quantity in the context of a "real" proton or else as purely a metric of sampling efficiency in the case of a "dummy" proton. Here we *minimize* $\tau$ as a functional of $U$ subject to the constraints that $U_2 \geq U_1$ and $\Delta F(U_1, U_2) - \Delta F_{\mathrm{FF}} = 0$, where $\Delta F_{\mathrm{FF}}$ is a constant from the force field that dictates the relative populations of the *syn* and *anti* states. The dummy atom model defined by these optimal parameters will have identical thermodynamic behavior to the physical model, but much faster kinetic behavior (in the sense of more rapidly changing state). For the case of acetic acid, the rate constant of each state increases by $\sim$4 orders of magnitude.

## 2.3 Force Field Parameters for CHARMM36

The dummy atom force field for CHARMM36 has been optimized to improve sampling where possible. This is most evident for aspartate and glutamate, but also resulted in completely new parameters for histidine and lysine. In some cases, the criteria for improvement are ambiguous, in which case the parameters are essentially unmodified from the analogous proton parameters and all other interactions are set to zero. This itself can be fairly arbitrary since the only constraint is that the requisite coordinate transformation, which is often non-unique, is well-defined. It is possible that "better" choices of coordinates and parameters exist.

# 3 WHAM Titration Curves: Special Cases

For simple systems where states never differ by more than one proton, the WHAM equations predict titration curves that are exactly sigmoidal. To see this, we begin with the WHAM equation for protonated fractions (Eqs. 22 and 23 in the main text) and note that $n_t$ can only have two values, zero or one, which will appear $N_0$ and $N_1$ times, respectively. Accordingly, there are only two possible values for the weights,

$$w_t^0 = e^{f(\mathrm{pH})}C_0^{-1} \quad \text{and} \quad w_t^1 = e^{f(\mathrm{pH})}10^{-\mathrm{pH}}C_1^{-1}, \tag{17}$$

Table S1: All modified and/or optimized dummy atom force field parameters are presented here. If an angle or torsion containing a one of the new atom types exists topologically, but is not defined, then the interaction is explicitly set to zero. Angle force constants ($k_\theta$ and $k_0$) are in kcal/mol-radian$^2$ while force centers are in degrees. Energy factors ($U_n$) are in kcal/mol and all distances are Å.

| type | description | in residues |
|------|-------------|-------------|
| HD | polar H | ASP, GLU, HIS |
| HCD | 1° amine H | LYS |
| HSD | thiol H | CYS |

| bond | $k_r$ | $r_0$ |
|------|-------|-------|
| HD–OB | 545.0 | 0.96 |
| HD–OC | 545.0 | 0.96 |
| HD–NR2 | 466.0 | 1.00 |
| HCD–NH2 | 403.0 | 1.04 |
| HSD–SS | 275.0 | 1.325 |

| angle | $k_\theta$ | $\theta_0$ |
|-------|-----------|-----------|
| HD–OC–CC | 55.0 | 115.0 |
| HD–OB–CD | 55.0 | 115.0 |
| HD–NR2–CPH1 | 50.0 | 126.0 |
| HCD–NH2–HC | 39.0 | 106.5 |
| HSD–SS–CS | 38.8 | 95.0 |

| dihedral | $k_0$ | $\phi_0$ | $U_1$ | $\phi_1$ | $U_2$ | $\phi_2$ | $U_3$ | $\phi_3$ |
|----------|-------|----------|-------|----------|-------|----------|-------|----------|
| CT2–CD–OB–HD | – | – | 0.88 | 0.0 | 0.88 | 180.0 | – | – |
| CT2–CC–OC–HD | – | – | 0.88 | 0.0 | 0.88 | 180.0 | – | – |
| CT2A–CC–OC–HD | – | – | 0.88 | 0.0 | 0.88 | 180.0 | – | – |
| CPH1–CPH1–NR2–HD | 12.0 | 180.0 | – | – | – | – | – | – |
| CT1–CS–SS–HSD | – | – | 0.20 | 0.0 | 0.65 | 0.0 | 0.22 | 0.0 |
| *NH2–HC–HC–HCD | 200.0 | 37.0 | – | – | – | – | – | – |

* – indicates improper dihedral

where $C_n \equiv (1/N) \sum_{k=1}^{M} N_k \exp \left[ f(\text{pH}_k) - n \ln 10 \text{pH}_k \right]$. As such, the protonated and deprotonated fractions $P_1(\text{pH})$ and $P_0(\text{pH})$ are just summations of identical terms:

$$P_1(\text{pH}) = N_1 w_t^1 = N_1 e^{f(\text{pH})} 10^{-\text{pH}} C_1^{-1} \tag{18a}$$

$$P_0(\text{pH}) = N_0 w_t^0 = N_0 e^{f(\text{pH})} C_0^{-1}. \tag{18b}$$

Using the fact that $P_1(\text{pH}) + P_0(\text{pH}) = 1$ it can be shown that

$$e^{f(\text{pH})} = \frac{C_0 C_1}{N_0 C_1 + N_1 C_0 10^{-\text{pH}}}, \tag{19}$$

which, when inserted into Eq. (18a), yields

$$P_1(\text{pH}) = \frac{1}{1 + \frac{N_0 C_1}{N_1 C_0} 10^{\text{pH}}}. \tag{20}$$

This is just the Henderson-Hasselbalch equation with the correspondence that

$$10^{-\text{p}K_a} = \frac{N_0}{N_1} \frac{\sum_{k=1}^{M} N_k e^{f(\text{pH}_k)} 10^{-\text{pH}_k}}{\sum_{k=1}^{M} N_k e^{f(\text{pH}_k)}}, \tag{21}$$

which is a constant with respect to pH as is required for the titration curve to be a pure sigmoid. One could use this equation to compute the $\text{p}K_a$ in this special case, but in practice a non-linear regression will yield the same result. Note that in the trivial case that $M = 1$ this reduces to a simple estimator for the $\text{p}K_a$ based on inversion of the Henderson-Hasselbalch equation:

$$\text{p}K_a = \text{pH} - \log \frac{N_0}{N_1} \tag{22}$$

The above holds for simple two state systems such as carboxylates, amines, and thiols. A similar derivation holds for more complicated three state systems (e.g., methyl imidazole). Such systems have two independent $\text{p}K_a$ values of the form:

$$10^{-\text{p}K_a} = \frac{N_{01}}{N_{11}} \frac{\sum_{k=1}^{M} N_k e^{f(\text{pH}_k)} 10^{-2\text{pH}_k}}{\sum_{k=1}^{M} N_k e^{f(\text{pH}_k)} 10^{-\text{pH}_k}}, \tag{23}$$

where $N_{ab}$ is the number of times the occupancy vector is observed with the vector $(a, b)$ (i.e., $N =$

$N_{01} + N_{10} + N_{11}$). The other p$K_a$ follows simply by changing $N_{01}$ to $N_{10}$.

# 4    Theoretically Optimal Sampling Efficiency

In a previous work[S1] it was shown that an optimal switch time can be obtained by maximizing the mean transition rate, $k$, between a pair of states. An analytic expression for $k$ can be obtained within nonequilibrium linear response and the assumption of a simple exponential form for the force autocorrelation function. The resulting expression involves two intrinsic physical properties of the system: 1) the magnitude of force fluctuations at equilibrium $\sigma_0^2$ and 2) the "molecular" time scale $\tau_m$ at which these fluctuations vary. Here we note that simple relationships between these quantities and the optimal switch time can be extracted from simple numerical analysis. In particular, it is noted that maximizing the rate leads to a simple condition on its first derivative. In the case that the switch protocol is linear one can obtain

$$\sqrt{2\pi}\sigma(\xi_{\mathrm{opt}})\,\mathrm{erfcx}\left(\frac{\sigma(\xi_{\mathrm{opt}})}{2\sqrt{2}}\right) \leq \frac{\sigma_0^2}{\xi_{\mathrm{opt}}}, \tag{24}$$

where $\xi_{\mathrm{opt}} \equiv \tau_{\mathrm{opt}}/\tau_m$, $\tau_{\mathrm{opt}}$ is the optimal switch time and erfcx $x$ is the scaled complementary error function. The inequality becomes an equality for large $\xi_{\mathrm{opt}}$, but the error is already less than 10% for $\xi_{\mathrm{opt}}$ greater than $\sim 10$. This is because $\sigma(\xi_{\mathrm{opt}})$ rapidly decays to an asymptote of 2.3805 $k_B T$ on this same interval such that the left hand side of Eq. (24) is essentially constant with a value of 2.83475. Inverting Eq. (24) yields

$$\tau_{\mathrm{opt}} \leq \frac{\sigma_0^2 \tau_m}{2.83475}. \tag{25}$$

The optimal mean acceptance probability can also be estimated as

$$\overline{P_{\mathrm{opt}}} \leq \mathrm{erfc}\left(\frac{\sigma(\xi_{\mathrm{opt}})}{2\sqrt{2}}\right) = 23.39\%, \tag{26}$$

such that the optimal mean transition rate is

$$k_{\mathrm{opt}} = \frac{\overline{P_{\mathrm{opt}}}}{\tau_{\mathrm{opt}}} \geq \frac{0.66318}{\sigma_0^2 \tau_m}. \tag{27}$$

However, this limit is only reached under the assumption that $pK_a^{(i)}$ is chosen as the exact $pK_a$ (see the main text and previous work by Radak and Roux[S1]). In principle, $\sigma_0^2$ and $\tau_m$ can be estimated to yield a guess for $\tau_{opt}$, although the simple collinear dependence means that errors in each quantity will compensate in an unhelpful manner. It may be easier and more stable to simply increase the switch time until the mean acceptance probability nears the ideal value while also updating the estimate for $pK_a^{(i)}$ based on a WHAM calculation as in the main text.

# 5 Detailed Fitting Data for SNase

Table S2: Apparent $pK_a$ values for SNase are tabulated from Hill equation fits to the data presented in the main text. Comparison values, where available, are given from both theory and experiment. Error bars have been adjusted to represent 95% confidence intervals. Errors from Huang et al.[S2] reported as zero were assumed to be 0.1 units before rescaling. Note that, although Huang et al.[S2] reported using Hill coefficients during their procedure, the values were not published.

| residue | | this work | | $\lambda$-dynamics[S2] | expt.[S3] | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $pK_a$ | $n$ | $pK_a$ | $pK_a$ | $n$ |
| | 10 | 3.23 (0.60) | 1.12 (0.06) | 3.20 (0.25) | 2.82 (0.22) | 0.85 (0.05) |
| | 43 | 4.44 (0.07) | 1.26 (0.10) | 4.10 (0.25) | 4.32 (0.10) | 0.69 (0.03) |
| | 52 | 5.01 (0.26) | 1.07 (0.15) | 4.70 (0.50) | 3.93 (0.20) | 0.65 (0.07) |
| | 57 | 4.85 (0.33) | 1.07 (0.04) | 4.10 (0.75) | 3.49 (0.22) | 0.83 (0.07) |
| | 67 | 4.23 (0.80) | 1.16 (0.21) | 4.00 (0.50) | 3.76 (0.18) | 0.99 (0.07) |
| GLU | 73 | 3.48 (0.92) | 1.08 (0.18) | 3.60 (0.25) | 3.31 (0.03) | 0.92 (0.03) |
| | 75 | 2.98 (1.31) | 1.10 (0.21) | 2.70 (1.00) | 3.26 (0.12) | 0.79 (0.10) |
| | 101 | 4.55 (0.45) | 1.10 (0.13) | 4.70 (0.50) | 3.81 (0.25) | 0.82 (0.05) |
| | 122 | 3.90 (0.64) | 1.31 (0.31) | 4.40 (0.25) | 3.89 (0.22) | 0.78 (0.07) |
| | 129 | 5.08 (0.61) | 0.91 (0.14) | 5.50 (0.25) | 3.75 (0.22) | 0.66 (0.07) |
| | 135 | 3.35 (0.48) | 1.29 (0.25) | 2.90 (0.25) | 3.76 (0.20) | 0.82 (0.03) |
| | 19 | 2.77 (0.76) | 1.40 (0.29) | 3.30 (1.50) | 2.21 (0.18) | – |
| | 21 | 6.78 (0.99) | 0.68 (0.10) | 6.00 (0.75) | 6.54 (0.05) | – |
| | 40 | 3.32 (0.52) | 1.18 (0.31) | 2.90 (0.25) | 3.87 (0.22) | 0.57 (0.05) |
| ASP | 77 | 0.82 (0.50) | 1.02 (0.92) | <-1.00 | <2.20 | – |
| | 83 | 1.97 (0.72) | 2.43 (1.08) | <0.00 | <2.20 | – |
| | 95 | 2.74 (0.39) | 1.39 (0.22) | 3.00 (0.25) | 2.16 (0.18) | 0.78 (0.05) |
| | 143 | 4.41 (0.64) | 1.13 (0.07) | n/a | 3.80 (0.25) | 0.77 (0.10) |
| | 146 | 4.01 (0.34) | 1.13 (0.12) | n/a | 3.86 (0.12) | 0.75 (0.03) |
| LYS | 24 | 8.43 (0.45) | 1.17 (0.11) | n/a | n/a | n/a |
| HIS | 8 | 6.66 (0.56) | 1.05 (0.05) | n/a | n/a | n/a |
| | 121 | 5.36 (0.50) | 1.03 (0.16) | n/a | n/a | n/a |

# References

(S1)  Radak, B. K.; Roux, B. *J. Chem. Phys.* **2016**, *145*, 134109.

(S2)  Huang, Y.; Chen, W.; Wallace, J. A.; Shen, J. *J. Chem. Theory Comput.* **2016**, *12*, 5411–5421.

(S3)  Castañeda, C. A.; Fitch, C. A.; Majumdar, A.; Khangulov, V.; Schlessman, J. L.; García-Moreno, B. E. *Proteins* **2009**, *77*, 570–588.