

GPU-Accelerated All-atom Particle-Mesh Ewald Continuous Constant pH Molecular Dynamics in Amber

Julie A. Harris,^{†,||} Ruibin Liu,^{†,‡} Vinicius Martins de Oliveira,[†] Erik Vaquez
Montelongo,^{†,¶} Jack A. Henderson,^{†,§} and Jana Shen^{*,†}

[†] *Department of Pharmaceutical Sciences, University of Maryland School of Pharmacy,
Baltimore, MD*

[‡] *Joint first author*

[¶] *Current address: The University of Texas Southwestern Medical Center, Arlington, TX*

[§] *Current address: Scorpion Therapeutics, Boston, MA*

^{||} *Current address: ComputChem LLC, Baltimore, MD*

E-mail: jana.shen@rx.umaryland.edu

Abstract

Constant pH molecular dynamics (MD) simulations sample protonation states on the fly according to the conformational environment and user specified pH condition; however, the current accuracy is limited due to the use of implicit-solvent models or a hybrid solvent scheme. Here we report the first GPU-accelerated implementation, parameterization, and validation of the all-atom continuous constant pH MD (CpHMD) method with particle-mesh Ewald (PME) electrostatics in the Amber22 *pmemd.cuda* engine. The titration parameters for Asp, Glu, His, Cys, and Lys were derived for the CHARMM c22 and Amber ff14sb and ff19sb force fields. We then evaluated the PME-CpHMD method using the asynchronous pH replica-exchange titration simulations with the c22 force field for six benchmark proteins, including BBL, hen egg white lysozyme (HEWL), staphylococcal nuclease (SNase), thioredoxin, ribonuclease A (RNaseA), and human muscle creatine kinase (HMCK). The root-mean-square deviation from the experimental pK_a 's of Asp, Glu, His, and Cys is 0.76 pH units, and the Pearson's correlation coefficient for the pK_a shifts with respect to model values is 0.80. We demonstrated that a finite-size correction or much enlarged simulation box size can remove a systematic error of the calculated pK_a 's and improve agreement with experiment. Importantly, the simulations captured the relevant biology in several challenging cases, e.g., the titration order of the catalytic dyad Glu35/Asp52 in HEWL and the coupled residues Asp19/Asp21 in SNase, the large pK_a upshift of the deeply buried catalytic Asp26 in thioredoxin, and the large pK_a downshift of the deeply buried catalytic Cys283 in HMCK. We anticipate that PME-CpHMD offers proper pH control to improve the accuracies of MD simulations and enables mechanistic studies of proton-coupled dynamical processes that are ubiquitous in biology but remain poorly understood due to the lack of experimental tools and limitation of current MD simulations.

1 Introduction

Accurate and efficient molecular modeling of proton-coupled dynamic processes is important, as biological functions and material properties often depend on protonation and deprotonation. For example, many secondary active membrane transporters utilize pH gradient and proton coupling to drive the conformation transitions for function.¹ Many enzymes have pH-dependent catalytic activities, e.g., the active site of SARS-CoV-2 main protease collapses upon protonation of a conserved histidine residue.^{2,3} Well known examples of pH-dependent materials include aminopolysaccharide chitosan which self assembles into hydrogels in response to a small increase in solution pH.⁴ The ability to model proton-coupled dynamic processes is also important for studying protein-ligand binding, as upon protein-ligand association, the protonation state of the protein and the ligand may change.^{5,6}

Unlike the conventional molecular dynamics (MD) that assumes fixed protonation states, constant pH MD allows protonation states to evolve with time according to the conformational environment and a preset solution pH. Currently, perhaps the most popular constant pH approaches are based on λ dynamics and the hybrid Monte-Carlo (MC)/MD scheme (also known as the stochastic titration method⁷). The former^{8–12} uses continuous titration coordinates to propagate protonation states based on an extended Lagrange approach called λ dynamics,¹³ while the latter^{7,14–16} combines MD with periodic MC sampling of discrete protonated and deprotonated states. Hereafter, we will refer to the former as the continuous and the latter as the discrete constant pH methods. The details of these techniques can be found in the recent reviews.^{17–19} Although the first (discrete) constant pH method⁷ is based on a hybrid-solvent scheme (see later discussion), the early implementations of the constant pH methods are solely based on the implicit-solvent generalized Born (GB) models, i.e., both conformational and protonation state sampling is conducted in implicit solvent.^{8,9,14} The use of the implicit-solvent models significantly reduces sys-

tem size and allows faster sampling of solute conformational states relative to simulations with explicit water models.¹⁷ However, for many biologically relevant systems, e.g., transmembrane proteins (with heterogeneous dielectric environment), nucleic acids (highly charged), and protein-ligand and protein-protein bound states, implicit-solvent models are not sufficiently accurate. This has motivated the development of explicit-solvent constant pH methods, which include the hybrid-solvent scheme and the all-atom approaches.

In a hybrid-solvent constant pH scheme, the MC sampling of protonation states or λ dynamics propagation of titration coordinates is conducted in implicit solvent, while MD is conducted in explicit water. The first hybrid-solvent constant pH method was developed by Baptista and Soares, who combined MD in explicit solvent with MC sampling based on Poisson-Boltzmann (PB) calculations.⁷ This method was first implemented in GROMOS96²⁰ and later improved and implemented²¹ in GROMACS.²² Following the aforementioned work and making use of the state-of-the-art GB models, the hybrid-solvent continuous²³ and discrete¹⁵ constant pH methods were developed and implemented in CHARMM²⁴ and Amber.²⁵ Compared to the purely GB based constant pH methods, the hybrid-solvent approaches demonstrated significantly improved accuracy for conformational dynamics and consequently better agreement with the experimental pK_a 's.^{15,23,26,27} Importantly, the hybrid-solvent approach allowed the investigations of pH-dependent mechanisms of a variety of systems that are (due to inaccuracy) unfeasible to model with implicit-solvent models, e.g., proteins in mixed solvent,²⁸ phase transition of surfactants,²⁹ polysaccharides,⁴ and lipid bilayers,³⁰ proteins at the water-membrane interface,³¹ as well as transmembrane proteins³² and peptides inserted in the membrane.³³ Nonetheless, a drawback is that the Hamiltonian cannot be expressed in a hybrid scheme (semigrand canonical ensemble), and thus energy conservation is not proven to hold. In terms of applications, hybrid-solvent simulations of protein-ligand complexes are challenging, as the implicit-solvent description for ligand is not very accurate and the effects of explicit water and ions which play significantly roles cannot be fully

modeled.³⁴

To overcome the limitations of the hybrid-solvent scheme, much effort has been made in the development of all-atom constant pH methods over the past decade. The CHARMM program²⁴ contains the CPU implementations of the all-atom continuous constant pH method with generalized reaction field^{35,36} or particle-mesh Ewald (PME) electrostatics for λ dynamics,¹¹ and the multiple site λ dynamics (MS λ D)³⁷ based constant pH method.^{10,38} These methods have been validated using pK_a calculations for a number of proteins^{10,11,36} as well as RNAs.³⁸ The λ dynamics based constant pH method was also implemented in the GROMACS program,²² although only the single-site titratable model was considered and performance for proteins remains to be demonstrated.¹² The NAMD program³⁹ contains an implementation of the all-atom constant pH method based on a non-equilibrium MD-MC approach,¹⁶ which overcomes the issue of low acceptance of MC moves due to a large energy change resulting from a sudden switch in protonation state, as in the aforementioned hybrid MC/MD constant pH approaches.^{15,21}

The aforementioned all-atom continuous constant pH methods^{10,11,35,38} are promising; however, the CPU implementations limit the simulation time scale and system size that can be studied. Recently, the Brooks group developed the basic lambda dynamics engine (BLaDE) which enables GPU acceleration for MS λ D based alchemical free energy calculations and constant pH simulations.⁴⁰ In this work, we report the development and validation of the GPU all-atom continuous constant pH method in the *pmemd.cuda* engine of Amber program (version 2020).⁴¹ Following the discussion of the model parameterization and validation, we present the data from the pK_a calculations of benchmark proteins, including BBL, HEWL, SNase, RNase A, BACE 1, thioredoxin, and HMCK. In addition to comparison to experimental pK_a values, we will discuss the coupled titration of catalytic residues, pH-dependent response of solvent exposure, titration of deeply buried sites. Finally, we will discuss the finite-size effects and project future directions.

2 Methods and Implementation

The all-atom PME continuous constant pH MD (CpHMD) method. In contrast to the conventional MD, the continuous constant pH MD (CpHMD) method treats the protonation states of titratable sites as dynamic variables $\{\lambda_\alpha\}$ and propagates them simultaneously with the spatial coordinates using an extended Hamiltonian,^{8,42}

$$\begin{aligned} \mathcal{H}(\{\mathbf{r}_i\}, \{\lambda_\alpha\}) = & \frac{1}{2} \sum_i m_i \dot{\mathbf{r}}_i^2 + \frac{1}{2} \sum_\alpha m_\alpha \dot{\lambda}_\alpha^2 + U^{\text{ind}}(\{\mathbf{r}_i\}) \\ & + U^{\text{hybr}}(\{\mathbf{r}_i\}, \{\lambda_\alpha\}) + \sum_\alpha U^*(\lambda_\alpha), \end{aligned} \quad (1)$$

where $\{\mathbf{r}_i\}$ and $\{\lambda_\alpha\}$ refer to the spatial and titration coordinates, respectively. A deprotonated state is represented by the λ values close to 1 ($\lambda > 0.8$ in this work), whereas a protonated state is represented by the λ values close to 0 ($\lambda \leq 0.2$ in this work). In order to impose the boundaries 0 and 1 for λ , we express it as^{11,23,42}

$$\lambda = \sin^2 \theta, \quad (2)$$

where the θ variable is allowed to assume any real value, as with the spatial coordinates. Therefore, θ is the actual coordinate in the integrator. However, for the convenience of discussion, we will write all equations in terms of λ .

The two first terms in the Hamiltonian (Eq. 1) describe the kinetic energies of the real atoms and λ particles. The λ particles are assigned a fictitious mass, which is similar to a heavy atom (10 amu). U^{ind} represent the λ -independent bonded energies (see later discussion) and non-bonded energies. For the all-atom CpHMD method, U^{hybr} is a sum of the λ -dependent Lennard Jones and electrostatic energies.¹¹ The last term U^* in the Hamiltonian (Eq. 1) represents three biasing potentials that are only dependent on λ ,

$$U^*(\lambda_\alpha) = -U^{\text{mod}}(\lambda_\alpha) + U^{\text{barr}}(\lambda_\alpha) + U^{\text{pH}}(\lambda_\alpha). \quad (3)$$

U^{mod} represents the potential of mean force (PMF) for titrating a model compound or peptide in solution, which can be obtained from the traditional free energy simulations such as thermodynamic integration (TI). U^{barr} is a quadratic barrier potential centered in the middle of the λ coordinate to prolong the residence times of the end states (λ close to 0 or 1):

$$U^{\text{barr}}(\lambda_\alpha) = 4\beta(\lambda_\alpha - 1/2)^2, \quad (4)$$

where β is a parameter affecting the barrier height. In the current implementation, it is set to 2.0 kcal/mol for all types of residues, similar to our previous work.¹¹ U^{pH} represents the free energy added to the deprotonation reaction due to a change in solution (infinite proton bath) pH

$$U^{\text{pH}}(\lambda_\alpha) = \ln 10 \cdot k_B T (\text{pH} - \text{p}K_a^{\text{mod}}) \lambda_\alpha \quad (5)$$

where k_B is the Boltzmann constant, T is the system temperature, and $\text{p}K_a^{\text{mod}}$ is reference $\text{p}K_a$ of the model compound or peptide.

When $\lambda = 0$, the proton is present and fully interacts with its environment, and when $\lambda = 1$, it is treated as a ghost particle without non-bonded interactions with its environment. The partial charges on the titratable residue are linearly scaled between the protonated and deprotonated states, as in the original CpHMD framework.^{8,42} This differs from the multi-site λ dynamics³⁷ dynamics based MS λ D CpHMD method,^{10,38} which scales potential energies. Formally, the λ -dependent Lennard-Jones interaction energy between a titratable hydrogen i and another non-titratable atom j is given by

$$\tilde{U}_{ij}^{LJ}(\lambda_i) = (1 - \lambda_i) U_{ij}^{LJ}, \quad (6)$$

where λ_i is the titration variable associated with the titratable hydrogen, and U_{ij}^{LJ} is the Lennard-Jones interaction energy between atoms i and j when the hydrogen is present. Similarly, the Lennard-Jones interaction energy between two titratable hydrogens is given

by

$$\tilde{U}_{ij}^{LJ}(\lambda_i, \lambda_j) = (1 - \lambda_i)(1 - \lambda_j) U_{ij}^{LJ}. \quad (7)$$

The charge of atom j in the titratable residue α is given by

$$q_j(\lambda_\alpha) = (1 - \lambda_\alpha) q_j(0) + \lambda_\alpha q_j(1), \quad (8)$$

where $q_j(0)$ is the charge appropriate to the protonated form of the residue, and $q_j(1)$ is the charge appropriate to the deprotonated form.

The implementation presented in this paper uses fully explicit water molecules and treats the nonbonded electrostatic interaction energy between atoms with particle-mesh Ewald (PME) electrostatics. Because the λ values are treated as dynamic coordinates of the system, the derivatives of the energy with respect to the λ values are required. In the *pmemd* implementation²⁵ of PME electrostatics, this interaction energy is separated into several terms,

$$V_{\text{Coulomb}} = V_{\text{direct}} + V_{\text{reciprocal}} + V_{\text{self}} + V_{\text{plasma}}, \quad (9)$$

where

$$V_{\text{direct}} = \frac{1}{2} \sum_{\mathbf{n}} \sum_{i,j}^{n_{\text{atoms}}} q_i q_j \frac{\text{erfc}(\alpha r_{ij,\mathbf{n}})}{r_{ij,\mathbf{n}}} \quad (10)$$

is the short-range component of the electrostatic energy, where \mathbf{n} enumerates the copies of each atom from the neighboring periodic cells. This summation is performed only over those atom pairs i, j for which r_{ij} falls within a small cutoff distance.

$$V_{\text{reciprocal}} = \frac{1}{2\pi\nu} \sum_{\mathbf{m} \neq 0} \frac{\exp(-(\pi\mathbf{m}/\alpha)^2)}{\mathbf{m}^2} S(\mathbf{m}) S(-\mathbf{m}) \quad (11)$$

is the reciprocal space energy, where ν is the volume of the unit cell, \mathbf{m} is reciprocal lattice vector, and $S(\mathbf{m})$ is the structure factor,

$$S(\mathbf{m}) = \sum_{i=1}^{n_{atoms}} q_i \exp(2\pi i \mathbf{m} \times \mathbf{r}_i), \quad (12)$$

which can be approximated by

$$\begin{aligned} S(\mathbf{m}) &\approx \sum_{k_1, k_2, k_3} Q(k_1, k_2, k_3) \exp\left(2\pi i \left(\frac{m_1 k_1}{K_1} + \frac{m_2 k_2}{K_2} + \frac{m_3 k_3}{K_3}\right)\right) \\ &= F(Q)(m_1, m_2, m_3), \end{aligned} \quad (13)$$

where $Q(k_1, k_2, k_3)$ is a matrix constructed by interpolating the charge distribution in the simulation cell to a grid with the same dimensions k_1, k_2, k_3 , and $F(Q)(m_1, m_2, m_3)$ is the fast Fourier transform of the Q matrix.

The $V_{\text{reciprocal}}$ can then be written as

$$\frac{1}{2\pi\nu} \sum_{\mathbf{m} \neq 0} \frac{\exp(-(\pi \mathbf{m}/\alpha)^2)}{\mathbf{m}^2} F(Q)(\mathbf{m}) F(Q)(-\mathbf{m}), \quad (14)$$

$$V_{\text{self}} = \frac{-\alpha}{\sqrt{\pi}} \sum_{i=1}^{n_{atoms}} q_i^2, \quad (15)$$

is a term that removes the self-interaction energies contained in $V_{\text{reciprocal}}$, and

$$V_{\text{plasma}} = -\frac{\pi}{2V\alpha^2} \left(\sum_i q_i \right)^2, \quad (16)$$

where V is the volume of the unit cell is a term that counterbalances any net charge on the system.

Implementation in the *pmemd.cuda* engine. As in our previous CPU implementation of the PME-CpHMD method in CHARMM,¹¹ the derivatives of $\tilde{U}_{ij}^{LJ}(\lambda_i)$ with respect

to the titration variables can be derived from Eqs. 6 and 7, and computing them requires changes to be made to the Lennard-Jones forces between titratable atoms. In the present implementation, these modifications were made by making appropriate changes to the direct-force CUDA kernel in *pmemd.cuda* where the Lennard-Jones forces are computed. This kernel was also modified to compute the Lennard-Jones contributions to the forces on the λ titration variables. The electrostatic spatial forces on the atoms can be made to depend on the λ values by using the normal force calculations with the charges given in Eq. 8 according to the instantaneous values of the λ titration variables. Since V_{self} and V_{plasma} are independent of the spatial coordinates of the atoms, they are not computed during standard MD runs in *pmemd.cuda*. However, because these energies do depend on the λ titration coordinates, their derivatives with respect to the titration coordinates are required in CpHMD. The calculation of these derivatives was added to the kernel that interpolates the λ -dependent atomic partial charges, which was previously implemented by us for the generalized Born based CpHMD method.⁴³ This kernel otherwise required minimal changes for the present implementation. The derivatives of V_{direct} with respect to λ are computed through appropriate changes to the direct-force kernel in *pmemd.cuda*. The derivatives of $V_{\text{reciprocal}}$ with respect to λ are computed with a new kernel that computes the derivatives given by Eq. 14 using the same method as outlined in our previous CPU implementation in CHARMM.¹¹

Modification of the force field parameters. The current constant pH methods are based on single topology, i.e., titration is represented by switching on/off the charge and Lennard-Jones interactions of the dummy hydrogen as well as by transforming between the protonated and deprotonated forms of sidechain partial charges.¹⁹ The latter is straightforward to implement for the CHARMM force fields,⁴⁴ as the backbone partial charges are independent of the side chain. This is however not the case for the Amber force fields,^{45,46} in which the backbone charges are dependent on the side chain protonation

state. Due to the 1-4 interactions between the backbone and adjacent sidechain, this dependency makes it impossible to use a single reference scheme, i.e., one model for one type of sidechain. To circumvent this problem, we adopted the scheme used in the discrete constant pH implementation in Amber¹⁴ by fixing the backbone charges to the values of one protonation state (charged Asp/Glu and neutral His in our implementation) and absorbing the residual change in charge (ranging from 0.10 to 0.14 e for Asp, Glu and His) onto the $C\beta$ atom. Such a scheme is not ideal and might introduce potential artifacts to conformational dynamics; thus, we only adopted it for titration dynamics. For conformational dynamics, the partial charges are unmodified and the charge interpolation between the protonated and deprotonated states is made to both backbone and sidechain atoms. Here we note that in our approach the conformational dynamics and titration dynamics are treated on an equal footing, with both sets of coordinates propagating together. We don't separate these into separate phases of the simulation. Simply, we use different charge sets for the forces on the titration and spatial coordinates. By doing so, the conformational dynamics from the optimized force field is preserved.

Another compromise and approximation we made is in the bonded terms, which are not scaled between two protonation states as in the early CpHMD implementations.^{8,42} For Asp and Glu, the bonded parameters for the protonated and deprotonated forms are different in both CHARMM and Amber force fields.^{44,46} The parameters of the deprotonated forms (which are most common at physiological pH) were used except for those related to the dummy hydrogens, which were taken from the protonated forms. For His, the bonded parameters for the protonated and deprotonated forms are the same in the Amber ff14SB⁴⁶ and CHARMM c22⁴⁴ force fields. Including the bonded terms in the calculation would require significant modifications to the code, and as such is deferred to the future work. However, from a large number of application studies we have conducted, no artifacts have been observed, which may be due to our choice of adopting the dominant form (e.g., charged Asp/Glu). Therefore, we believe that the improvement with adding the

bonded term perturbation may be minimal.

Finally, the use of two dummy hydrogens for Asp/Glu introduces an issue, namely, once an uncharged (ghost) dummy hydrogen rotates to the *anti* configuration, it loses the ability to titrate. This is because a ghost proton in the *anti* position is unfavorable to protonate, and due to zero force it is unlikely to move until it is protonated, as noted in the early developments of constant pH methods.^{8,14} Following our previous work,⁴² the rotation barrier of the C-O bond in the carboxylate group of Asp/Glu was increased to 6 kcal/mol to keep the dummy hydrogens in the *syn* configuration. This is a limitation, as the *anti* configuration might become favorable in some protein, although it is unfavorable in the peptide.⁴⁷ One solution is to include both *anti*- and *syn*- positions for each oxygen as implemented in the discrete constant pH methods in Amber.^{14,15} This solution however is difficult to implement for CpHMD methods, as it would add additional variables which makes the analytic form of the model PMF impossible to derive (Eq. 18). To complicate the case, experimental evidence of *syn* vs. *anti* configuration is lacking. This is a topic that warrants future investigation.

Potential of mean force functions for model titration. The linear response theory states that the charging free energy of an ion in polar solvent is quadratic in the charge perturbation.⁴⁸ Thus, the PMF for protonation/deprotonation of a single-site titratable group (e.g., Cys and Lys) in explicit solvent can be approximated as a quadratic function in terms of λ .^{11,36}

$$U_{\text{single}}^{\text{mod}}(\lambda) = A(\lambda - B)^2. \quad (17)$$

Following our previous work,⁴² for residues with two titratable sites such as carboxylic acids and histidines an additional variable x is introduced to represent the tautomer states. The underlying variable θ_x which is defined in analogy to θ (Eq. 2) is dynamically propagated on the same footing as θ . For carboxylic acids Asp and Glu which have two equivalent protonation sites (carboxyl oxygens), the model PMF function can be written

as^{9,42,49}

$$U_{\text{Asp/Glu}}^{\text{mod}}(\lambda, x) = (R_1\lambda^2 + R_2\lambda + R_3)(x + R_4)^2 + R_5\lambda^2 + R_6\lambda \quad (18)$$

where R_1, \dots, R_6 are parameters that can be determined by one-dimensional fitting of the corresponding mean forces $(\partial U / \partial \theta)|_{\theta_x}$ and $(\partial U / \partial \theta_x)|_{\theta}$ calculated using thermodynamic integration (TI) at different combinations of θ and θ_x values. The detailed derivation and protocol are given in Ref.^{9,49}

The model PMF function for His titration can be written as⁹

$$\begin{aligned} U^{\text{mod}} = & A_{10}\lambda^2x^2 + 2(A_1B_1 - A_0B_0)\lambda x \\ & + 2(A_0B_0 - A_1B_1 - A_{10}B_{10})\lambda^2x \\ & + A_1\lambda^2 - 2A_1B_1\lambda. \end{aligned} \quad (19)$$

The parameters in Eq. 19 are those in the one-dimensional PMF functions, where either λ or x is fixed at one of the end points (1 or 0).⁹

$$U_{\text{His}}^{\text{mod}}(\lambda, 0) = A_0(\lambda - B_0)^2 \quad (20)$$

$$U_{\text{His}}^{\text{mod}}(\lambda, 1) = A_1(\lambda - B_1)^2 \quad (21)$$

$$U_{\text{His}}^{\text{mod}}(1, x) = A_{10}(x - B_{10})^2 \quad (22)$$

Detailed protocols for obtaining the parameters are given in Ref.^{9,49}

Finite-size corrections to the calculated pK_a values for proteins. In our previous work,¹¹ we proposed a correction for the pK_a 's calculated from the all-atom PME constant pH simulations under periodic boundary conditions. According to the analysis of Hünenberger and colleagues, the finite-size errors for the ligand charging free energies arise from four physical effects, among which the discrete solvent effect dominates when

the protein's net charge is neutralized by counter-ions.⁵⁰ The discrete solvent effect arises from a homogeneous, constant potential that is applied to offset the potential generated by isotropically tumbling solvent molecules so that the average potential over the simulation box is zero.⁵⁰ This “offset” potential is positive for typical three-site water models, and needs to be corrected when calculating ligand charging free energies.⁵⁰ Hünenberger and colleagues developed an analytic correction to the ligand charging free energy⁵⁰

$$\Delta G^{\text{corr}}(\text{charging}) = -\frac{2\pi}{3} k \gamma^s Q \frac{N^s}{V}, \quad (23)$$

where k is the electrostatic constant, γ^s is the quadrupole moment trace of the solvent model relative to a van der Waals interaction site. γ^s is calculated as $0.764 \text{ e} \cdot \text{\AA}^2$ for TIP3P water model.¹¹ Q is the charge (-1 for charging to -1 e and +1 for charging to +1 e), N^s is the number of solvent molecules, and V is the simulation box volume. We note, the correction in Eq. 23 is very similar to that proposed by Roux and coworkers.⁵¹

Now we consider the deprotonation reactions of protein titratable residues, which refers to the charging process of an acidic sidechain, e.g., aspartic acid, $\text{Asp} \longrightarrow \text{Asp}^-$, or the discharging process of a basic sidechain, e.g., histidine, $\text{His}^+ \longrightarrow \text{His}$. Based on Eq. 23 (correction for the charging free energy), we obtain the correction for the deprotonation free energy of a titratable residue in a protein in reference to a model system

$$\Delta \Delta G^{\text{corr}}(\text{deprot}) = \frac{2\pi}{3} k \gamma^s \left(\frac{N_p^s}{V_p} - \frac{N_m^s}{V_m} \right), \quad (24)$$

where N_p^s/V_p and N_m^s/V_m refer to the solvent number density in the protein and model systems, respectively. Note, the minus sign in Eq. 23 and Q are absorbed due to the fact that Q is -1 for acidic residues, and for basic residues $\Delta G(\text{deprot}) = -\Delta G(\text{charging})$. The corresponding $\text{p}K_a$ correction is

$$\Delta \text{p}K_a^{\text{corr}}(\text{deprot}) = \frac{\Delta \Delta G^{\text{corr}}(\text{deprot})}{\ln(10)RT} \quad (25)$$

where R is the ideal gas constant, T is the temperature. Since the solvent number density is higher in the model system than in the protein system, the pK_a correction is negative for both acidic and basic sites.

3 Simulation Protocols

Preparation of model peptide systems. Capped pentapeptides ACE-AAXAA-NH₂ (X = Asp, Glu, His, Cys, or Lys) were used to parameterize and validate the model PMF functions. First, each peptide structure was generated and placed in a cubic water box using CHARMM scripts (version c38b2).²⁴ The minimum distance between the heavy atoms of the peptide and the edges of the box was set as 10 Å. Next, to neutralize the system at pH 7.5, one Cl⁻ counterion was added to the Lys pentapeptide system, and one Na⁺ counterion was added to the Asp and Glu pentapeptide systems. The peptides were represented by the CHARMM c22,⁴⁴ Amber ff14SB,⁴⁶ or Amber ff19SB⁵² force field. Water was represented by the TIP3P water model.⁵³

Thermodynamic integration and titration simulations of the model peptides. We carried out an energy minimization in each pentapeptides system applying a force constant of 100 kcal mol⁻¹Å⁻² to the peptide heavy atoms for 200 steps of SD followed by 300 steps of conjugate gradient method. Then, the system was heated from 100 to 300K using Langevin thermostat and a force constant of 5 kcal mol⁻¹Å⁻² on the heavy atoms. After heating, three stages of equilibration were performed with 250 ps each, whereby the force constant was 2 and 1, and 0 kcal mol⁻¹Å⁻². Finally, thermodynamic integration (TI) simulations were conducted for the model pentapeptides under constant NPT conditions at fixed θ or θ_x values of 0.2, 0.4, 0.6, 0.7854, 1.0, 1.2, and 1.4. Each simulation lasted 10 ns. The TI simulations gave the mean forces, $\langle \partial U / \partial \theta |_{\theta_x} \rangle$ and $\langle \partial U / \partial \theta_x |_{\theta} \rangle$, which were used to obtain parameters in the PMF functions (Eq. 17, 18, and 19). The detailed

protocols are given in a recent tutorial.⁴⁹

As validation of the model parameters, titration simulations were conducted for the model peptides at independent pH conditions, which were placed at 0.5-pH intervals in the range of 2–5.5 for Asp, 2.5–6 for Glu, 4.5–8 for His, 6.5–10 for Cys, and 8–11.5 for Lys model peptides. The equilibration and production runs of the peptide systems followed the same protocols as the protein simulations (see latter discussion). The production run at each pH lasted 20 ns and was repeated three times. With the CHARMM c22 force field,⁴⁴ we also performed pH replica-exchange simulations of 10 ns/replica for Asp, His, and Lys model peptides with the same pH conditions. Additional pH replica-exchange simulations were also performed with the hydrogen mass repartition scheme^{54,55} and 4-fs timestep. The simulation length was 10 ns/replica.

Preparation of the protein systems. For protein simulations, the following PDB files were downloaded: 1W4H (peripheral subunit-binding domain protein BBL),⁵⁶ 2LZT (hen egg white lysozyme or HEWL),⁵⁷ 3BDC (Staphylococcus nuclease or SNase),⁵⁸ 7RSA (ribonucleas A or RNaseA),⁵⁹ 1ERU (thioredoxin),⁶⁰ and 1I0E (human muscle creatine kinase or HMCK).⁶¹ The coordinates were first processed using the convpdb.pl script from MMTSB Toolset⁶² to remove hetero atoms, ions, water, ligands, and hydrogen atoms. The CHARMM c22 protein force field and CHARMM modified TIP3P water model were used to represent the protein and water, respectively.⁴⁴ The following steps were performed using the CHARMM package (c38b2).²⁴ The proteins were embedded in a pre-equilibrated cubic TIP3P water box with at least 10 Å cushion between the protein heavy atoms and the edges of the box. Sodium and chloride ions were added to neutralize the systems (assuming model pK_a 's and pH 7.5) and to provide a physiological (0.15 M) or experimental salt concentration (0.1 M for SNase, 0.5 M for thioredoxin, and 0.06 M for RNase A). Using the HBUILD facility, missing hydrogens were added, and a custom CHARMM script is used to add two dummy hydrogens on the carboxylate oxygens.⁹

The protein structures were energy minimized using 50 steps of steepest descent (SD) method with a harmonic force constant of $50 \text{ kcal} \cdot \text{mol}^{-1} \text{Å}^{-2}$ on the heavy atoms followed by 100 steps of adoptive basis Newton-Raphson (ABNR) method.

Equilibration of the protein systems at independent pH conditions. The CHARMM22 topology and parameter files were converted to the Amber compatible format with the command `chamber` in ParmEd.⁶³ With the Amber input files prepared, a last round of minimization was performed in Amber22,²⁵ using 200 steps of SD followed by 300 steps of conjugate gradient method, whereby a force constant of $100 \text{ kcal mol}^{-1} \text{Å}^{-2}$ was applied to the protein heavy atoms. Keeping the same restraint and with a time step of 1 fs, the system was then heated for 100 ps from the initial temperature of 100 K to 300 K using Langevin thermostat). Following heating, two stages of equilibration was performed. The first stage consisted of two runs of 250 ps each performed at pH 7, whereby the harmonic force constant was 100 and $10 \text{ kcal} \cdot \text{mol}^{-1} \text{Å}^{-2}$. The second stage of equilibration was performed at the individual pH conditions of the replica-exchange simulations. Here, four runs of 500 ns were performed using a time step of 2 fs. The heavy-atom force constant was gradually reduced from 10.0 to 1.0, 0.1, and $0.0 \text{ kcal mol}^{-1} \text{Å}^{-2}$.

Production CpHMD simulations of proteins with pH replica-exchange. For CpHMD production runs, the asynchronous pH replica exchange algorithm⁶⁴ was employed to accelerate sampling convergence of conformational and protonation states and accelerate pK_a calculations.²³ 2 NVIDIA GTX 2080 Ti GPU cards were used. The pH range of the protein simulations was extended at least 1 pH unit below or above the lowest or highest experimental pK_a values, and the pH spacing was 0.5 pH unit. Additional pH replica at 0.25 pH units were added in some cases to increase the probabilities of replica exchange. The exchanges between adjacent pH replicas were attempted every 2 ps (1000 MD steps). Each replica in the simulations of BBL, HEWL, SNase, thioredoxin, RNase A, and HMCK was run for 34, 40, 40, 50, 40, and 30 ns, respectively. The simulation

length was sufficient to converge the pK_a 's of all titratable sites (for HMCK we were only interested in Cys283). For SNase, additional simulations with larger box sizes were carried out. In these systems, the distance between the protein and edges of the water box was increased from the default 10 Å to 12, 14, and 18 Å, and the corresponding simulations lasted 20, 20, 60, and 75 ns per replica, respectively. All settings in the CpHMD are identical to our previous work.^{11,43}

Settings in the MD. Unless otherwise noted, the integration timestep in the production runs was 2 fs. Lennard Jones energies and forces were smoothly switched off over the range of 10–12 Å. For long-range electrostatics, the PME method was used with a real-space cutoff of 12 Å and grid spacing of 1 Å. Each pH replica simulations was performed under constant NPT conditions, where the pressure was maintained at 1 atm by the Berendsen barostat with a relaxation time of 0.1 ps and the temperature was maintained at 300 K by the Langevin thermostat with a collision frequency of 1.0 ps^{-1} .²⁵

pK_a calculation. λ coordinates from the titration simulations were post-processed to calculate pK_a values. Following our previous definition of protonated and deprotonated states,⁴³ $\lambda \leq 0.2$ and $\lambda \geq 0.8$ represent the protonated and unprotonated states, respectively, while $0.2 < \lambda < 0.8$ is considered unphysical and discarded. The unprotonated fractions S at all simulation pH conditions were collected and the data were fit to the Hill (or generalized Henderson-Hasselbalch) equation

$$S = \frac{1}{1 + 10^{n(pK_a - \text{pH})}}, \quad (26)$$

where pK_a and n are the fitting parameters, and S is defined as $S = N_{\text{unprot}} / (N_{\text{unprot}} + N_{\text{prot}})$, where N_{unprot} and N_{prot} refer to the number of λ values representing the unprotonated and protonated states, respectively.

For two residues experiencing linked titration, the average number of protons bound to

the two residues ($\langle P \rangle$) are calculated at all simulation pH, and fit to the following coupled titration model to determine the macroscopic stepwise pK_a 's,^{35,65}

$$\langle P \rangle = \frac{10^{pK_{a2}-pH} + 2 \cdot 10^{pK_{a1}+pK_{a2}-2pH}}{1 + 10^{pK_{a2}-pH} + 10^{pK_{a1}+pK_{a2}-2pH}} \quad (27)$$

where pK_{a1} and pK_{a2} are the two stepwise pK_a 's.

Finite-size corrections. A finite-size correction (Eq. 25) was applied to the calculated pK_a 's. For the pK_a 's in Table 2 (i.e., a minimum of 10 Å distance between the protein and the edge of the water box), the corrections are: BBL (Asp: -0.33, Glu: -0.39, His: -0.30); HEWL (Asp: -0.63, Glu: -0.70, His: -0.61); SNase (Asp: -0.70, Glu: -0.77, His: -0.67); thioredoxin (Asp: -0.96, Glu: -1.02, His: -0.93); RnaseA (Asp: -0.66, Glu: -0.72, His: -0.63); creatine kinase (Cys: -0.64). Corrections for the simulations with larger box sizes (Table 3) are given in Table S1. At a first glance, it may seem odd that these corrections differ by residue type. This is because the corrections for the model pK_a 's are different. In the future, these differences can be eliminated by using larger solvent boxes for the model simulations. Additionally, the reference pK_a 's can be adjusted to account for the pK_a corrections which can be calculated at the simulation set up by using lattice parameters.¹¹

4 Results and Discussion

4.1 Model parameterization and validation

Parameterization of the model potential of mean functions for titrating model pentapeptides. First, TI simulations of model pentapeptides $\text{CH}_3\text{CO-Ala-Ala-X-Ala-Ala-CONH}_2$ (X=Asp, Glu, His, Cys or Lys) were performed to obtain the mean forces, $\langle \partial U / \partial \theta \rangle |_{\theta_x}$ and/or $\langle \partial U / \partial \theta_x \rangle |_{\theta}$, which were then fit to the analytic functions (derivatives of Eqs.17, 18,

and 19 expressed in θ) to obtain the parameters. The fitting was generally very good (see an example fitting of His in Fig. 1), suggesting that the linear response theory holds, consistent with the results of both the GRF-based and PME-CpHMD in CHARMM.^{11,36} Integration of the mean forces followed by coordinate transformation gives the PMF as a function of λ (see examples in Ref³⁶). We note, the parameters are more accurate when they are derived from fitting the mean forces (as in our early work^{23,42} rather than the PMF (as in the PME-CpHMD implementation in CHARMM¹¹).

Table S2 gives the parameters in the model PMF functions of Asp, Glu, and His (Eq. 18 and 19) for the CHARMM c22,⁴⁴ Amber ff14sb,⁴⁶ and ff19sb⁵² force fields. The model PMF parameters for Cys titration were also obtained for the CHARMM c22⁴⁴ and Amber ff14sb⁴⁶ force fields (Table S2). In the rest of the paper, we focused on the CHARMM c22 force field⁶⁶ to facilitate comparisons with our previous PME-CpHMD¹¹ and the Brooks' lab's MS λ D CpHMD implementations^{10,37} in CHARMM.²⁴ A force field comparison study will be conducted in the near future.

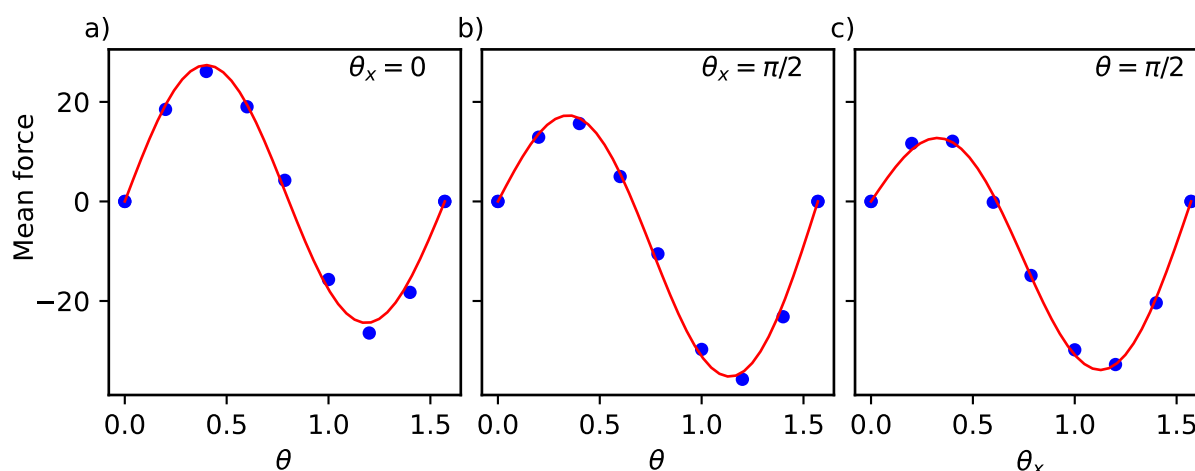


Figure 1: Nonlinear fitting of the mean forces to obtain the model PMF parameters for His titration. a) and b) Fitting $\langle \partial U / \partial \theta \rangle$ at $\theta_x = 0$ (a) or at $\theta_x = \pi/2$ (b) to $2A_0(\sin^2\theta - B_0)\sin 2\theta$ gives A_0 and B_0 (a) or A_1 and B_1 (b), respectively. c) Fitting $\langle \partial U / \partial \theta_x \rangle$ at $\theta = \pi/2$ to $2A_{10}(\sin^2\theta_x - B_{10})\sin 2\theta_x$ gives A_{10} and B_{10} . The fitting equations are derivatives of Eqs. 20, 21, and 22. The red curves are the best fits.

Independent pH and replica-exchange simulations of model pentapeptides. The PMF function obtained from the TI simulations describes the free energy change along λ , and the difference between the two end points ($\lambda=1$ and 0) gives the deprotonation free energy. If the latter is reproduced by the CpHMD simulation, λ should sample two end (protonated and deprotonated) states with equal probabilities when pH is set to the reference pK_a value. In other words, the pK_a calculated from the titration simulation should be the same as the reference pK_a . To test it, we carried out titration simulations of model pentapeptides at 8 independent pH conditions. Three replica runs of 20 ns each were performed at each pH. The unprotonated fractions at all pH conditions are converged (see time series analysis in Fig. S1). Fitting the unprotonated fractions to the Henderson-Hasselbalch equation (Eq. 26) gave the pK_a 's of 3.4 ± 0.04 , 4.2 ± 0.02 , 6.5 ± 0.12 , 8.4 ± 0.03 , and 10.3 ± 0.01 for Asp, Glu, His, Cys, and Lys, respectively (Fig. 2). Fitting to the generalized Henderson-Hasselbalch equation gave identical pK_a 's and error estimates, but revealed a small underestimation of the Hill coefficient for all but Cys model peptides. Except for Asp, the calculated pK_a 's are within 0.1 unit of the target experimental values (Table 1). His has two titratable nitrogens and hence three protonation states: the doubly protonated Hip (charge +1) and two neutral tautomers, with a proton on either $N\delta$ or $N\epsilon$. These tautomer are respectively named Hid and Hie in Amber²⁵ or HSD and HSE in CHARMM.²⁴ The calculated pK_a 's of $N\delta$ (Hip \rightleftharpoons Hie) and $N\epsilon$ (Hip \rightleftharpoons Hid) are 7.0 ± 0.11 and 6.7 ± 0.12 , respectively. These values are also within 0.1 units from the values estimated by Tanokura based on NMR data of a model compound.⁶⁷ The titration of Asp and His is noisier than Glu, Cys, and Lys, as evident from the larger uncertainties of the unprotonated fractions at pH conditions near the pK_a value, consistent with the larger bootstrap errors (0.09 and 0.12, see Table 1). Trajectory analysis showed that the Asp and His sidechains can form hydrogen bonds (h-bonds) with the neighboring backbone group, resulting in meta-stable states. The carboxylate group of Asp is stabilized by h-bonding with the neighboring backbone amide, which contributes to the 0.3-unit under-

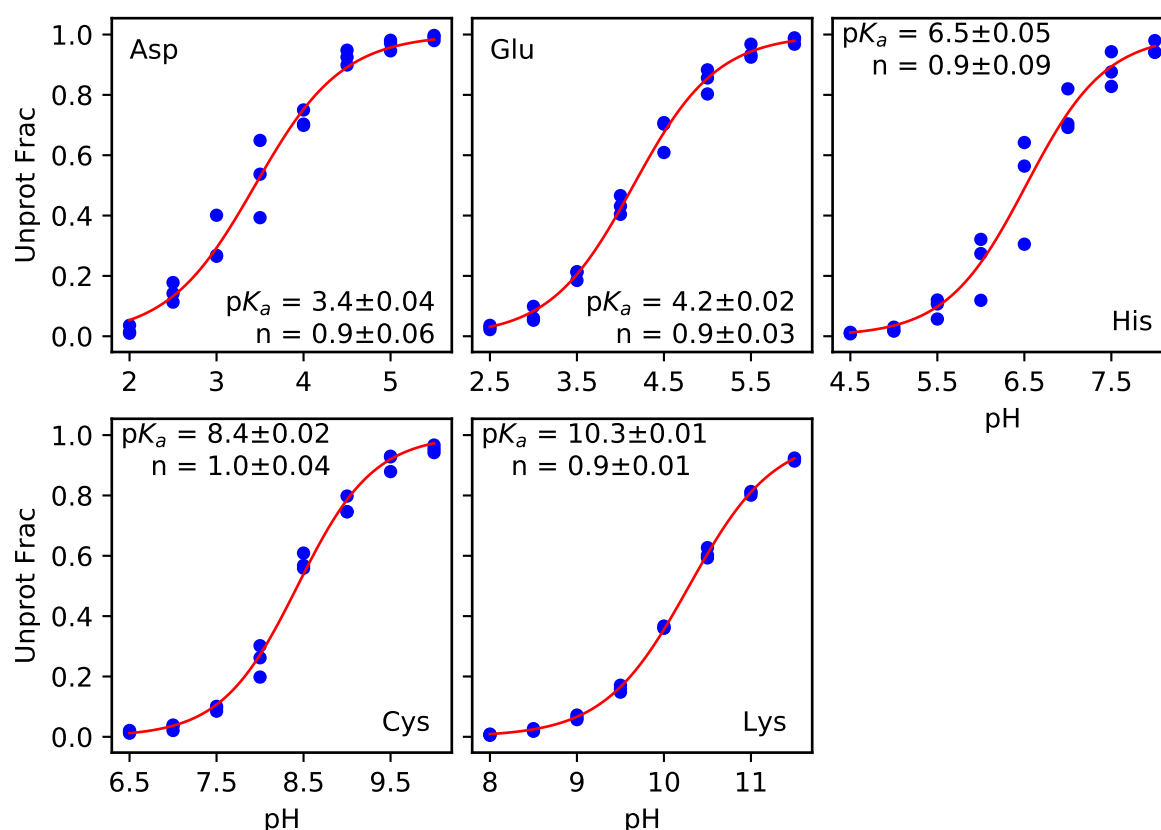


Figure 2: **Simulated titration plots of model peptides ACEAAXAANH₂ (X=Asp, Glu, His, Cys, and Lys) at independent pH conditions.** Top panel: unprotonated fractions of Asp, Glu, and His at different pH. Bottom panel: unprotonated fractions of Cys and Lys at different pH. At each pH, three simulation runs were performed starting from different initial velocity seeds. The pK_a, Hill coefficient (*n*), and fitting error are given. The boot strap errors are given in Table 2. The fitting was performed on all data points using the generalized Henderson-Hasselbalch equation. Performing the fits against the Henderson-Hasselbalch equation yields identical pK_a values and error estimates.

estimation of the target pK_a value. This behavior was previously observed in both the GB and PME-CpHMD simulations in CHARMM.^{9,11}

To investigate if the proton-coupled conformational dynamics is adequately sampled for Asp and His in the independent pH simulations, we compared the results with those from three sets of pH replica-exchange simulations. The latter were conducted with the asynchronous pH replica-exchange scheme that was recently implemented for Amber simulations.⁶⁴ The previous work of us^{11,23,43} and others^{10,15} demonstrated that the pH

replica-exchange protocol significantly accelerates protonation state and conformational sampling. Indeed, the pK_a convergence is significantly accelerated; the unprotonated fractions generally plateau after about 5 ns, compared to more than 10 ns in the independent pH simulations (Fig. S2 and S3). Interestingly, the resulting pK_a 's (3.3 and 6.5) of Asp and His are very similar to those from the independent pH titration, which suggests that sampling is sufficient in the latter (Fig. S2 and S3). Note, to account for the (~ 0.3 unit) difference between the calculated and target pK_a 's of Asp pentapeptide, we changed the Asp reference pK_a (from the experimental value of 3.7 to 4.0 in the CpHMD parameter file (charmm22_pme.parm) for protein simulations.

Table 1: Calculated and target experimental pK_a values of model pentapeptides

	Calc (IN) ^a	Calc (REX) ^b	Calc (HMR) ^c	Expt ^d
Asp	3.4 \pm 0.09	3.3 \pm 0.10	3.6 \pm 0.07	3.7
Glu	4.2 \pm 0.04		4.3 \pm 0.02	4.2
His	6.5 \pm 0.12	6.5 \pm 0.02	6.3 \pm 0.03	6.5
Hie ^e	7.0 \pm 0.11	7.1 \pm 0.02	6.9 \pm 0.01	7.0
Hid ^f	6.7 \pm 0.12	6.7 \pm 0.03	6.4 \pm 0.03	6.6
Cys	8.4 \pm 0.03		8.6 \pm 0.01	8.5
Lys	10.3 \pm 0.01	10.3 \pm 0.01	10.0 \pm 0.01	10.4

^aIndependent pH simulations, whereby each simulation was conducted for 20 ns and repeated three times. ^bThree sets of pH replica-exchange simulations of 10 ns/replica. ^cThree sets of pH replica-exchange simulations of 10 ns/replica with the HMR scheme and 4-fs timestep. All pK_a 's and errors were calculated from bootstrap. ^dExpt refers to the NMR derived pK_a 's of the model pentapeptides from Thurlkill et al.⁶⁸ The His tautomer pK_a 's are those estimated by Tanokura based on the NMR data of a model compound.⁶⁷ ^eHie refers to the pK_a associated with Hip \rightleftharpoons Hid. ^fHid refers to the pK_a associated with Hip \rightleftharpoons Hie.

In order to further accelerate simulations, we tested the sensitivities of pK_a 's for 4-fs timestep in conjunction with the hydrogen mass repartitioning (HMR) scheme.^{54,55} Three sets of pH replica-exchange simulations of 10 ns/replica were conducted for the five model peptides with HMR/4-fs timestep. All simulations converged within 5-10 ns/replica, representing a twice speed up relative to the standard 2-fs simulations. The calculated pK_a is 3.6 \pm 0.07 for Asp, 4.3 \pm 0.02 for Glu, 6.3 \pm 0.03 for His, 8.6 \pm 0.01 for Cys, and 10.0 \pm 0.01

for Lys. These pK_a 's deviate from the 2-fs simulations by 0.1–0.3 units. Notably, the pK_a 's of the basic residues are lower, by 0.2 units for His and 0.3 units for Lys. The latter is surprising, given the rapid convergence (less than 5 ns/replica) and small random error (bootstrap error of 0.01). Trajectory analysis showed that the solvent accessible surface area (SASA) of the Lys sidechain with HMR has similar pH response, i.e., slightly decreasing with pH; however, the value for all pH conditions are higher by about 4% compared to the 2-fs simulations (data not shown). This might be related to the slightly increased diffusion constant and decreased order parameter with the 4-fs timestep, as demonstrated by a recent benchmark study.⁶⁹ We note, evaluation of the 4-fs/HMR scheme for CpHMD simulations of proteins is not in the scope of the present work and will be conducted in the near future.

4.2 Titration simulations of proteins

Overall comparison of the calculated and experimental pK_a values. To test the accuracy of the PME-CpHMD method for modeling protonation states of proteins, we calculated the pK_a 's of Asp, Glu, His, and Cys residues in BBL, HEWL, SNase, RNase A, thioredoxin, and creatine kinase (HMCK) proteins, which have been previously used to benchmark CpHMD methods.^{11,43,70,71} For a total of 67 residues, the root mean square error (RMSE) and the mean unsigned error (MUE) of the calculated pK_a 's are 0.76 and 0.61, respectively, while the Pearson's correlation coefficient r is 0.85 (Figure 3). A more stringent test of the pK_a prediction accuracy is to correlate the calculated and experimental pK_a shifts (ΔpK_a) with respect to model values, as the ΔpK_a range is much smaller than the pK_a range, exposing potentially problematic cases. Encouragingly, the r value for ΔpK_a 's is 0.80, similar to the r value for absolute pK_a 's, suggesting that a good correlation with experimental is achieved and consistent for different residue types (see later discussion).

Comparison of the calculated pK_a 's with the all-atom CpHMD implementations in

CHARMM. The pK_a 's of BBL, HEWL, and SNase have been previously calculated using the all-atom PME-CpHMD implementation in CHARMM (Table S4).¹¹ The r value of ΔpK_a 's (from correlation with experiment) for these proteins from the present work is 0.80, which is nearly identical to that (0.78) using the CHARMM PME-CpHMD titration.¹¹ A comparison between the individual pK_a 's shows that most pK_a values agree within 0.2–0.3 units (Table S4); the agreement is especially remarkable for the pK_a 's of the catalytic dyad in HEWL, which differ by 0.2 units for Glu35 and are identical for Asp52.¹¹ Note, the previous CHARMM PME-CpHMD simulations were run for 10 ns/replica,¹¹ whereas the present Amber PME-CpHMD simulations were run until full convergence for 30–40 ns/replica. This analysis suggests that the pK_a drift is small over time and the replica-exchange CpHMD simulations offer pK_a calculations with good precision, consistent with our previous observations.^{11,23} We further compared the calculated pK_a 's of BBL and HEWL, which were previously reported with the MS λ D method in CHARMM (Table S4).¹⁰ Note, the MS λ D simulations in Ref. used a force-based cutoff for long-range electrostatics in λ dynamics and therefore we did not apply a finite-size correction for the pK_a 's.¹⁰ For BBL, the order of the two His pK_a 's are in agreement between the MS λ D and CpHMD results. Although the pK_a 's from MS λ D are 0.6–0.8 units higher, in better agreement with experiment, the simulation length was only 5 ns/replica and no finite-size correction was applied (which downshifts the pK_a 's). As to HEWL, the RMSE from MS λ D¹⁰ is nearly identical to the current work.

Comparison of the calculated and experimental pK_a 's of different residue types.

The Asp pK_a 's vary the most in this dataset. Encompassing both down- and upshifted pK_a 's, the experimental pK_a range for Asp is 1.2 to 8.1, similar to the calculated range of 1.5 to 7.7 (Fig. 3a, magenta). The experimental pK_a 's of Glu also display both the down- and upshifts, but the range is smaller than Asp, from 2.6 to 6.1, compared to the

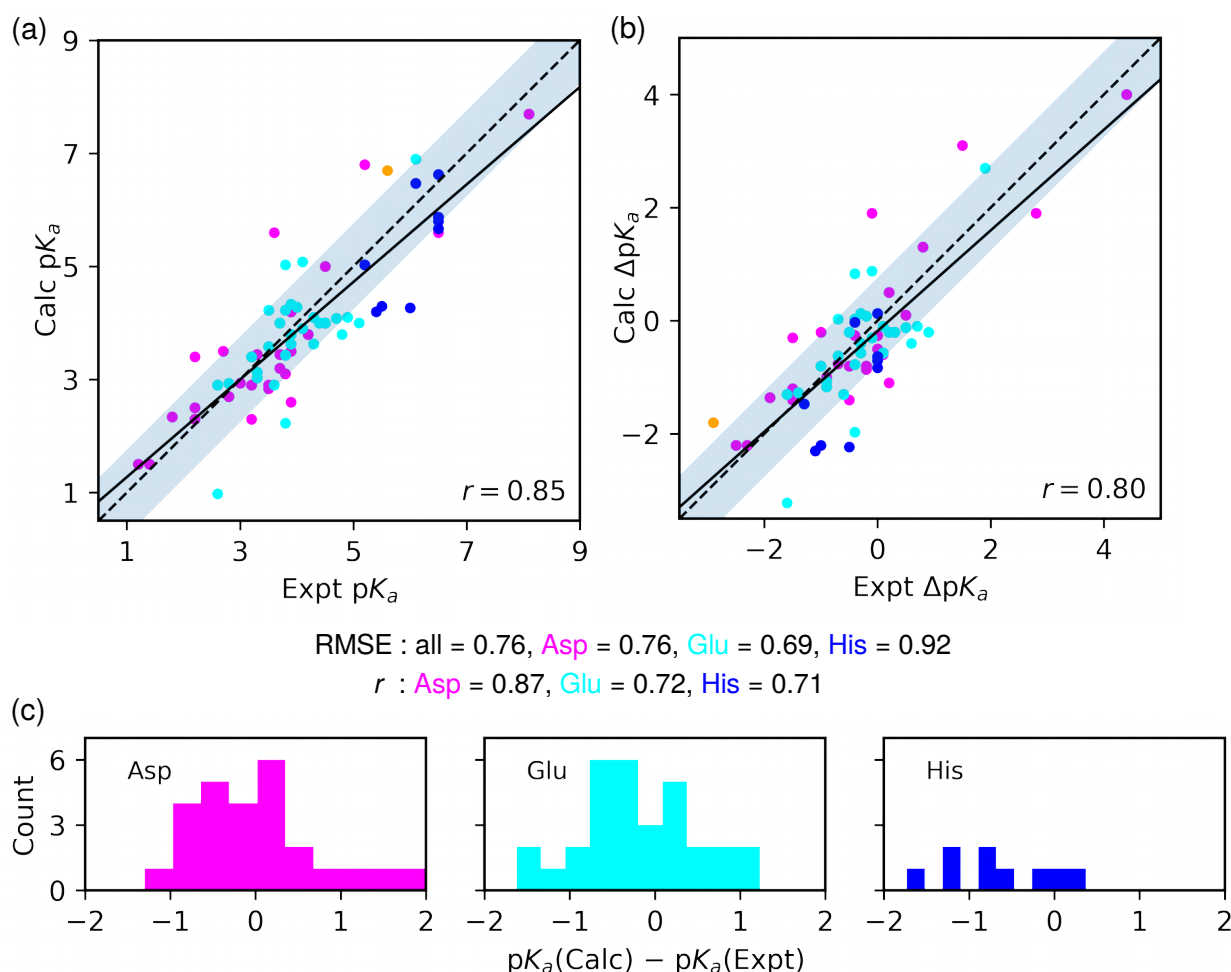


Figure 3: Comparison between the calculated and experimental pK_a 's and pK_a shifts of the benchmark proteins. a) Calculated pK_a 's vs. experimental pK_a 's. b) Calculated vs. experimental pK_a shifts with respect to the experimental model peptide pK_a 's (Table 1). The data for Asp, Glu, His, and Cys are shown in magenta, cyan, blue, and orange, respectively. Pearson's correlation coefficient (r) and RMSE are given. The solid black lines represent the linear regression. The shaded region indicates the calculated pK_a 's within the overall RMSE (0.76 units) of the experimental values. To guide the eye, the dashed diagonal line ($x=y$) is shown. c) Histograms of the deviations between the calculated and experimental pK_a 's for Asp (left), Glu (middle), and His (right) residues.

Table 2: Calculated and experimental pK_a 's of benchmark proteins^a

Residue	Expt	Calc	CHM	Residue	Expt	Calc	CHM	Residue	Expt	Calc	Residue	Expt	Calc
BBL				SNase				Thioredoxin				RNase A	
Asp129	3.9	3.5	3.7	His8	6.5	6.6	n.d.	Glu6	4.9	4.1	Glu2	2.6	1.0
Glu141	4.5	4.0	4.3	Glu10	2.8	2.9	3.2	Glu13	4.4	4.0	Glu9	4.0	4.3
His142	6.5	5.8	5.4	Asp19*	2.2	2.5	3.3	Asp16	4.2	3.8	His12	6.0	4.3
Asp145	3.7	3.2	3.4	Asp21*	6.5	5.6	6.0	Asp20	3.8	3.1	Asp14	1.8	2.3
Glu161	3.7	4.0	4.0	Asp40	3.9	2.6	2.9	Asp26	8.1	7.7	Asp38	3.5	2.8
Asp162	3.2	2.9	2.7	Glu43	4.3	3.6	4.1	His43	n/d	5.4	His48	6.1	6.5
Glu164	4.5	4.0	4.3	Glu52	3.9	4.3	4.7	Glu47	4.3	4.1	Glu49	4.7	4.1
His166	5.4	4.2	4.1	Glu57	3.5	4.2	4.1	Glu56	3.2	3.4	Asp53	3.7	3.4
RMSE		0.62	0.66	Glu67	3.8	3.4	4.0	Asp58	5.2	6.8	Asp83	3.3	3.4
HEWL				Glu73	3.3	3.0	3.6	Asp60	2.7	3.5	Glu86	4.1	5.1
Glu7	2.6	2.9	3.2	Glu75	3.3	3.1	2.7	Asp61	3.9	4.2	His105	6.5	5.9
His15	5.5	4.3	4.0	Asp77	<2.2	-0.2	<-0.0	Asp64	3.2	2.3	Glu111	3.5	3.6
Asp18	2.8	2.7	2.9	Asp83	<2.2	0.3	0.0	Glu68	5.1	4.0	His119	6.5	5.7
Glu35	6.1	6.9	7.1	Asp95	2.2	3.4	3.0	Glu70	4.8	3.8	Asp121	3.0	2.9
Asp48	1.4	1.5	0.9	Glu101	3.8	4.2	4.7	Glu88	3.6	2.9	RMSE		0.81
Asp52	3.6	5.6	5.6	His121	5.2	5.0	/	Glu95	4.1	3.9	HMCK		
Asp66	1.2	1.5	1.1	Glu122	3.9	3.6	4.4	Glu98	3.9	3.8	Cys283	5.6	6.7
Asp87	2.2	2.3	2.3	Glu129	3.8	5.0	5.5	Glu103	4.5	4.0			
Asp101	4.5	5.0	5.2	Glu135	3.8	2.2	2.9	RMSE		0.71			
Asp119	3.5	2.9	3.5	RMSE		0.76	0.80	All		Max	RMSE	MUE	
RMSE		0.83	0.92						2.0	0.76	0.61		

^aCHM column contains the pK_a 's from the CHARMM PME-CpHMD simulations with finite-size corrections for Asp, Glu, and His (~ -0.5 for BBL; ~ -0.9 for HEWL/SNase).¹¹ Experimental data are taken from ref^{72,73} for BBL, ref⁷⁴ for HEWL, ref⁵⁸ for SNase, ref⁷⁵ for thioredoxin, ref⁷⁶ for RNaseA, and ref⁷⁷ for HMCK. Glu56 of thioredoxin has two reported pK_a 's, 3.2 and 5.1, and former was used to calculate the error. For coupled residues (indicated by an asterisk), the macroscopic stepwise pK_a 's (from experiment and simulations) are listed.

calculated range of 1.0 to 6.9 (Fig. 3a, cyan). The overall accuracy of the pK_a calculation for Asp is slightly worse than Glu (RMSE of 0.76 and 0.69 respectively), but the r value for Asp (0.87) is somewhat larger than for Glu (0.72), which may be attributed to the larger pK_a range. There are only 9 experimental pK_a 's for His in the current dataset, which has a range of 5.2–6.5 and do not include upshifted values (Fig. 3a, blue). The calculated His pK_a range is 4.2–6.6 (Fig. 3a, blue), and there is a trend of systematic overestimation of pK_a downshifts (Fig. 3c, right). In contrast, there is no clear trend for the pK_a errors of Asp and Glu (Fig. 3c, left and middle). The RMSE (0.92) for His pK_a 's is larger than those for Asp (0.76) or Glu (0.9).

pK_a calculation for BBL: pH-dependent solvent exposure of His166. BBL is a mini-protein with 45 residues and 8 titratable sites. The RMSE of the calculated pK_a 's is 0.62 units, with His166 showing the largest error of 1.2 units, representing an overestimation to the experimentally observed pK_a downshift (Table 2). The pK_a downshift of His166 can be attributed to solvent exclusion and lack of hydrogen bonding (h-bonding) or electrostatic interactions (Fig. 4a). As pH decreases from 6 to 4, His166 undergoes a sigmoidal transition (Fig. 4b) from the deprotonated fraction of 1 (singly protonated neutral state) to 0 (doubly protonated charged state). As expected, the fraction of the fully buried state also decreases (i.e., solvent exposure increases); however, the decrease does not appear to be sufficient, i.e., at low pH values the buried fraction does not plateau (Fig. 4c). This analysis suggests that while the PME-CpHMD method is able to reproduce the experimental pK_a downshift by capturing the pH-induced decrease in solvent exclusion or increase in solvent exposure of His166, sampling of the exposed state at low pH may be insufficient, which contributes to the pK_a underestimation. Another potential contributor is an overestimation of the desolvation penalty by the CHARMM c22 force field.⁴⁴ Note, our previous CHARMM PME-CpHMD simulations gave a similarly underestimated pK_a for His166 (by 1.3 units),¹¹ and the MS λ D simulations underestimated the pK_a by 0.6 units (analysis of conformational sampling was not given).¹⁰

pK_a calculation for HEWL: titration order of the catalytic dyad. HEWL is a small protein with 129 residues and 10 titratable sites; it is a popular test system for pK_a prediction methods due to the abundance of experimental data.⁷⁴ The RMSE of the calculated pK_a 's is 0.83 units. The two largest errors are for Asp52 and His15; the calculated pK_a 's are 2 units over- and 1 unit underestimated, respectively (Table 2). Despite the overestimation, the calculated pK_a of Asp52 is 1.4 units lower than that of Glu35 (the second catalytic residue, Fig. 5), which indicates that Glu35 is a general acid and Asp52 is a general base in catalysis, in agreement with experiment (Table 2). Consistent with the

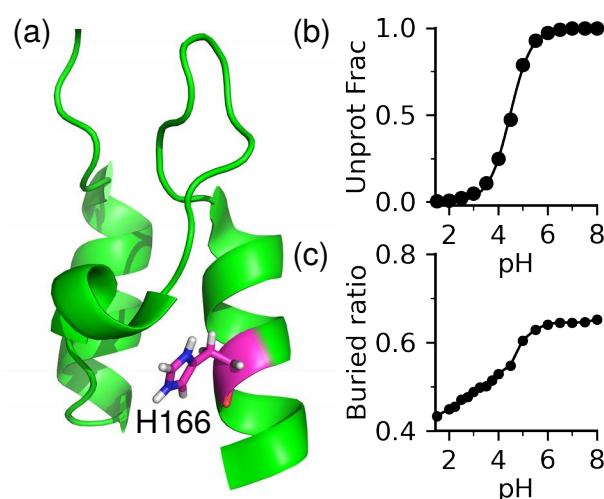


Figure 4: Protonation of His166 in BBL is correlated with the pH-dependent increase in solvent exposure. a) Deprotonated fraction of His166 at different pH. b) Buried ratio of His166 at different pH, defined as $1 - fSASA$. $fSASA$ (fraction of solvent accessible surface area) was calculated as $SASA$ of the sidechain atoms relative to that in the model pentapeptide.

CHARMM PME-CpHMD as well as the hybrid-solvent CpHMD simulations,⁷⁸ the titration events of Glu35 and Asp52 are uncoupled, as evident from the nearly identical stepwise pK_a 's (7.0 and 5.5) from fitting to the two-proton titration model (Eq. 27; figure not shown). The lack of coupled titration is due to the relatively large distance between the carboxylate sidechains (>6.5 Å between the nearest carboxylate oxygens at any pH). Note, the calculated dyad pK_a 's are nearly the same as the values from the CHARMM PME-CpHMD simulations.¹¹ The $MS\lambda D$ method in CHARMM gave a nearly identical pK_a for Glu35, but a 1.1-unit lower pK_a for Asp52.¹⁰

To understand the pK_a order of the catalytic dyad Glu35/Asp52 and the possible factors for the pK_a overestimation of the latter, we compared the pH profiles of the solvent exposure, h-bonding and electrostatic interactions of the dyad residues with the pH-dependent titration curves (Fig. 5b–d). Both residues are partially buried. As Glu35 switches from being fully unprotonated to fully protonated in the pH range 7 to 10, the solvent exposure decreases from about 30% to just under 20% (Fig. 5b and c, blue). Asp52

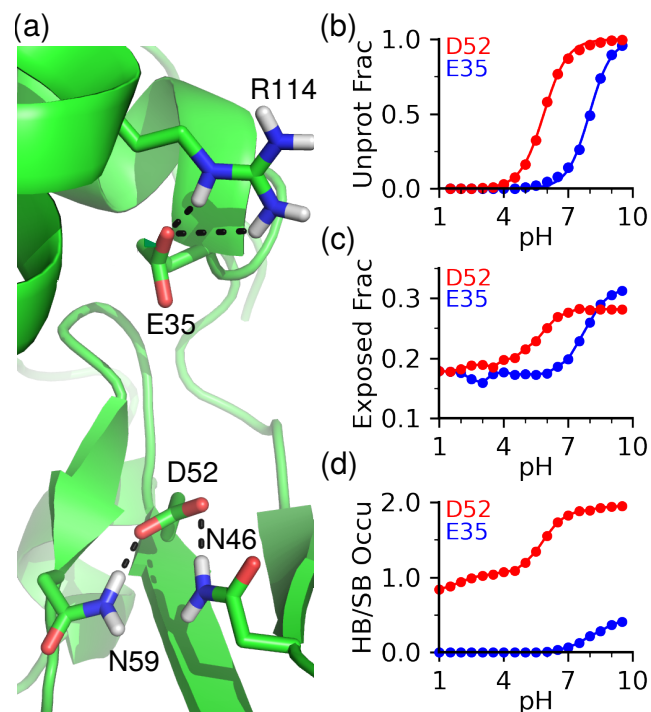


Figure 5: **Factors influencing the pK_a 's of the catalytic dyad in HEWL.** (a) A representative snapshot at pH 7.5 showing the h-bonding and salt bridge environment of Glu35 and Asp52. (b) The unprotonation fractions of Glu35 and Asp52 at different pH. (c) Fractions of the sidechain solvent exposure (SASA value relative to the model pentapeptide) of Glu35 and Asp52 at different pH. (d) The h-bond and salt bridge occupancies of Glu35 and Asp52 at different pH.

has a similar behavior, except that the titration and change in solvent exposure are shifted to a lower pH range 7 to 4 (Fig. 5b and c, red). Now we turn to h-bonding and electrostatic interactions that are also physical determinants of pK_a shifts.⁷⁹ Glu35 does not form h-bonds below pH 7.5, and above pH 7.5, occasional salt-bridge interaction with Arg114 was observed, with an occupancy less than 20% (Fig. 5d, blue). In contrast, the carboxylate group of Asp52 can accept h-bonds from the sidechains of Asn46 and Asn59, and the occupancy increases to nearly 2 with increased deprotonation of Asp52 (Fig. 5d, red). The analysis of solvent exposure and h-bond suggests that the latter is the major determinant for the lower pK_a value of Asp52 relative to Glu35, in agreement with our previous work using the hybrid-solvent as well as the PME-CpHMD in CHARMM.⁷⁸ It is noteworthy that despite the correlation between charging of Asp52 and the increase in solvent exposure and h-bond formation, the pH profiles of solvent exposure and h-bond occupancy are more gradual than the titration curve, which may indicate that the pH-dependent conformational changes might be undersampled, contributing the overestimation of the pK_a of Asp52.

pK_a calculation for SNase: partially buried residues and coupled titration of Asp19 and Asp21. SNase has a large number of engineered mutants, which are popular model systems for testing pK_a prediction methods.⁷⁹ Here we used the hyperstable, acid-resistant form of SNase Δ +PHS (hereafter referred to as SNase)⁵⁸ which is only slightly (14 residues) larger than HEWL, but has 9 more titratable sites. The pK_a 's of SNase are more challenging to predict than HEWL due to the fact that most titratable sites are partially buried.²³ The calculated pK_a 's have a RMSE of 0.76 units, similar to the RMSE of 0.80 from the CHARMM PME-CpHMD simulations.¹¹ The largest error is for Glu129, for which the experimental pK_a is 0.4 units down- and the calculated pK_a is 0.8 units upshifted relative to the model value of 4.2. Curiously, simulation also fails to reproduce the direction of the experimental pK_a shifts of Glu52, Glu57, and Glu101, although the

magnitude of the errors is smaller (0.4, 0.7, and 0.4 respectively). Analysis showed that all these residues are partially buried, suggesting that desolvation penalty contributes to the pK_a upshift. Based on the analysis of BBL's His166 and HEWL's Asp52, we hypothesized that simulation overestimates the desolvation penalty due to inadequate sampling of the solvent exposed state. To test this, we plotted the fractional SASA values vs. pH for Glu52, Glu57, Glu101, and Glu129 (Fig. S4). Deprotonation of glutamic acid is expected to induce larger solvent exposure. This is indeed the case for the more exposed residues Glu52 and Glu57 (fractional SASA about 60% at low pH), although the degree of increase is small. However, solvent exposure change with pH for the more buried residues Glu101 and Glu129, for which the fractional SASA values remain at about 40 and 20% throughout the entire pH range (Fig. S4). These data support the hypothesis that the solvent exposed state may be inadequately sampled, contributing to the desolvation related pK_a upshift for Asp and Glu.

The NMR data⁵⁸ as well as our previous work⁷⁸ based on the hybrid-solvent and PME-CpHMD simulations in CHARMM suggest that the titration Asp19 and Asp21 is coupled. The current simulations confirmed the strong coupling as a result of h-bond formation between the two residues (Fig. 6a). Fitting the titration data to a two-proton coupled equation (Eq. 27) gives the stepwise macroscopic pK_a 's of 2.5 and 5.6 (Fig. 6b), which are in good agreement with the experimental values of 2.2 and 6.5.⁵⁸ To assign the stepwise pK_a 's to individual residues, we examine the pH-dependent probabilities of four microscopic states, doubly protonated (HH), singly protonated with proton on D19 (H-) or Asp21 (-H), and doubly deprotonated (—) states (Fig. 6c). Above pH 7, Asp19/Asp21 are in the — state (Fig. 6c, blue). As pH decreases from 7 to 5, the probability of — decreases, while that of the -H or H- increases. Since the -H state (cyan) is more probable than the H- state (magenta) as protonation first occurs, Asp21 receives a proton first, which means the higher pK_a should be assigned to Asp21. As pH further decreases from 5 to 2, both -H and H- states are possible; however, their combined probability

decreases, while the probability of the HH state increases (Fig. 6c, red). Below pH 2, the latter state dominates. Analysis of h-bonding and electrostatic interactions (Fig. 6a) shows a network among Asp19, Asp21, Thr22, Thr41, and Arg35, consistent with the CHARMM hybrid-solvent and PME-CpHMD simulations.⁴³

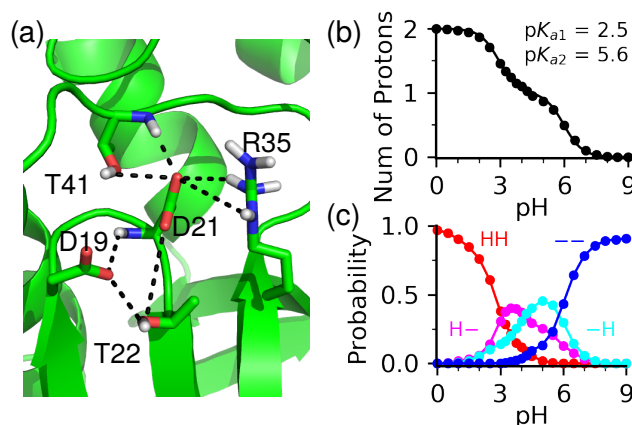


Figure 6: **Linked titration of Asp19 and Asp21 in SNase.** (a) A snapshot from the pH 4 simulation showing the h-bonding environment of Asp19 and Asp21. (b) Total number of protons of Asp19/Asp21 at different pH. The stepwise pK_a 's are obtained from the best fit (black curve) to the two-proton coupled equation (Eq. 27). (c) The pH-dependent probabilities of four microscopic states: two protons (HH, red); proton on Asp19 (H⁻, magenta) or Asp21 (⁻H, cyan); zero proton (⁻, blue).

pK_a calculation for thioredoxin: the deeply buried Asp26 and Asp58. Thioredoxin has 105 residues with 18 titratable sites. The RMSE of the calculated pK_a 's is 0.71. We first consider Asp26, which has one of the highest measured pK_a 's of any carboxylic acids in proteins. Encouragingly, the calculated pK_a of Asp26 is 7.7, in excellent agreement with the NMR-derived value of 8.1. The large pK_a upshift of nearly 4 units can be attributed to the extremely low fraction of solvent exposure (below 7% at all pH, Fig. S5). Trajectory analysis showed that Asp26 does not form h-bonds with nearby residues. The only factor that may stabilize the deprotonated form is the salt-bridge formation with Lys39; however, the salt bridge is only formed above pH 8 and the solvent exposure is very low (<20% at pH 8, Fig. S5). A previous experimental study⁷⁵ suggested that the protonated Asp26

may be stabilized by donating a h-bond to the nearby Ser28; however, in the simulation the average distance from the hydroxyl oxygen of Ser28 to the nearest carboxylate oxygen of Asp26 is about 4.9 Å at pH 7, similar to the distance of 4.7 Å in the X-ray structure (PDB 1ERU). Thus, the dry environment, along with lack of polar interactions, results in the very large pK_a upshift of Asp26.

The largest error in the calculated pK_a 's of thioredoxin is for Asp58, whose direction of pK_a shift is reproduced but the magnitude is 1.6 units too large (Table 2). Analysis showed that the Asp58 is also deeply buried, with ~20% solvent exposure below pH 5, which explains the pK_a upshift relative to the model. However, the solvent exposure only slightly increases to ~30% at pH 8 before increasing steeply to over 50% at pH 10 (Fig. S6). H-bond analysis showed that the deprotonated Asp58 can accept h-bonds from the backbones of neighboring Asp60 and Asp61, which can stabilize the deprotonated state; however, the pH profile of the h-bond occupancy is irregular, showing a nearly 50% decreased occupancy in the pH range 4–8 (Fig. S6). The latter indicates a sampling issue, which explains the overestimation of the pK_a of Asp58.

pK_a calculation for RNase A: the deeply buried His12. RNase A has 124 residues with 14 titratable sidechains. The RMSE for the calculated pK_a 's is 0.81, and the largest error is for His12 (Fig. 7a). The experimental pK_a of His12 is 0.5 units downshifted relative to the model, and the simulation overestimated the downshift by 1.7 units (Table 2). Analysis showed that His12 titrates over the pH range 3 to 6 (Fig. 7b), and the titration is correlated with two physical determinants, an increase in solvent exposure (decreased buried fraction) at lower pH (Fig. 7c) and an increase in h-bond formation at higher pH (Fig. 7d). However, unlike in the previous GB- or hybrid-solvent CpHMD simulations,^{19,43} the pH profiles of the buried fraction and the h-bond occupancy do not fully match the titration curve. Above pH 6, His12 is over 90% buried, and the decrease in the buried fraction at lower pH does not follow an expected sigmoidal curve. The buried fraction

decreases by about 5% as pH decreases from 9 to 5 and remains constant between pH 4.5 and 2, before further decreasing to 80% at pH 0 (Fig. 7c). A major h-bond partner is the neighboring Asn11, which can donate a h-bond from its carboxamide group to the ϵ nitrogen of His12 to stabilize its deprotonated form (Fig. 7a). As pH increases from 3 to 6, the occupancy of the h-bond increases from zero to about 60%, and it further increases to nearly 100% at pH 8. Based on the above analysis, we suggest that the h-bond formation and solvent exposure of His12 may be insufficiently sampled in the pH range 3–6. The under-sampling of the solvent exposed state (buried fraction remains unchanged between pH 2 and 4.5) is particularly evident, which may be a major factor for the overestimation of pK_a downshift of His12.

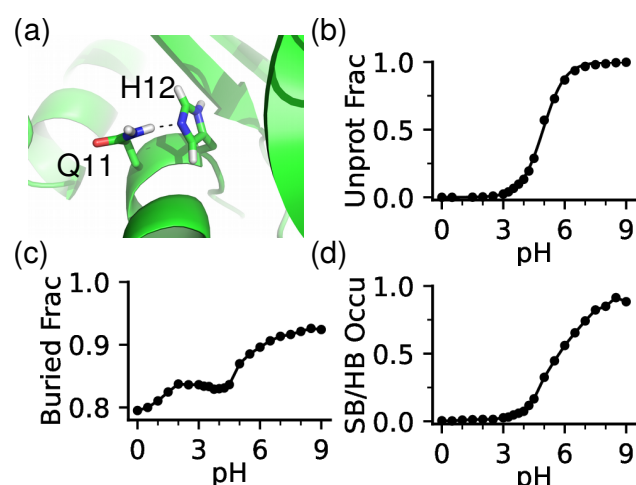


Figure 7: Protonation of His12 in RNaseA is correlated with the decreased solvent exclusion and hydrogen bonding. a) A zoomed-in view of the hydrogen bonding between His12 and Asn11 in RNase A. The snapshot was taken from the simulation at pH 7. b) Unprotonated fraction of His12 at different pH. c) Buried fraction of His12 at different pH. Definition of the buried fraction is given in the caption of Fig. 4. d) Occupancy of the h-bond between His12 and Asn11 at different pH.

pK_a calculation for HMCK: a buried active-site cysteine. To test the accuracy of Cys titration, we calculated the pK_a of Cys283 in the active site of HMCK, which has a NMR measured pK_a of 5.6,⁷⁷ one of the lowest in the literature.⁸⁰ Our simulations correctly reproduced the direction of the pK_a shift relative to the model; however the downshift is

1.1 units underestimated compared to the experiment (Table 2). Analysis showed that Cys283 is buried and does not have nearby cationic residues; however, once deprotonated it can accept h-bonds from the sidechains and backbones of Ser285 and Asn286 (Fig. 8a and b), consistent with the GB-based CpHMD titration simulation.⁷¹ Based on the structural analysis of the thioredoxin family of proteins,⁸¹ Roos and Messens hypothesized that hydrogen bonding rather than electrostatics plays a major role in stabilizing Cys thiolates. Our current data and recent GB-based CpHMD simulations of a large number of proteins^{80,82,83} are in support of this hypothesis.

As Cys283 becomes deprotonated in the pH range 6 to 9, the total h-bond occupancy increases and plateaus at 1; however, the exposed fraction does not increase and instead remains at about 40% (Fig. 8c and d). Since solvent exposure promotes the charged thiolate state and decreases the pK_a , we suggest that insufficient sampling of the solvent-exposed conformations may contribute to the overestimation the pK_a of Cys283.

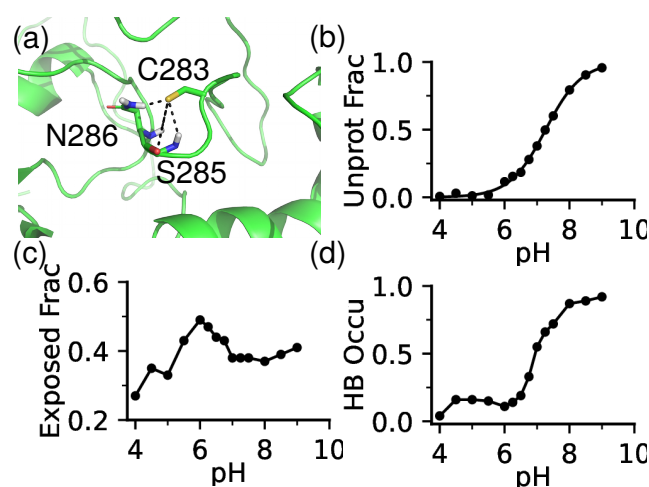


Figure 8: Factors influencing the pK_a of Cys283 in HMCK. (a) A zoomed-in view of the h-bond environment of Cys283 (from simulation at pH 7.5). Cys283 thiolate can form h-bonds with the backbones and sidechains of Ser285 and Asn286. (b) Unprotonated fraction of Cys283 at different pH. (c) Exposure fraction (SASA relative to that of the model pentapeptide) at different pH. (d) Occupancy of the total h-bond formation of Cys283 thiolate with Ser285 and Asn286 at different pH.

Finite-size effect and corrections. Following the work of Hünenberger and colleagues,⁵⁰ we previously proposed an analytical pK_a correction (Eq. 25) to correct for the effect of an offset potential introduced in PME simulations under periodic boundaries.¹¹ For the current simulations, the pK_a corrections for Asp, Glu, His, and Cys are between -0.3 and -1.0 pH units (see Methods). To assess the effectiveness of the corrections and better understand the finite-size effect, we performed additional titration simulations of SNase with increased box sizes, i.e., adding more water to the simulation system. Table 3 summarizes the raw and corrected pK_a 's using four different boxes, which have 10 (default), 12, 14, or 18 Å cushion space between the protein and edges of the box (minimum distance between the heavy atoms of protein and water oxygens on the box edges). The corresponding cubic box lengths are 68, 71, 76, and 84 Å, respectively.

We first examine the raw calculated pK_a 's from simulations with different box sizes. As expected, with increasing box size the raw pK_a 's decrease for all but four residues (Fig. 9a). Increasing box size also leads to better agreement with the experimental pK_a 's; the RMSEs of the raw pK_a 's are 1.0, 1.0, 0.97 and 0.76 for boxes with 10, 12, 14, and 18 Å cushion space, respectively (Table 3 and Fig. 9a). The MUE also decreases from 0.86 (10 Å cushion) to 0.81 (12 Å cushion), 0.80 (14 Å cushion), and 0.62 (18 Å cushion) (Table 3). Comparison of the raw pK_a 's between the smallest (10 Å cushion) and largest (18 Å cushion) boxes shows that the pK_a changes due to box size increase vary (Table 3, last column). Excluding the four residues (Asp19, Asp21, H121, and Glu135) that show very little pK_a changes, the pK_a 's mostly decrease by 0.3 to 0.7 units, as compared to the finite-size corrections of -0.70 to -0.8 units for the smallest box. The effect of box size is not clear for the coupled residues Asp19/Asp21, which have the raw calculated pK_a 's of 3.2/6.3 with the smallest box; however, the pK_a 's increase to 3.6/6.5 and 3.4/7.1 with the larger boxes (12 and 14 Å cushion space), and then decrease back to 3.1/6.4 with the largest box (18 Å cushion). Increasing box size has negligible effect on the downshifted pK_a of His121. With the increasing box sizes, its raw pK_a changes from 5.7 to 6.1,

5.7, and 5.8. Box size also shows little effect on the pK_a of Glu135, which has the raw calculated pK_a 's of 3.0, 3.4, 3.0, and 3.0 with the increasing box sizes.

Now we examine the pK_a corrections for the different simulation boxes. It is apparent that application of the finite-size correction removes the systematic overshift error (Fig. 9b). As the box size increases, the solvent number density increases and therefore the magnitude of the correction decreases (Eq. 24). The magnitude of the corrections decreases by about 0.4 pH units going from the smallest to the largest box. Interestingly, this difference is roughly the same as the average difference between the raw pK_a 's (of all but the aforementioned four residues) calculated with the smallest and largest box (Table 3, last column), which suggests that the finite-size correction is valid. Another interesting observation is that the increasing box size does not significantly reduce the RMSE of the finite-size corrected pK_a 's. The RMSE's are 0.76, 0.80, 0.80, and 0.70 with the increasing box sizes (Fig. 9b), which is another piece of evidence supporting the validity of the finite-size corrections.

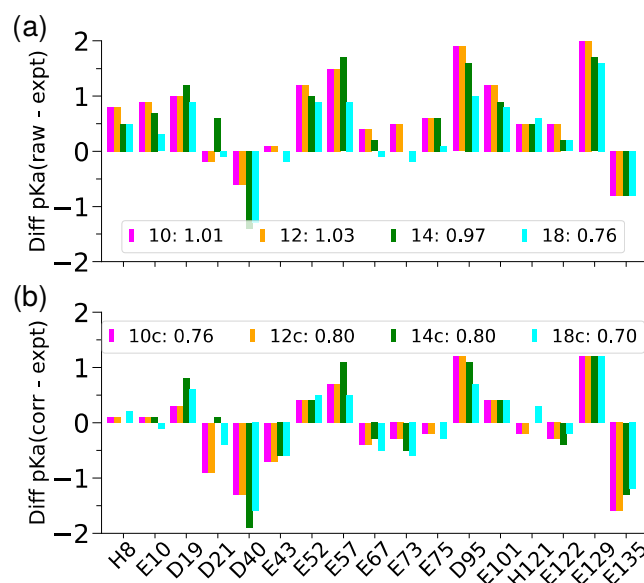


Figure 9: **Effect of box size on the calculated pK_a 's.** The errors of the raw (a) and finite-size corrected (b) pK_a of SNase with different solvent cushion spaces, 10 (magenta), 12 (orange), 14 (green), and 18 Å (cyan). The corresponding RMSE values are shown next to the legends.

Table 3: Effect of simulation box size on the calculated pK_a values of SNase^a

Residue	Expt	Raw	Corr	Raw	Corr	Raw	Corr	Raw	Corr	Δ Box
L_{cus} (Å)		10		12		14		18		
L_{box} (Å)		67.6		70.6		75.5		83.5		
H8	6.5	7.3	6.6	7.1	6.5	7.0	6.5	7.0	6.7	-0.3
E10	2.8	3.7	2.9	3.5	2.8	3.5	2.9	3.1	2.7	-0.6
D19	2.2	3.2	2.5	3.6	3.0	3.4	3.0	3.1	2.8	-0.1
D21	6.5	6.3	5.6	6.5	5.9	7.1	6.6	6.4	6.1	0.1
D40	3.9	3.3	2.6	2.9	2.3	2.5	2.0	2.6	2.3	-0.7
E43	4.3	4.4	3.6	4.4	3.7	4.3	3.7	4.1	3.7	-0.3
E52	3.9	5.1	4.3	5.2	4.5	4.9	4.3	4.8	4.4	-0.3
E57	3.5	5.0	4.2	5.1	4.4	5.2	4.6	4.4	4.0	-0.6
E67	3.8	4.2	3.4	4.1	3.4	4.0	3.5	3.7	3.3	-0.5
E73	3.3	3.8	3.0	3.6	3.0	3.3	2.8	3.1	2.7	-0.7
E75	3.3	3.9	3.1	3.3	2.6	3.9	3.3	3.4	3.0	-0.5
D77	<2.2	0.5	-0.2	0.6	0.0	0.5	-0.1	0.4	0.1	-
D83	<2.2	1.0	0.3	2.1	1.5	<0.0	<0.0	<0.0	<0.0	-
D95	2.2	4.1	3.4	4.1	3.5	3.8	3.3	3.2	2.9	-0.9
E101	3.8	5.0	4.2	4.9	4.2	4.7	4.2	4.6	4.2	-0.4
H121	5.2	5.7	5.0	6.1	5.6	5.7	5.2	5.8	5.5	0.1
E122	3.9	4.4	3.6	4.1	3.4	4.1	3.5	4.1	3.7	-0.3
E129	3.8	5.8	5.0	5.8	5.1	5.5	5.0	5.4	5.0	-0.4
E135	3.8	3.0	2.2	3.4	2.7	3.0	2.5	3.0	2.6	0.0
RMSE		1.0	0.76	1.0	0.80	0.97	0.80	0.76	0.70	
MUE		0.86	0.61	0.81	0.67	0.80	0.60	0.62	0.58	

^a Box size is represented by the solvent cushion space (L_{cus}), i.e., minimum distance (10, 12, 14, and 18 Å) between the protein heavy atoms and edges of the water box, and the length of a cubic box (L_{box}) converted from the average system volume. The columns Raw and Corr refer to the pK_a 's before and after the finite-size corrections (see Methods). The column Δ Box refers to the raw pK_a difference between the largest and smallest boxes.

5 Concluding Discussion

We presented the first implementation, parameterization, and validation of the GPU-accelerated continuous constant pH particle-mesh Ewald molecular dynamics method in Amber22 (hereafter referred to as Amber PME-CpHMD). Titration parameters for three force fields (CHARMM c22,⁴⁴ Amber ff14SB,⁴⁶ and ff19SB⁵²) were derived and validated using model pentapeptides AAXAA, where X represents Asp, Glu, His, Cys, or Lys. To benchmark the performance and accuracy for constant pH simulations of proteins, we carried out titration simulations with the c22 force field for 6 proteins, including BBL, HEWL, SNase, RNase A, thioredoxin, and HMCK, which have NMR derived pK_a values of Asp, Glu, His, and Cys residues. The asynchronous pH replica-exchange algorithm⁶⁴ was employed to enhance sampling of protonation and conformational states. The simulations were run for 30-50 ns per pH replica until all pK_a 's were converged. The resulting RMSE and MUE with respect to the experimental pK_a 's are 0.76 and 0.61, respectively, and the largest pK_a deviation is 2 units. The Pearson's correlation coefficients for the calculated vs. experimental pK_a 's and pK_a shifts are 0.85 or 0.80, respectively. Importantly, the titration simulations quantitatively reproduced the experiment pK_a orders of the catalytic dyad in HEWL and the coupled residues in SNase. Simulations also quantitatively captured one of the largest upshifted pK_a 's of a deeply buried Asp in thioredoxin as well as the downshifted pK_a of an active-site Cys in HMCK.

We compared the current validation data with those based on the CHARMM²⁴ CPU all-atom PME-CpHMD¹¹ and MS λ D¹⁰ simulations with the same c22 force field.⁴⁴ The Asp, Glu, and His pK_a 's calculated from the CHARMM PME-CpHMD simulations of 10 ns per replica (much shorter than the present work) are in close agreement with the present work, suggesting that the pK_a drifts over prolonged simulation time are small. Comparing to the calculated pK_a 's of HEWL and the two His residues in BBL based on the MS λ D simulations of 5-20 ns per pH replica (with a 12-Å electrostatic cutoff),¹⁰ the overall RMSE

is similar, and the pK_a orders of the catalytic Glu35/Asp52 in HEWL and His142/His166 in BBL are consistent with the present simulations.

In agreement with the previous CHARMM PME-CpHMD simulations¹¹ the present data demonstrated that the finite-size effect needs to be taken into account for the accurate calculation of titration free energies with lattice sum methods under periodic boundary conditions. Applying the pK_a correction¹¹ to account for a positive offset potential due to TIP3P water in periodic boxes, a systematic pK_a upshift error in the calculated pK_a 's was removed, and the overall agreement with experiment was improved. We note, in the revision stage of the current paper, the work from the Roux group⁸⁴ was published which used a similar pK_a correction to account for the (Gavani) offset potential.⁵¹

To further examine the finite-size effect and the validity of the correction, the pK_a 's of SNase were calculated from simulations with four different box sizes. Consistent with the negative sign of the correction, increasing box size lowers the raw pK_a 's of all but four residues that do not show significant changes. The RMSE of the raw pK_a 's decreases from 1.0 with the smallest box to 0.76 with the largest box; the latter is identical to the RMSE (0.76) of the corrected pK_a 's obtained from the simulation with the smallest box. The quantitative validity of the correction is also supported by a good agreement between the change in the finite-size correction and the average change of the raw pK_a 's going from the smallest to the largest box size. As expected, the finite-size correction decreases with increasing box size, and consistently, the reduction in RMSE due to the correction also decreases. Using the 18-Å cushion space, the correction is 0.3–0.4, and the RMSE (0.70) of the corrected pK_a 's is only slightly smaller than the RMSE (0.76) of the raw pK_a 's. This suggests that the box size effect may start to become negligible with this size of water box.

Although the overall box-size dependent trend of pK_a 's is consistent with the positive offset potential being the dominant factor,¹¹ there are exceptions. The simulations of SNase showed that box size has negligible effect on the coupled pK_a 's of Asp19 and

Asp21 as well as the downshifted pK_a 's of His121 and Glu135. We note that the effect of the offset potential and the corresponding pK_a correction deal with an ideal situation in which a single residue titrates in a neutral background. Thus, it is possible that the correction is not valid for coupled pK_a 's. However, with regards to the pK_a 's of His121 and Glu135, the cause for the box size independence is difficult to speculate. An alternative approach to the finite-size pK_a correction is to enforce system charge neutrality i.e., by including titratable water as in our previous work.¹¹ We tested this approach on the BBL protein; however, due to the slower convergence and small pK_a differences compared to the simulations without titratable water, studies of other proteins were not pursued. We defer a more thorough investigation of the finite-size effects to a future work.

We analyzed the pH-dependent solvent exposure and formation of hydrogen bonds as well as electrostatic interactions of catalytic residues and those that exhibit larger pK_a deviations from experiment. These analyses suggested while PME-CpHMD captures the proton-coupled conformational rearrangements, charging-induced increase of solvent exposure for buried residues is inadequate. This may be a major contributor to the pK_a errors, including the overestimated pK_a downshifts for buried His residues, e.g., His166 in BBL and His12 in RNaseA; the overestimated pK_a upshifts for buried carboxyl residues, e.g., Glu57 in SNase and Asp58 in thioredoxin; and the overestimated pK_a upshift for buried Cys, e.g., Cys283 in HMCK. Undersampling of the solvent-exposed state may also be related to the combination of c22/TIP3P force field,⁴⁴ which slightly biased solute-solute over solute-solvent interactions.⁸⁵ Overestimation of desolvation penalty may also be a source of error, which can be attributed to the low dielectric constant in the protein interior as a result of the lack of polarization in simulations with additive force fields.⁸⁶ Lack of polarization in the interior of protein may also lead to overly strong salt bridges, which may explain the overestimation of the pK_a downshifts of Asp140 and Glu135 in SNase. While the use of polarizable force field for both protein and water is desirable, it may not be currently feasible due to speed. One intriguing idea worth exploring is to

mix a polarizable water model such as OPC3-pol⁸⁷ with an additive force field to improve solute-solvent interactions. The present study did not examine the potential dependence on the additive force field. The force field related topics as well as the evaluation of PME-CpHMD for model proton-coupled conformational dynamics of catalytic residues in larger proteins (e.g., BACE1⁷⁰) will be explored in a future work.

By removing the reliance on the implicit-solvent model, the PME-CpHMD method can be applied to any system that has a force field representation. We anticipate the GPU accelerated PME-CpHMD to become a powerful tool for the investigation of a variety of proton-coupled dynamical phenomena that are poorly understood due to the current limitations in experimental and MD techniques, for example, secondary transport of ions/substrates across membrane transporter proteins and pH-dependent self-assembly of materials. Another important application of PME-CpHMD is to offer proper pH control, for example, by allowing protein and ligand to titrate while binding and unbinding,^{5,6} or allowing His residues to fluctuate among the doubly protonated and two singly protonated tautomer states, which has been shown to affect the ligand binding mechanism and kinetics.^{88,89}

Supporting Information Available

Supporting Information contains additional tables and figures.

Acknowledgements

The authors acknowledge National Institutes of Health (R01GM098818) for funding.

References

- (1) Schuldiner, S. Competition as a Way of Life for H⁺-Coupled Antiporters. *J. Mol. Biol.* **2014**, 426.
- (2) Tan, J.; Verschueren, K. H.; Anand, K.; Shen, J.; Yang, M.; Xu, Y.; Rao, Z.; Bigalke, J.; Heisen, B.; Mesters, J. R.; Chen, K.; Shen, X.; Jiang, H.; Hilgenfeld, R. pH-dependent Conformational Flexibility of the SARS-CoV Main Proteinase (Mpro) Dimer: Molecular Dynamics Simulations and Multiple X-ray Structure Analyses. *J. Mol. Biol.* **2005**, 354, 25–40.
- (3) Verma, N.; Henderson, J. A.; Shen, J. Proton-Coupled Conformational Activation of SARS Coronavirus Main Proteases and Opportunity for Designing Small-Molecule Broad-Spectrum Targeted Covalent Inhibitors. *J. Am. Chem. Soc.* **2020**, jacs.0c10770.
- (4) Morrow, B. H.; Payne, G. F.; Shen, J. pH-Responsive Self-Assembly of Polysaccharide through a Rugged Energy Landscape. *J. Am. Chem. Soc.* **2015**, 137, 13024–13030.
- (5) Harris, R. C.; Tsai, C.-C.; Ellis, C. R.; Shen, J. Proton-Coupled Conformational Allostery Modulates the Inhibitor Selectivity for β -Secretase. *J. Phys. Chem. Lett.* **2017**, 8, 4832–4837.
- (6) Henderson, J. A.; Harris, R. C.; Tsai, C.-C.; Shen, J. How Ligand Protonation State Controls Water in Protein–Ligand Binding. *J. Phys. Chem. Lett.* **2018**, 9, 5440–5444.
- (7) Baptista, A. M.; Teixeira, V. H.; Soares, C. M. Constant-*p* H Molecular Dynamics Using Stochastic Titration. *J. Chem. Phys.* **2002**, 117, 4184–4200.
- (8) Lee, M. S.; Salsbury, F. R.; Brooks III, C. L. Constant-pH Molecular Dynamics Using Continuous Titration Coordinates. *Proteins* **2004**, 56, 738–752.

- (9) Khandogin, J.; Brooks III, C. L. Constant pH Molecular Dynamics with Proton Tautomerism. *Biophys. J.* **2005**, *89*, 141–157.
- (10) Goh, G. B.; Hulbert, B. S.; Zhou, H.; Brooks III, C. L. Constant pH Molecular Dynamics of Proteins in Explicit Solvent with Proton Tautomerism: Explicit Solvent CPHMD of Proteins. *Proteins* **2014**, *82*, 1319–1331.
- (11) Huang, Y.; Chen, W.; Wallace, J. A.; Shen, J. All-Atom Continuous Constant pH Molecular Dynamics With Particle Mesh Ewald and Titratable Water. *J. Chem. Theory Comput.* **2016**, *12*, 5411–5421.
- (12) Donnini, S.; Tegeler, F.; Groenhof, G.; Grubmüller, H. Constant pH Molecular Dynamics in Explicit Solvent with λ -Dynamics. *J. Chem. Theory Comput.* **2011**, *7*, 1962–1978.
- (13) Kong, X.; Brooks III, C. L. Lambda-Dynamics: A New Approach to Free Energy Calculations. *J. Chem. Phys.* **1996**, *105*, 10.
- (14) Mongan, J.; Case, D. A.; McCammon, J. A. Constant pH Molecular Dynamics in Generalized Born Implicit Solvent. *J. Comput. Chem.* **2004**, *25*, 2038–2048.
- (15) Swails, J. M.; York, D. M.; Roitberg, A. E. Constant pH Replica Exchange Molecular Dynamics in Explicit Solvent Using Discrete Protonation States: Implementation, Testing, and Validation. *J. Chem. Theory Comput.* **2014**, *10*, 1341–1352.
- (16) Chen, Y.; Roux, B. Constant-pH Hybrid Nonequilibrium Molecular Dynamics–Monte Carlo Simulation Method. *J. Chem. Theory Comput.* **2015**, *11*, 3919–3931.
- (17) Chen, J.; Brooks III, C. L.; Khandogin, J. Recent Advances in Implicit Solvent-Based Methods for Biomolecular Simulations. *Curr. Opin. Struct. Biol.* **2008**, *18*, 140–148.
- (18) Wallace, J. A.; Shen, J. K. *Methods Enzymol.*; Elsevier, 2009; Vol. 466; pp 455–475.

- (19) Chen, W.; Morrow, B. H.; Shi, C.; Shen, J. K. Recent Development and Application of Constant pH Molecular Dynamics. *Mol. Simulat.* **2014**, *40*, 830–838.
- (20) van Gunsteren, W. F.; Billeter, S. R.; A. A. Eising, P. H. H.; Krüger, P.; Mark, A. E.; Scott, W. R. P.; Tironi, I. G. Biomolecular Simulation: The GROMOS96 Manual and User Guide. 1996.
- (21) Machuqueiro, M.; Baptista, A. M. Constant-pH Molecular Dynamics with Ionic Strength Effects: Protonation-Conformation Coupling in Decalysine. *J. Phys. Chem. B* **2006**, *110*, 2927–2933.
- (22) Abraham, M. J.; Murtola, T.; Schulz, R.; Páll, S.; Smith, J. C.; Hess, B.; Lindahl, E. GROMACS: High Performance Molecular Simulations through Multi-Level Parallelism from Laptops to Supercomputers. *SoftwareX* **2015**, *1–2*, 19–25.
- (23) Wallace, J. A.; Shen, J. K. Continuous Constant pH Molecular Dynamics in Explicit Solvent with pH-Based Replica Exchange. *J. Chem. Theory Comput.* **2011**, *7*, 2617–2629.
- (24) Brooks, B.; Brooks III, C.; MacKerell, A.; Nilsson, L.; Petrella, R.; Roux, B.; Won, Y.; Archontis, G.; Bartels, C.; Boresch, S.; Caflisch, A.; Caves, L.; Cui, Q.; Dinner, A.; Feig, M.; Fischer, S.; Gao, J.; Hodoscek, M.; Im, W.; Kuczera, K.; Lazaridis, T.; Ma, J.; Ovchinnikov, V.; Paci, E.; Pastor, R.; Post, C.; Pu, J.; Schaefer, M.; Tidor, B.; Venable, R. M.; Woodcock, H. L.; Wu, X.; Yang, W.; York, D.; Karplus, M. CHARMM: The Biomolecular Simulation Program. *J. Comput. Chem.* **2009**, *30*, 1545–1614.
- (25) Case, D.; Aktulga, H.; Belfon, K.; Ben-Shalom, I.; Berryman, J.; Brozell, S.; Cerutti, D.; Cheatham, T., III; Cisneros, G.; Cruzeiro, V.; Darden, T.; Duke, R.; Giambasu, G.; Gilson, M.; Gohlke, H.; Goetz, A.; Harris, R.; Izadi, S.; Izmailov, S.; Kasavajhala, K.; Kaymak, M.; King, E.; Kovalenko, A.; Kurtzman, T.; Lee, T.; LeGrand, S.; Li, P.; Lin, C.; Liu, J.; Luchko, T.; Luo, R.; Machado, M.; Man, V;

- Manathunga, M.; Merz, K.; Miao, Y.; Mikhailovskii, O.; Monard, G.; Nguyen, H.; O'Hearn, K.; Onufriev, A.; Pan, F.; Pantano, S.; Qi, R.; Rahnamoun, A.; Roe, D.; Roitberg, A.; Sagui, C.; Schott-Verdugo, S.; Shajan, A.; Shen, J.; Simmerling, C.; Skrynnikov, N.; Smith, J.; Swails, J.; Walker, R.; Wang, J.; Wang, J.; Wei, H.; Wolf, R.; Wu, X.; Xiong, Y.; Xue, Y.; York, D.; Zhao, S.; Kollman, P. AMBER 2022. 2022.
- (26) Shi, C.; Wallace, J.; Shen, J. Thermodynamic Coupling of Protonation and Conformational Equilibria in Proteins: Theory and Simulation. *Biophys. J.* **2012**, *102*, 1590–1597.
- (27) Hofer, F.; Kraml, J.; Kahler, U.; Kamenik, A. S.; Liedl, K. R. Catalytic Site pK_a Values of Aspartic, Cysteine, and Serine Proteases: Constant pH MD Simulations. *J. Chem. Inf. Model.* **2020**, *60*, 3030–3042.
- (28) Carvalheda, C. A.; Campos, S. R. R.; Machuqueiro, M.; Baptista, A. M. Structural Effects of pH and Deacylation on Surfactant Protein C in an Organic Solvent Mixture: A Constant-pH MD Study. *J. Chem. Inf. Model.* **2013**, *53*, 2979–2989.
- (29) Morrow, B. H.; Koenig, P. H.; Shen, J. K. Atomistic simulations of pH-dependent self-assembly of micelle and bilayer from fatty acids. *J. Chem. Phys.* **2012**, *137*, 194902.
- (30) Santos, H. A. F.; Vila-Viçosa, D.; Teixeira, V. H.; Baptista, A. M.; Machuqueiro, M. Constant-pH MD Simulations of DMPA/DMPC Lipid Bilayers. *J. Chem. Theory Comput.* **2015**, *11*, 5973–5979.
- (31) Teixeira, V. H.; Vila-Viçosa, D.; Reis, P. B. P. S.; Machuqueiro, M. pK_a Values of Titrable Amino Acids at the Water/Membrane Interface. *J. Chem. Theory Comput.* **2016**, *12*, 930–934.
- (32) Yue, Z.; Chen, W.; Zgurskaya, H. I.; Shen, J. Constant pH Molecular Dynamics Reveals How Proton Release Drives the Conformational Transition of a Transmembrane Efflux Pump. *J. Chem. Theory Comput.* **2017**, *13*, 6405–6414.

- (33) Vila-Viçosa, D.; Silva, T. F. D.; Slaybaugh, G.; Reshetnyak, Y. K.; Andreev, O. A.; Machuqueiro, M. The Membrane-Induced pKa Shifts in Wt-pHLIP and Its L16H Variant. *J. Chem. Theory Comput.* **2018**, *14*, 3289–3297.
- (34) Huang, Y.; Chen, W.; Dotson, D. L.; Beckstein, O.; Shen, J. Mechanism of pH-dependent activation of the sodium-proton antiporter NhaA. *Nat. Commun.* **2016**, *7*, 12940.
- (35) Wallace, J. A.; Shen, J. K. Charge-Leveling and Proper Treatment of Long-Range Electrostatics in All-Atom Molecular Dynamics at Constant pH. *J. Chem. Phys.* **2012**, *137*, 184105.
- (36) Chen, W.; Wallace, J. A.; Yue, Z.; Shen, J. K. Introducing Titratable Water to All-Atom Molecular Dynamics at Constant pH. *Biophys. J.* **2013**, *105*, L15–L17.
- (37) Knight, J. L.; Brooks III, C. L. Multisite λ Dynamics for Simulated Structure–Activity Relationship Studies. *J. Chem. Theory Comput.* **2011**, *7*, 2728–2739.
- (38) Goh, G. B.; Knight, J. L.; Brooks III, C. L. Constant pH Molecular Dynamics Simulations of Nucleic Acids in Explicit Solvent. *J. Chem. Theory Comput.* **2012**, *8*, 36–46.
- (39) Phillips, J. C.; Braun, R.; Wang, W.; Gumbart, J.; Tajkhorshid, E.; Villa, E.; Chipot, C.; Skeel, R. D.; Kalé, L.; Schulten, K. Scalable Molecular Dynamics with NAMD. *J. Comput. Chem.* **2005**, *26*, 1781–1802.
- (40) Hayes, R. L.; Buckner, J.; Brooks III, C. L. BLaDE: A Basic Lambda Dynamics Engine for GPU-Accelerated Molecular Dynamics Free Energy Calculations. *J. Chem. Theory Comput.* **2021**, *17*, 6799–6807.
- (41) Case, D. A.; Ben-Shalom, I. Y.; Brozell, S. R.; Cerutti, D. S.; Cheatham, T., III; Cruzeiro, V. W. D.; Darden, T. A.; Duke, R. E.; Ghoreishi, D.; Gilson, M. K.; Gohlke, H.; Goetz, A. W.; Greene, D.; Harris, R.; Homeyer, N.; Huang, Y.; Izadi, S.;

- Kovalenko, A.; Kurtzman, T.; Lee, T. S.; LeGrand, S.; Li, P.; Lin, C.; Liu, J.; Luchko, T.; Luo, R.; Mermelstein, D. J.; Merz, K. M.; Miao, Y.; Monard, G.; Nguyen, C.; Nguyen, H.; Omelyan, I.; Onufriev, A.; Pan, F.; Qi, R.; Roe, D. R.; Roitberg, A.; Sagui, C.; Schott-Verdugo, S.; Shen, J.; Simmerling, C. L.; Smith, J.; Salomon-Ferrer, R.; Swails, J.; Walker, R. C.; Wang, J.; Wei, H.; Wolf, R. M.; Wu, X.; Xiao, L.; York, D. M.; Kollman, P. A. AMBER 2020. 2018.
- (42) Khandogin, J.; Brooks III, C. L. Constant pH Molecular Dynamics with Proton Tautomerism. *Biophys. J.* **2005**, *89*, 141–157.
- (43) Harris, R. C.; Shen, J. GPU-Accelerated Implementation of Continuous Constant pH Molecular Dynamics in Amber: pKa Predictions with Single-pH Simulations. *J. Chem. Inf. Model.* **2019**, *59*, 4821–4832.
- (44) MacKerell, A. D. J.; Bashford, D.; Bellott, M.; Dunbrack, R. L.; Evanseck, J. D.; Field, M. J.; Fischer, S.; Gao, J.; Guo, H.; Ha, S.; Joseph-McCarthy, D.; Kuchnir, L.; Kuczera, K.; Lau, F. T. K.; Mattos, C.; Michnick, S.; Ngo, T.; Nguyen, D. T.; Prodhom, B.; Reiher, W. E.; Roux, B.; Schlenkrich, M.; Smith, J. C.; Stote, R.; Straub, J.; Watanabe, M.; Wiórkiewicz-Kuczera, J.; Yin, D.; Karplus, M. All-Atom Empirical Potential for Molecular Modeling and Dynamics Studies of Proteins. *J. Phys. Chem. B* **1998**, *102*, 3586–3616.
- (45) Cornell, W. D.; Cieplak, P.; Bayly, C. I.; Gould, I. R.; Merz, K. M.; Ferguson, D. M.; Spellmeyer, D. C.; Fox, T.; Caldwell, J. W.; Kollman, P. A. A Second Generation Force Field for the Simulation of Proteins, Nucleic Acids, and Organic Molecules. *J. Am. Chem. Soc.* **1995**, *117*, 5179–5197.
- (46) Maier, J. A.; Martinez, C.; Kasavajhala, K.; Wickstrom, L.; Hauser, K. E.; Simmerling, C. ff14SB: Improving the Accuracy of Protein Side Chain and Backbone Parameters from ff99SB. *J. Chem. Theory Comput.* **2015**, *11*, 3696–3713.

- (47) Gu, Z.; Ridenour, C. F.; Bronnimann, C. E.; Iwashita, T.; McDermott, A. Hydrogen Bonding and Distance Studies of Amino Acids and Peptides Using Solid State 2D ^1H - ^{13}C Heteronuclear Correlation Spectra. *J. Am. Chem. Soc.* **1996**, *118*, 822–829.
- (48) Levy, R. M.; Gallicchio, E. Computer Simulations with Explicit Solvent: Recent Progress in the Thermodynamic Decomposition of Free Energies and in Modeling Electrostatic Effects. *Annu. Rev. Phys. Chem.* **1998**, *49*, 531–567.
- (49) Henderson, J. A.; Liu, R.; Harris, J. A.; Huang, Y.; de Oliveira, V. M.; Shen, J. A Guide to the Continuous Constant pH Molecular Dynamics Methods in Amber and CHARMM v1.0. *Liv. J. Comput. Mol. Sci.* **2022**, *4*, 1563.
- (50) Rocklin, G. J.; Mobley, D. L.; Dill, K. A.; Hünenberger, P. H. Calculating the Binding Free Energies of Charged Species Based on Explicit-Solvent Simulations Employing Lattice-Sum Methods: An Accurate Correction Scheme for Electrostatic Finite-Size Effects. *J. Chem. Phys.* **2013**, *139*, 184103.
- (51) Lin, Y.-L.; Aleksandrov, A.; Simonson, T.; Roux, B. An Overview of Electrostatic Free Energy Computations for Solutions and Proteins. *J. Chem. Theory Comput.* **2014**, *10*, 2690–2709.
- (52) Tian, C.; Kasavajhala, K.; Belfon, K. A. A.; Raguette, L.; Huang, H.; Miguës, A. N.; Bickel, J.; Wang, Y.; Pincay, J.; Wu, Q.; Simmerling, C. ff19SB: Amino-Acid-Specific Protein Backbone Parameters Trained against Quantum Mechanics Energy Surfaces in Solution. *J. Chem. Theory Comput.* **2020**, *16*, 528–552.
- (53) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (54) Feenstra, K. A.; Hess, B.; Berendsen, H. J. C. Improving Efficiency of Large Time-

- Scale Molecular Dynamics Simulations of Hydrogen-Rich Systems. *J. Comput. Chem.* **1999**, *20*, 786–798.
- (55) Hopkins, C. W.; Le Grand, S.; Walker, R. C.; Roitberg, A. E. Long-Time-Step Molecular Dynamics through Hydrogen Mass Repartitioning. *J. Chem. Theory Comput.* **2015**, *11*, 1864–1874.
- (56) Ferguson, N.; Sharpe, T. D.; Schartau, P. J.; Sato, S.; Allen, M. D.; Johnson, C. M.; Rutherford, T. J.; Fersht, A. R. Ultra-Fast Barrier-limited Folding in the Peripheral Subunit-binding Domain Family. *J. Mol. Biol.* **2005**, *353*, 427–446.
- (57) Ramanadham, M.; Sieker, L. C.; Jensen, L. H. Refinement of Triclinic Lysozyme: II. The Method of Stereochemically Restrained Least Squares. *Acta Cryst. B* **1990**, *46*, 63–69.
- (58) Castañeda, C. A.; Fitch, C. A.; Majumdar, A.; Khangulov, V.; Schlessman, J. L.; García-Moreno, B. E. Molecular Determinants of the pK_a Values of Asp and Glu Residues in Staphylococcal Nuclease: pK_a Values of Asp and Glu in SNase. *Proteins* **2009**, *77*, 570–588.
- (59) Wlodawer, A.; Svensson, L. A.; Sjoelin, L.; Gilliland, G. L. Structure of Phosphate-Free Ribonuclease A Refined at 1.26 Å. *Biochemistry* **1988**, *27*, 2705–2717.
- (60) Weichsel, A.; Gasdaska, J. R.; Powis, G.; Montfort, W. R. Crystal Structures of Reduced, Oxidized, and Mutated Human Thioredoxins: Evidence for a Regulatory Homodimer. *Structure* **1996**, *4*, 735–751.
- (61) Shen, Y.-q.; Tang, L.; Zhou, H.-m.; Lin, Z.-j. Structure of Human Muscle Creatine Kinase. *Acta Cryst. D* **2001**, *57*, 1196–1200.
- (62) Feig, M.; Karanicolas, J.; Brooks III, C. L. MMTSB Tool Set: Enhanced Sampling and

- Multiscale Modeling Methods for Applications in Structural Biology. *J. Mol. Graph. Model.* **2004**, *22*, 377—395.
- (63) Shirts, M. R.; Klein, C.; Swails, J. M.; Yin, J.; Gilson, M. K.; Mobley, D. L.; Case, D. A.; Zhong, E. D. Lessons learned from comparing molecular dynamics engines on the SAMPL5 dataset. *J. Comput. Aided Mol. Des.* **2017**, *31*, 147–161.
- (64) Henderson, J. A.; Verma, N.; Harris, R. C.; Liu, R.; Shen, J. Assessment of Proton-Coupled Conformational Dynamics of SARS and MERS Coronavirus Papain-like Proteases: Implication for Designing Broad-Spectrum Antiviral Inhibitors. *J. Chem. Phys.* **2020**, *153*, 115101.
- (65) Ullmann, G. M. Relations between Protonation Constants and Titration Curves in Polyprotic Acids: A Critical View. *J. Phys. Chem. B* **2003**, *107*, 1263–1271.
- (66) MacKerell, A. D.; Wiorkiewicz-Kuczera, J.; Karplus, M. An All-Atom Empirical Energy Function for the Simulation of Nucleic Acids. *J. Am. Chem. Soc.* **1995**, *117*, 11946–11975.
- (67) Tanokura, M. ¹H-NMR study on the tautomerism of the imidazole ring of histidine residues. I. Microscopic pK values and molar ratios of tautomers in histidine-containing peptides. *Biochim. Biophys. Acta* **1983**, *742*, 576–585.
- (68) Thurlkill, R. L.; Grimsley, G. R.; Scholtz, J. M.; Pace, C. N. pK Values of the Ionizable Groups of Proteins. *Protein Sci.* **2006**, *15*, 1214–1218.
- (69) Balusek, C.; Hwang, H.; Lau, C. H.; Lundquist, K.; Hazel, A.; Pavlova, A.; Lynch, D. L.; Reggio, P. H.; Wang, Y.; Gumbart, J. C. Accelerating Membrane Simulations with Hydrogen Mass Repartitioning. *J. Chem. Theory Comput.* **2019**, *15*, 4673–4686.

- (70) Ellis, C. R.; Shen, J. pH-Dependent Population Shift Regulates BACE1 Activity and Inhibition. *J. Am. Chem. Soc.* **2015**, *137*, 9543–9546.
- (71) Liu, R.; Yue, Z.; Tsai, C.-C.; Shen, J. Assessing Lysine and Cysteine Reactivities for Designing Targeted Covalent Kinase Inhibitors. *J. Am. Chem. Soc.* **2019**, *141*, 6553–6560.
- (72) Arbely, E.; Rutherford, T. J.; Sharpe, T. D.; Ferguson, N.; Fersht, A. R. Downhill versus Barrier-Limited Folding of BBL 1: Energetic and Structural Perturbation Effects upon Protonation of a Histidine of Unusually Low pKa. *J. Mol. Biol.* **2009**, *387*, 986–992.
- (73) Arbely, E.; Rutherford, T. J.; Neuweiler, H.; Sharpe, T. D.; Ferguson, N.; Fersht, A. R. Carboxyl pKa Values and Acid Denaturation of BBL. *J. Mol. Biol.* **2010**, *403*, 313–327.
- (74) Webb, H.; Tynan-Connolly, B. M.; Lee, G. M.; Farrell, D.; O'Meara, F.; Søndergaard, C. R.; Teilum, K.; Hewage, C.; McIntosh, L. P.; Nielsen, J. E. Re-measuring HEWL pKa Values by NMR Spectroscopy: Methods, Analysis, Accuracy, and Implications for Theoretical pKa Calculations. *Proteins* **2011**, *79*, 685–702.
- (75) Qin, J.; Clore, G. M.; Gronenborn, A. M. Ionization Equilibria for Side-Chain Carboxyl Groups in Oxidized and Reduced Human Thioredoxin and in the Complex with Its Target Peptide from the Transcription Factor NF κ B. *Biochemistry* **1996**, *35*, 7–13.
- (76) Baker, W. R.; Kintanar, A. Characterization of the pH Titration Shifts of Ribonuclease A by One- and Two-Dimensional Nuclear Magnetic Resonance Spectroscopy. *Arch. Biochem. Biophys.* **1996**, *327*, 189–199.
- (77) Wang, P.-F.; McLeish, M. J.; Kneen, M. M.; Lee, G.; Kenyon, G. L. An Unusually Low p K_a for Cys282 in the Active Site of Human Muscle Creatine Kinase[†]. *Biochemistry* **2001**, *40*, 11698–11705.

- (78) Huang, Y.; Yue, Z.; Tsai, C.-C.; Henderson, J. A.; Shen, J. Predicting Catalytic Proton Donors and Nucleophiles in Enzymes: How Adding Dynamics Helps Elucidate the Structure-Function Relationships. *J. Phys. Chem. Lett.* **2018**, 26.
- (79) Alexov, E.; Mehler, E. L.; Baker, N.; M. Baptista, A.; Huang, Y.; Milletti, F.; Erik Nielsen, J.; Farrell, D.; Carstensen, T.; Olsson, M. H. M.; Shen, J. K.; Warwicker, J.; Williams, S.; Word, J. M. Progress in the Prediction of pK_a Values in Proteins: Prediction of pK_a Values in Proteins. *Proteins* **2011**, 79, 3260–3275.
- (80) Harris, R. C.; Liu, R.; Shen, J. Predicting Reactive Cysteines with Implicit-Solvent-Based Continuous Constant pH Molecular Dynamics in Amber. *J. Chem. Theory Comput.* **2020**, 16, 3689–3698.
- (81) Roos, G.; Foloppe, N.; Messens, J. Understanding the pK_a of Redox Cysteines: The Key Role of Hydrogen Bonding. *Antioxid. Redox Signal.* **2013**, 18, 94–127.
- (82) Liu, R.; Zhan, S.; Che, Y.; Shen, J. Reactivities of the Front Pocket N-Terminal Cap Cysteines in Human Kinases. *J. Med. Chem.* **2022**, 65, 1525–1535.
- (83) Liu, R.; Verma, N.; Henderson, J. A.; Zhan, S.; Shen, J. Profiling MAP kinase cysteines for targeted covalent inhibitor design. *RSC Med. Chem.* **2022**, 13, 54–63.
- (84) Bignucolo, O.; Chipot, C.; Kellenberger, S.; Roux, B. Galvani Offset Potential and Constant-pH Simulations of Membrane Proteins. *J. Phys. Chem. B* **2022**, acs.jpcc.2c04593.
- (85) Best, R. B.; Zhu, X.; Shim, J.; Lopes, P. E. M.; Mittal, J.; Feig, M.; MacKerell, A. D. Optimization of the Additive CHARMM All-Atom Protein Force Field Targeting Improved Sampling of the Backbone ϕ , ψ and Side-Chain χ_1 and χ_2 Dihedral Angles. *J. Chem. Theory Comput.* **2012**, 8, 3257–3273.

- (86) Huang, J.; Lopes, P. E. M.; Roux, B.; MacKerell, A. D. Recent Advances in Polarizable Force Fields for Macromolecules: Microsecond Simulations of Proteins Using the Classical Drude Oscillator Model. *J. Phys. Chem. Lett.* **2014**, *5*, 3144–3150.
- (87) andSaeed Izadi, Y. X.; Onufriev, A. A Fast Polarizable Water model for Atomistic Simulations. *Chem RXiv* **2022**,
- (88) Vo, Q. N.; Mahinthichaichan, P.; Shen, J.; Ellis, C. R. How μ -Opioid Receptor Recognizes Fentanyl. *Nat. Commun.* **2021**, *12*, 984.
- (89) Mahinthichaichan, P.; Vo, Q. N.; Ellis, C. R.; Shen, J. Kinetics and Mechanism of Fentanyl Dissociation from the μ -Opioid Receptor. *JACS Au* **2021**, *1*, 2208–2215.