


 Cite this: *Phys. Chem. Chem. Phys.*, 2025, 27, 9794

# Determinants of hydrogen bond distances in proteins†

 Masaki Tsujimura, \*<sup>a</sup> Hiroshi Ishikita <sup>bc</sup> and Keisuke Saito \*<sup>bc</sup>

Hydrogen bonds (H-bonds) between oxygen atoms, with the O–H bond donated to the acceptor O atom ( $O_{\text{donor}}-\text{H}\cdots O_{\text{acceptor}}$ ), are essential for stabilizing protein structures and facilitating enzymatic reactions. The dielectric and electrostatic environment of proteins, as well as structural constraints imposed by protein folding, influence the nature of H-bonds. In this study, we investigated how these factors affect H-bond distances in proteins. Analysis of 906 high-resolution protein structures ( $\leq 1.2$  Å) from the Protein Data Bank revealed that H-bond distances for H-bonds with the same donor and acceptor groups are distributed around a value primarily determined by the  $pK_a$  difference between these groups ( $\Delta pK_a$ ) in water, with lower  $\Delta pK_a$  values leading to shorter distances. This correlation arises from enhanced electron redistribution from the H-bond acceptor to the donor in lower  $\Delta pK_a$  H-bonds, which increases the covalent character of the H-bond and decreases the  $\text{H}\cdots O_{\text{acceptor}}$  distance. In contrast, H-bond distances are largely unaffected by whether the H-bond is buried in the protein interior or exposed to bulk water, as the strength of the electrostatic interaction between the donor and acceptor groups plays a minor role in determining distances. Furthermore, analysis of H-bonds in microbial rhodopsins using a quantum mechanical/molecular mechanical approach demonstrates that the protein environment primarily influences H-bond distances electrostatically by altering the  $\Delta pK_a$  of the H-bond, while structural constraints impose a secondary influence by altering  $O_{\text{donor}}-\text{H}\cdots O_{\text{acceptor}}$  angles or  $\text{H}\cdots O_{\text{acceptor}}$  distances without changing  $\Delta pK_a$ .

 Received 8th February 2025,  
 Accepted 8th April 2025

DOI: 10.1039/d5cp00511f

rsc.li/pccp

## Introduction

Hydrogen bonds (H-bonds) are abundant in proteins and play a crucial role in stabilizing protein structures and facilitating enzymatic reactions.<sup>1</sup> The H-bond donor and acceptor can be defined as the groups donating and accepting the O–H bond (or the N–H bond), respectively (*e.g.*,  $O_{\text{donor}}-\text{H}\cdots O_{\text{acceptor}}$ ). H-bonds are characterized by properties such as distances, stretching vibrational frequencies, and <sup>1</sup>H NMR chemical shifts. Over 200 000 structures in the Protein Data Bank (PDB)<sup>2</sup> provide extensive information about H-bond distances.<sup>3–9</sup> These distances reflect the chemical properties of the donor and acceptor groups as well as the influence of the protein environment on the nature of the H-bond.

An H-bond can be considered an interaction between a Brønsted acid (H-bond donor) and a Brønsted base (H-bond acceptor).<sup>10</sup>  $pK_a$  is an indicator of the Brønsted acidity. Therefore, the  $pK_a$  difference between the donor and acceptor groups ( $\Delta pK_a$ )<sup>10</sup> reflects the characteristics of the H-bond. A lower  $\Delta pK_a$  leads to a higher binding enthalpy,<sup>11</sup> a lower  $O_{\text{donor}}-\text{H}$  stretching vibrational frequency,<sup>12,13</sup> and a higher <sup>1</sup>H NMR chemical shift.<sup>12–16</sup>

A positive correlation between  $O\cdots O$  distances and  $\Delta pK_a$  values was reported for 68 H-bonds of small compounds in neutron diffraction structures, with a coefficient of determination  $R^2 = 0.86$ .<sup>16,17</sup> Similarly, positive correlations between  $\text{N}\cdots\text{O}$  or  $\text{N}\cdots\text{N}$  distances and  $\Delta pK_a$  values were reported ( $R^2 = 0.74$  and  $0.69$  for 86 and 29 H-bonds, respectively).<sup>17</sup> Deviations from these correlations mainly arise from additional H-bonds of the donor and acceptor groups,<sup>17,18</sup> as well as geometric constraints in crystals.<sup>17</sup>  $O\cdots O$  distances are almost independent of the solvent dielectric environment, as the <sup>1</sup>H NMR chemical shift for the same H-bond remains largely unchanged across different solvents (chloroform, acetone, or water), and the calculated equilibrium  $O\cdots O$  distance of the  $[\text{HOH}\cdots\text{OH}]$  H-bond remains nearly constant among solvents with dielectric constants ranging from 5 to 78.<sup>16,17</sup>

The present study aims to elucidate the determinants of H-bond distances in proteins, focusing primarily on  $O\cdots O$

<sup>a</sup> Department of Advanced Interdisciplinary Studies, The University of Tokyo, 4-6-1 Komaba, Meguro-ku, Tokyo 153-8904, Japan.

E-mail: masaki.tsujimura@riken.jp; Fax: +81-3-5452-5083; Tel: +81-3-5452-5056

<sup>b</sup> Department of Applied Chemistry, The University of Tokyo, 7-3-1 Hongo,

Bunkyo-ku, Tokyo 113-8654, Japan. E-mail: ksaito@appchem.t.u-tokyo.ac.jp

<sup>c</sup> Research Center for Advanced Science and Technology, The University of Tokyo, 4-6-1 Komaba, Meguro-ku, Tokyo 153-8904, Japan

 † Electronic supplementary information (ESI) available. See DOI: <https://doi.org/10.1039/d5cp00511f>


H-bonds, but also examining N···O and N···N H-bonds. We begin by investigating small-compound H-bonds relevant to those in proteins to clarify the basis for the correlation between H-bond distances and  $\Delta pK_a$  in the absence of the protein environment.

However, in contrast to H-bonds in solution under comparable conditions, systematic analysis of H-bond distances in protein environments remains limited. A major challenge lies in the heterogeneity of the PDB, which contains structures determined under inconsistent conditions and resolutions. For instance, many structures resolved at  $>2.5$  Å lack sufficient clarity to assign side-chain orientations or identify interacting water molecules. This heterogeneity has hindered meaningful comparisons of H-bond properties among proteins.

To overcome these limitations, we next focus exclusively on 906 high-resolution protein crystal structures ( $\leq 1.2$  Å), in which both side chains and interacting water molecules can be reliably modeled based on electron density. By limiting our dataset to these high-resolution structures, we are able to clarify which factors in the protein environment determine H-bond distances. This classification further enables us to compare buried and solvent-exposed H-bonds and to elucidate how the dielectric properties of the protein environment influence H-bond distances. Finally, based on the insights obtained, we analyze H-bonds in microbial rhodopsins using a quantum mechanical/molecular mechanical (QM/MM) approach.

## Methods

### Quantum chemical calculations of small-compound H-bonds

H-bonds involving water ( $pK_a = 15.74$ ), methanol ( $15.5^{19}$ ), phenol ( $9.99^{19}$ ), and protonated acetic acid ( $4.76^{19}$ ) as donors, and water ( $-1.74$ ), methanol ( $-2^{20}$ ), phenol ( $-6^{20}$ ), deprotonated acetic acid ( $4.76^{19}$ ), protonated acetic acid ( $-6^{20}$ ), and *N*-methylacetamide ( $-1^{20}$ ) as acceptors, were analyzed. H-bond donors and acceptors were solvated with explicit water molecules, and additional solvent effects were modeled using the polarizable continuum model (PCM) (Fig. S1, ESI†). This approach, which solvates H-bonded compounds using several explicit water molecules and an implicit solvent model, has been shown to reproduce the experimentally observed vibrational<sup>21</sup> and excitation<sup>22</sup> energies of H-bonded compounds in aqueous solutions. Geometries were optimized using the second-order Møller-Plesset perturbation theory (MP2) method. H-bond interaction energies were calculated using the two-body fragment molecular orbital (FMO) method<sup>23</sup> combined with MP2,<sup>24</sup> and decomposed into electrostatic, exchange, charge-transfer, and dispersion terms using pair interaction energy decomposition analysis (PIEDA).<sup>25</sup> Although the FMO method is commonly applied to large molecular systems, it was specifically used here due to its capability to calculate and decompose H-bond interaction energies in the presence of explicit solvent water molecules. To ensure meaningful fragmentation, each molecule in the system was treated as a separate fragment. Dimer corrections were included in the

electrostatic potential calculations in PCM (*i.e.*, PCM[2]).<sup>26</sup> The electron redistribution amount was obtained from the Mulliken charge distribution in the dimer. The 6-31G\* basis set was used. All calculations were performed using the GAMESS program.<sup>27</sup>

Natural bond orbital (NBO)<sup>28,29</sup> energies were calculated for (i) water, methanol, phenol, and protonated acetic acid molecules donating an H-bond to a water molecule, and (ii) water, methanol, phenol, deprotonated acetic acid, protonated acetic acid, and *N*-methylacetamide molecules accepting an H-bond from a water molecule (Fig. S2, ESI†). NBO energies were obtained following geometry optimization using the density functional theory (DFT) method. To include long-range corrections in the DFT functional,<sup>30,31</sup> the CAM-B3LYP functional<sup>32</sup> was employed with the 6-31G\*\*+ basis set. The CAM-B3LYP-related parameters  $\alpha$ ,  $\beta$ , and  $\mu$  were set to the standard values of 0.19, 0.46, and 0.33, respectively.<sup>32</sup> All calculations were performed using the NBO 5.0 program<sup>33</sup> implemented in Jaguar.<sup>34</sup>

### Distributions of H-bond distances in high-resolution protein structures

A dataset of 906 protein structures was obtained from ref. 9. This dataset comprises X-ray crystal structures of proteins with resolutions of  $\leq 1.2$  Å, sequence identities of  $\leq 90\%$ , and *R*-factors of  $\leq 0.20$ .<sup>9</sup> PDB IDs of the structures included in the dataset are listed in Supporting Data 1 (ESI†). For structures with several models, the first model was analyzed. Distances for all O···O, N···O, and N···N atom pairs satisfying the following criteria were obtained, excluding pairs within the same residue: (i) O···O distance  $< 3.0$  Å, N···O or N···N distance  $< 3.2$  Å; (ii) *B*-factor values for both atoms  $< 40$  Å<sup>2</sup>; and (iii) occupancies for both atoms equal to 1.00. For  $sp^2$ -hybridized N atoms in the backbone and in the side-chains of Arg, His, and Trp, a dihedral angle criterion was applied to exclude N···O and N···N atom pairs that are closer than 3.2 Å but do not form H-bonds (Fig. S3, ESI†). H-bonds involving Asn and Gln side-chains were excluded due to difficulties in unambiguously assigning the O and N atoms in these side-chains, even in high-resolution structures. The analysis was conducted using the Biopython package<sup>35,36</sup> in Python.

H-bonds with relative solvent accessibilities of  $< 16\%$  and  $\geq 16\%$  were classified as buried and exposed, respectively.<sup>37</sup> Solvent accessibilities were calculated in the absence of crystal water molecules, using the DSSP program.<sup>38,39</sup> Therefore, the absence of water molecules that are difficult to capture by crystallography does not affect this classification. Asymmetric units, the smallest portions of crystal structures, are deposited in the PDB. The entire crystal can be reconstructed by applying symmetry operations to the asymmetric unit. In the absence of neighboring asymmetric units,  $\sim 10\%$  of H-bonds that are buried in the entire crystal are calculated as exposed (Fig. S4 and Table S1, ESI†). Therefore, solvent accessibilities were calculated in the presence of neighboring asymmetric units. Residues of neighboring asymmetric units within 7 Å of the focusing asymmetric unit were included in the calculation of solvent accessibility. The 7 Å threshold was set to be longer



than the sum of the longest atomic diameter (3.74 Å for the C<sub>α</sub> atom) and the water probe diameter (2.80 Å) used in the DSSP program.<sup>38</sup>

### QM/MM calculations of microbial rhodopsins

Atomic coordinates were obtained from the X-ray crystal structures of (i) ground-state bacteriorhodopsin from *Halobacterium salinarum* (BR) (PDB ID: 5ZIM<sup>40</sup>), (ii) N'-state V49A mutant BR (1P8U<sup>41</sup>), (iii) halorhodopsin from *Natromonas pharaonis* (pHR) (3A7K<sup>42</sup>), (iv) sodium-pumping rhodopsin KR2 from *Krokinobacter eikastus* (6YC3<sup>43</sup>), and (v) sodium-pumping rhodopsin ErNaR from *Erythrobacter* sp. HL-111 (8QLF<sup>44</sup>). Hydrogen atom positions were optimized with the heavy atom positions fixed, using the CHARMM program.<sup>45</sup> Atomic charges and force field parameters were obtained from the CHARMM22 parameter set.<sup>46</sup>

The protonation pattern was determined using the electrostatic continuum model by solving the linear Poisson–Boltzmann equation with the MEAD program.<sup>47</sup> The experimentally measured pK<sub>a</sub> values employed as references were 12.0 for Arg, 4.0 for Asp, 9.5 for Cys, 4.4 for Glu, 10.4 for Lys, 9.6 for Tyr,<sup>48</sup> and 7.0 and 6.6 for the N<sub>ε</sub> and N<sub>δ</sub> atoms of His, respectively.<sup>49–51</sup> Dielectric constants were set to 4 for the protein interior and 80 for water. All calculations were performed at 300 K, pH 7.0, and with an ionic strength of 100 mM. The linear Poisson–Boltzmann equation was solved using a three-step grid-focusing procedure at resolutions of 2.5, 1.0, and 0.3 Å. Protonation patterns were sampled using the Monte Carlo method with the Karlsberg program.<sup>52</sup>

Geometries were optimized using a QM/MM approach. The restricted DFT method was employed with the B3LYP functional and the LACVP\*\* basis set, using the QSite program.<sup>53,54</sup> The QM region was defined as follows: (i) ground-state BR: retinal, side-chains of Lys216, Tyr57, Arg82, Asp85, Trp86, Thr89, Tyr185, and Asp212, and H<sub>2</sub>O-402, 401, and 406. (ii) N'-state BR: side-chains of Tyr57, Asp85, and Asp212, and H<sub>2</sub>O-401, 406, and 407. (iii) pHR: retinal, side-chains of Lys256, Ser78, Ser81, Tyr82, Arg123, Thr126, Trp127, Ser130, Tyr225, Asp252, and Tyr257, H<sub>2</sub>O-502, 503, and 504, and Cl<sup>-</sup>-401. (iv) KR2: retinal, side-chains of Lys255, Ser70, Arg109, Asn112, Trp113, Asp116, Tyr218, Asp251, and Ser254, and H<sub>2</sub>O-434, 437, 501, and 512. (v) ErNaR: retinal, side-chains of Lys246, Ser25, Ser60, Glu64, Arg98, Trp102, Asp105, Tyr215, Thr239, and Asp242, and H<sub>2</sub>O-503, 504, 532, and 542. All atomic coordinates were relaxed in the QM region. In the MM region, hydrogen atom positions were optimized using the

OPLS2005 force field,<sup>55</sup> while heavy atom positions were fixed. The protonation pattern of titratable residues in the MM region was implemented in the atomic partial charges. Vibrational frequencies were calculated at the same level of theory as the geometry optimizations. The calculated frequencies were scaled using a standard factor of 0.9614 for the B3LYP functional.<sup>56</sup>

## Results and discussion

In this study, we classified the donors and acceptors of O<sub>donor</sub>–H···O<sub>acceptor</sub> H-bonds in proteins (excluding ligand molecules) into the following five groups (Table 1): (i) water molecules; (ii) Ser and Thr side-chains, which have hydroxyl OH groups (denoted as C–OH in this study); (iii) Tyr side-chains (denoted as PhOH); (iv) Asp and Glu side-chains, which have carboxyl COOH groups (denoted as COOH); and (v) backbone and Asn/Gln side-chains, which have amide C=O groups (denoted as C=O). Relevant H-bonds involving (i) water molecules, (ii) alcohols (denoted as C–OH), (iii) phenols (denoted as PhOH), (iv) carboxylic acids (denoted as COOH), and (v) amide C=O groups (denoted as C=O) are abundant in the crystal structures of small compounds (Table 1).<sup>6,57,58</sup>

Even when the groups forming an H-bond are the same, several H-bond types can exist depending on (i) which group serves as the donor and (ii) the protonation states of the donor and acceptor groups. For instance, the donor/acceptor of an H-bond between water and a carboxyl group can be COOH/H<sub>2</sub>O, H<sub>2</sub>O/COO<sup>-</sup>, or H<sub>2</sub>O/COOH. Different H-bond types for the same pair can be distinguished for crystal structures of small compounds obtained from the Cambridge Structural Database (CSD).<sup>59</sup> This is not the case for H-bonds in proteins, where the hydrogen atom positions are mostly unidentified. To clarify why H-bond distances are primarily determined by ΔpK<sub>a</sub>, we first investigate small-compound H-bonds relevant to those in proteins.

### Basis for the correlation between H-bond distances and ΔpK<sub>a</sub>

H-bonds with H<sub>2</sub>O, C–OH, PhOH, and COOH groups as donors are investigated. Similarly, H-bonds with H<sub>2</sub>O, C–OH, PhOH, COO<sup>-</sup>, COOH, and C=O groups as acceptors are investigated (Table 1). To confirm the correlation between O···O distances and ΔpK<sub>a</sub>, the average O···O distance was compared with ΔpK<sub>a</sub> for 23 out of the 24 possible H-bond types involving these donor and acceptor groups (Table 2 and Fig. 1). Average O···O distances were obtained from small-compound H-bonds in

Table 1 Approximate pK<sub>a</sub> values in water of five groups

	Small compounds	In proteins	pK <sub>a</sub> as a donor (acid/conjugated base)	pK <sub>a</sub> as an acceptor (base/conjugated acid)
(i)	Water	Water	16 (H <sub>2</sub> O/OH <sup>-</sup> )	-2 (H <sub>2</sub> O/H <sub>3</sub> O <sup>+</sup> )
(ii)	Alcohol	Ser, Thr	16 (C–OH/C–O <sup>-</sup> ) <sup>a</sup>	-2 (C–OH/C–OH <sub>2</sub> <sup>+</sup> ) <sup>c</sup>
(iii)	Phenol	Tyr	10 (PhOH/PhO <sup>-</sup> ) <sup>b</sup>	-6 (PhOH/PhOH <sub>2</sub> <sup>+</sup> ) <sup>c</sup>
(iv)	Carboxylic acid	Asp, Glu	4 (COOH/COO <sup>-</sup> ) <sup>b</sup>	4 (COO <sup>-</sup> /COOH) <sup>b</sup> , -6 (COOH/COOH <sub>2</sub> <sup>+</sup> ) <sup>c</sup>
(v)	Amide C=O	Backbone, Asn, Gln	—	-1 (C=O/C=OH <sup>+</sup> ) <sup>c</sup>

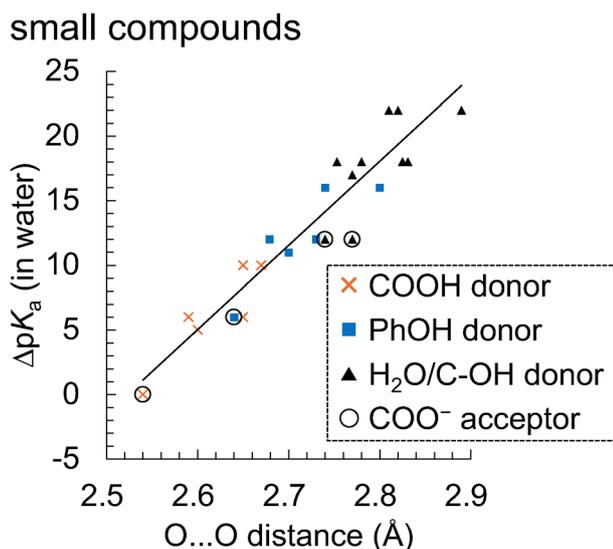
<sup>a</sup> pK<sub>a</sub> value of methanol.<sup>19</sup> <sup>b</sup> Ref. 48. <sup>c</sup> Ref. 20.



**Table 2** Average O...O distances and  $\Delta pK_a$  values for each H-bond type in crystal structures from the CSD

H-bond pair	Donor	Acceptor	$\Delta pK_a$ (in water)	$\overline{r_{O...O}}^a$ (Å)
Water...Water	H <sub>2</sub> O	H <sub>2</sub> O	18	2.83
Water...Alcohol	H <sub>2</sub> O	C-OH	18	2.83
	C-OH	H <sub>2</sub> O	18	2.75
Water...Phenol	PhOH	H <sub>2</sub> O	12	2.68
	H <sub>2</sub> O	PhOH	22	2.89
Water...Carboxylic acid	COOH	H <sub>2</sub> O	6	2.59
	H <sub>2</sub> O	COO <sup>-</sup>	12	2.77
	H <sub>2</sub> O	COOH	22	2.82
Alcohol...Alcohol	C-OH	C-OH	18	2.78
	C-OH	PhOH	22	2.82
Alcohol...Phenol	PhOH	C-OH	12	2.73
	C-OH	PhOH	22	2.82
Alcohol...Carboxylic acid	COOH	C-OH	6	2.65
	C-OH	COO <sup>-</sup>	12	2.74
	C-OH	COOH	22	2.81
	C-OH	C=O	17	2.77
Alcohol...Amide C=O	C-OH	C=O	17	2.77
Phenol...Phenol	PhOH	PhOH	16	2.80
	PhOH	COO <sup>-</sup>	6	2.64
Phenol...Carboxylic acid	COOH	PhOH	10	2.67
	PhOH	COOH	16	2.74
	PhOH	C=O	11	2.70
	COOH	COO <sup>-</sup>	0	2.54
Phenol...Amide C=O	PhOH	C=O	11	2.70
	COOH	COO <sup>-</sup>	0	2.54
Carboxylic acid...Carboxylic acid	COOH	COO <sup>-</sup>	0	2.54
	COOH	COOH	10	2.65
Carboxylic acid...Amide C=O	COOH	C=O	5	2.60

<sup>a</sup> Average O...O distances. Values for [water...water], [water...alcohol], and [water...phenol] pairs were taken from ref. 57. Values for the [water...carboxylic acid] pair were taken from ref. 58. Other values were taken from ref. 6.



**Fig. 1** Correlation between average O...O distances<sup>6,57,58</sup> and  $\Delta pK_a$  for each H-bond type in small-compound crystal structures from the CSD ( $R^2 = 0.89$ ). Orange crosses, blue squares, and black triangles indicate H-bonds involving COOH, PhOH, and H<sub>2</sub>O/C-OH groups as donors, respectively. H-bonds involving the COO<sup>-</sup> group as acceptors are surrounded by open circles.

the CSD, with each group containing from 18 to 4931 H-bonds.<sup>6,57,58</sup> The average O...O distance is highly correlated with  $\Delta pK_a$  for each H-bond type (Fig. 1,  $R^2 = 0.89$ ), confirming that O...O distances are primarily determined by  $\Delta pK_a$ .

To clarify the basis for the correlation between O...O distances and  $\Delta pK_a$ , we performed quantum chemical calculations of small-compound H-bonds. The O...O distance is correlated with  $\Delta pK_a$  for quantum-chemically optimized H-bonds in water solvent (Fig. S5a,  $R^2 = 0.73$ , ESI<sup>†</sup>), which aligns with the correlation between the average O...O distance from the CSD and  $\Delta pK_a$  (Fig. 1).

A lower  $\Delta pK_a$  leads to a shorter O...O distance because as  $\Delta pK_a$  decreases, (i) the O<sub>donor</sub>-H distance increases, and (ii) the H...O<sub>acceptor</sub> distance decreases more significantly than the increase in the O<sub>donor</sub>-H distance (Fig. 2a). Thus, the H...O<sub>acceptor</sub> distance serves as the primary limiting factor for the O...O distance. For example, the O...O distances for the H-bond between water and protonated acetic acid ( $\Delta pK_a = 22$ , Fig. 2b), and between protonated acetic acid and deprotonated acetic acid ( $\Delta pK_a = 0$ , Fig. 2c), are 2.91 and 2.59 Å, respectively. O<sub>donor</sub>-H distances are 0.98 and 1.04 Å, with the latter being 0.06 Å longer than the former. In contrast, H...O<sub>acceptor</sub> distances are 1.97 and 1.55 Å, with the latter being 0.42 Å shorter than the former.

Since  $pK_a$  is an indicator of proton affinity (Brønsted acidity), a lower  $\Delta pK_a$  indicates a closer proton affinity between the H-bond donor and acceptor. In low  $\Delta pK_a$  H-bonds, the proton is strongly attracted to the H-bond acceptor, resulting in an increased O<sub>donor</sub>-H distance and a significantly decreased H...O<sub>acceptor</sub> distance.

Relationships among H-bond geometries and  $\Delta pK_a$  can also be explained in terms of electron redistribution induced by H-bond formation. A lower  $\Delta pK_a$  leads to pronounced electron redistribution from the H-bond acceptor to the donor (Fig. 2a). For example, electron redistribution amounts for the H-bond between water and protonated acetic acid ( $\Delta pK_a = 22$ , Fig. 2b), and between protonated acetic acid and deprotonated acetic acid ( $\Delta pK_a = 0$ , Fig. 2c), are 0.02e and 0.10e, respectively. Analysis of NBO<sup>28</sup> showed that this electron redistribution arises from the hybridization between the lone pair orbital of the H-bond acceptor ( $n_O$ ) and the O<sub>donor</sub>-H antibonding orbital of the H-bond donor ( $\sigma_{O-H}^*$ ).<sup>29</sup>

$\Delta pK_a$  is related to the energy difference between (i) the orbital of the donor O atom forming the O<sub>donor</sub>-H bond, and (ii) the orbital of the acceptor O atom accepting the O<sub>donor</sub>-H bond (*i.e.*, the  $n_O$  orbital) (Fig. 2b and c). When  $\Delta pK_a$  is high, the orbital energy of the donor O atom is much higher than that of the acceptor O atom (Fig. 2b), resulting in the H atom forming a bond with the donor O atom. As  $\Delta pK_a$  decreases, the energy difference between these orbitals decreases. In the extreme case where  $\Delta pK_a \sim 0$ , the energies of these orbitals are nearly equal (Fig. 2c). In this case, the H atom can bind to either the donor or acceptor O atom with almost no energy difference (a low-barrier H-bond, LBHB<sup>65,66</sup>). In LBHBs, the O...O distances are as short as  $\sim 2.5$  Å.<sup>67</sup>

Therefore, a lower  $\Delta pK_a$  leads to a decreased energy difference between the  $n_O$  orbital of the H-bond acceptor and the  $\sigma_{O-H}^*$  orbital of the H-bond donor, which enhances their hybridization (Fig. 2b and c). Indeed,  $pK_a$  values are correlated with the  $n_O$  and  $\sigma_{O-H}^*$  orbital energies (Lewis acidity, Fig. S2, ESI<sup>†</sup>).



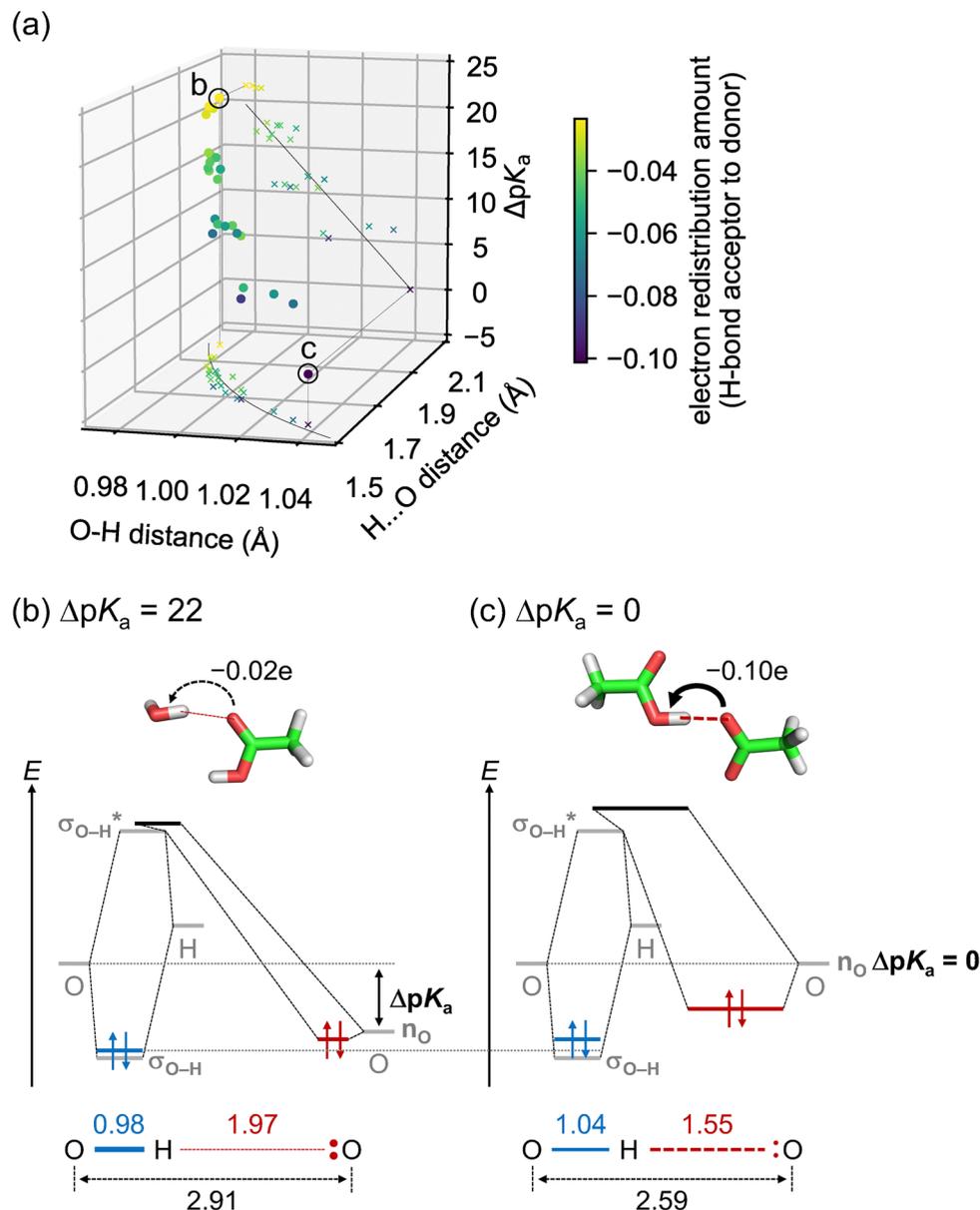


Fig. 2 Relationships among H-bond geometries and  $\Delta pK_a$ . (a)  $O_{\text{donor}}-H$  distance ( $r_{O-H}$ ),  $H \cdots O_{\text{acceptor}}$  distance ( $r_{H \cdots O}$ ), and  $\Delta pK_a$  for quantum-chemically optimized H-bonds in water solvent (circles). Relationships between  $r_{O-H}$  and  $r_{H \cdots O}$ , and between  $r_{O-H}$  and  $\Delta pK_a$ , are shown as crosses on the  $r_{O-H}-r_{H \cdots O}$  plane and the  $r_{O-H}-\Delta pK_a$  plane, respectively. Circles and crosses are color-scaled according to the electron redistribution amount from the H-bond acceptor to donor, calculated from the change in Mulliken charges due to H-bond formation. The curve on the  $r_{O-H}-r_{H \cdots O}$  plane indicates the correlation derived from the bond order model.<sup>60–64</sup> The parameters  $r_{O-H}^0$  and  $b$  in the model were set to 0.95 Å and 0.38 Å, respectively, to best fit the  $r_{O-H}$  versus  $r_{H \cdots O}$  relationship (see the discussion in ESI†). The line on the  $r_{O-H}-\Delta pK_a$  plane indicates the linear fit line ( $R^2 = 0.76$ ). (b) and (c) Energy diagrams of molecular orbitals for the H-bond: (b) between water and protonated acetic acid ( $\Delta pK_a = 22$ ), and (c) between protonated acetic acid and deprotonated acetic acid ( $\Delta pK_a = 0$ ). Gray bars indicate energies of the donor O atom and H atom orbitals forming the  $O_{\text{donor}}-H$  bond, as well as of the  $O_{\text{donor}}-H$  bonding ( $\sigma_{O-H}$ ),  $O_{\text{donor}}-H$  antibonding ( $\sigma_{O-H}^*$ ), and  $O_{\text{acceptor}}$  lone pair ( $n_O$ ) NBOs. Blue, black, and red bars indicate localized molecular orbitals formed by the three NBOs, which roughly correspond to the  $O_{\text{donor}}-H$  bonding orbital, the  $O_{\text{donor}}-H$  antibonding orbital, and the lone pair orbital of the acceptor O atom, respectively. Values in the lower panels indicate distances in Å.

This enhanced hybridization increases the  $n_O \rightarrow \sigma_{O-H}^*$  electron redistribution, leading to an increased  $\sigma_{O-H}^*$  orbital occupancy, a decreased  $O_{\text{donor}}-H$  bond order, and an increased  $O_{\text{donor}}-H$  distance. These correspond to the destabilization of the  $O_{\text{donor}}-H$  bonding orbital, resulting in a weaker  $O_{\text{donor}}-H$  bond (Fig. 2b and c). Furthermore, the enhanced hybridization between the  $n_O$  and  $\sigma_{O-H}^*$  orbitals, corresponding to the increased covalent

character of the  $H \cdots O_{\text{acceptor}}$  “bond”,<sup>64</sup> decreases the  $H \cdots O_{\text{acceptor}}$  distance more significantly than the  $O_{\text{donor}}-H$  distance increases. Thus, a lower  $\Delta pK_a$  results in a shorter  $O \cdots O$  distance.

Notably, the correlation between the  $O \cdots O$  distance and  $\Delta pK_a$  is mostly unaffected by the total charge of the H-bond (0 or  $-1$ ) (Fig. 1 and Fig. S5a, ESI†). Asp and Glu side-chains are



frequently involved in short H-bonds ( $O \cdots O$  distances  $< 2.7 \text{ \AA}$ ).<sup>4,8,9</sup> This is due to their tendency to form low  $\Delta pK_a$  H-bonds not only when (i) the deprotonated  $\text{COO}^-$  group serves as an acceptor, but also when (ii) the protonated  $\text{COOH}$  group serves as a donor (Fig. 1). The shorter  $O \cdots O$  distances are not caused by the stronger  $O_{\text{donor}}\text{-H} \cdots \text{O}_{\text{acceptor}}$  electrostatic interactions resulting from the negative charge of the deprotonated  $\text{COO}^-$  group. Indeed, the average  $O \cdots O$  distance of the charge-neutral  $[\text{COOH} \cdots \text{OH}_2]$  H-bonds with  $\Delta pK_a$  of  $\sim 6$  ( $2.59 \text{ \AA}$ ) is shorter than that of the negatively charged  $[\text{HOH} \cdots \text{OOC}]$  H-bonds with  $\Delta pK_a$  of  $\sim 12$  ( $2.77 \text{ \AA}$ ) for small compound H-bonds from the CSD (Table 2).<sup>58</sup> Although Asp/Glu side-chains are often assumed to be deprotonated,<sup>4,8</sup> the possibility of Asp/Glu side-chains serving as H-bond donors in their protonated forms should be considered, especially for short H-bonds.

### Distributions of H-bond distances involving the same donor/acceptor groups

Next, we analyze H-bonds in proteins to elucidate how the protein environment affects distances.  $O \cdots O$  distances for 313 680 O atom pairs with distances  $< 3.0 \text{ \AA}$  were obtained from 906 X-ray crystal structures. These structures were selected based on resolutions ( $\leq 1.2 \text{ \AA}$ ), sequence identities ( $\leq 90\%$ ), and  $R$ -factors ( $\leq 0.20$ ).<sup>9</sup> Among the 15 possible pairs of the five groups (Table 1),  $O \cdots O$  distance distributions for 14 pairs were analyzed, excluding the [backbone  $\cdots$  backbone] pair. Average  $O \cdots O$  distances for each pair were determined by fitting the distribution histograms (Fig. S6, ESI†) using a Gaussian function (Table 3).

The distributions of  $O \cdots O$  distances were compared with the  $\Delta pK_a$  values in water for each H-bond pair (Fig. 3). A lower  $\Delta pK_a$  leads to a shorter  $O \cdots O$  distance distribution. The  $O \cdots O$  distance distributions exhibit certain widths, with standard deviations ranging from 0.05 to 0.12  $\text{ \AA}$  (Table 3).

Several H-bond types with different  $\Delta pK_a$  values are possible for the [Asp/Glu  $\cdots$  Asp/Glu], [Tyr  $\cdots$  Asp/Glu], [Ser/Thr  $\cdots$  Asp/Glu], [water  $\cdots$  Asp/Glu], [Ser/Thr  $\cdots$  Tyr], and [water  $\cdots$  Tyr] pairs (Fig. 3). Among these six pairs, the predominant H-bond types for the [Asp/Glu  $\cdots$  Asp/Glu] and [water  $\cdots$  Asp/Glu] pairs can be deduced from their  $O \cdots O$  distance distributions. For the [Asp/Glu  $\cdots$  Asp/Glu] pair, the  $[\text{COOH} \cdots \text{OOC}]$  type with  $\Delta pK_a \sim 0$  is likely predominant. For the [water  $\cdots$  Asp/Glu] pair, the  $[\text{HOH} \cdots \text{OOC}]$  type with  $\Delta pK_a \sim 12$  is likely predominant (see the discussion in ESI†).

Average  $O \cdots O$  distances in proteins were compared with  $\Delta pK_a$  values in water. We compared the values of H-bond pairs whose  $\Delta pK_a$  values of the predominant H-bond types were inferred (indicated by the dotted squares in Fig. 3). The average  $O \cdots O$  distance ( $\overline{r_{O \cdots O}}$ ) is highly correlated with  $\Delta pK_a$  (Fig. 4a,  $R^2 = 0.91$ ), which is best described by the following equation:

$$\overline{r_{O \cdots O}} [\text{\AA}] = 0.0123\Delta pK_a + 2.55 \quad (1)$$

This high correlation indicates that  $O \cdots O$  distances with the same donor and acceptor groups are distributed around a value primarily determined by  $\Delta pK_a$  in water. This  $\Delta pK_a$  value does

Table 3  $\Delta pK_a$  values and  $O \cdots O$  distance distributions for H-bonds in proteins

H-Bond pair	Donor	Acceptor	$\Delta pK_a$ (in water)	$\overline{r_{O \cdots O}}^a$ ( $\text{\AA}$ )	$N^b$
Water $\cdots$ Water	H <sub>2</sub> O	H <sub>2</sub> O	18	$2.77 \pm 0.12$	149 940
Water $\cdots$ Ser/Thr	H <sub>2</sub> O	C-OH	18	$2.76 \pm 0.09$	17 411
		C-OH	18		
Water $\cdots$ Tyr	PhOH	H <sub>2</sub> O	12	$2.72 \pm 0.11$	5924
		PhOH	22		
Water $\cdots$ Asp/Glu	COOH	H <sub>2</sub> O	6	$2.73 \pm 0.10$	35 151
		COO <sup>-</sup>	12		
		COOH	22		
Water $\cdots$ Backbone	H <sub>2</sub> O	C=O	17	$2.79 \pm 0.09$	92 415
Ser/Thr $\cdots$ Ser/Thr	C-OH	C-OH	18	$2.75 \pm 0.08$	635
Ser/Thr $\cdots$ Tyr	PhOH	C-OH	12	$2.74 \pm 0.07$	329
		PhOH	22		
		COOH	22		
Ser/Thr $\cdots$ Asp/Glu	COOH	C-OH	6	$2.68 \pm 0.07$	2675
		COO <sup>-</sup>	12		
		COOH	22		
Ser/Thr $\cdots$ Backbone	C-OH	C=O	17	$2.75 \pm 0.09$	6220
Tyr $\cdots$ Tyr	PhOH	PhOH	16	$2.73 \pm 0.07$	72
Tyr $\cdots$ Asp/Glu	PhOH	COO <sup>-</sup>	6	$2.63 \pm 0.06$	1359
		PhOH	10		
Tyr $\cdots$ Backbone	PhOH	COOH	16	$2.69 \pm 0.06$	1313
		C=O	11		
Asp/Glu $\cdots$ Asp/Glu	COOH	COO <sup>-</sup>	0	$2.52 \pm 0.05^c$	254
		COOH	10		
Asp/Glu $\cdots$ Backbone	COOH	C=O	5	$2.64 \pm 0.05^c$	173

<sup>a</sup> Average  $O \cdots O$  distance and standard deviation, obtained by fitting the distribution histogram (Fig. S6, ESI) using a Gaussian function.

<sup>b</sup> Total number of O atom pairs. <sup>c</sup> Obtained from distributions of  $O \cdots O$  distances shorter than 2.75  $\text{ \AA}$ , as many of the O atom pairs with distances longer than 2.75  $\text{ \AA}$  are unlikely forming H-bonds (Fig. S6, ESI).

not account for the influence of the protein environment. Therefore, deviations in  $O \cdots O$  distances from their average values, corresponding to the width of the  $O \cdots O$  distance distribution (Fig. 3), reflect the influence of the protein environment on the characteristics of the H-bond.

### H-Bond distances for buried and solvent-exposed H-bonds

For small compound H-bonds, the distances are largely unaffected by differences in solvent dielectric constants.<sup>16,17</sup> The environment of H-bonds buried in the protein interior is low-dielectric, whereas that of H-bonds exposed to bulk water is high-dielectric.<sup>69,70</sup> To investigate the influence of the dielectric properties of the protein environment on  $O \cdots O$  distances, we compared the average  $O \cdots O$  distances of buried and exposed H-bonds. H-bonds were classified as buried or exposed according to their relative solvent accessibility values.<sup>37</sup>

The average  $O \cdots O$  distances of buried and exposed H-bonds are nearly identical for H-bond pairs with inferred predominant types (Fig. 4b and c, root mean square distance [RMSD] = 0.03  $\text{ \AA}$ ). This result indicates that  $O \cdots O$  distances are unaffected by whether the H-bond is buried in the protein interior or exposed to bulk water. This suggests that differences in the dielectric properties of the protein environment do not influence  $O \cdots O$  distances, consistent with the previous reports for small compounds.<sup>16,17</sup>

Water solvent weakens the electrostatic interaction between the  $O_{\text{donor}}\text{-H}$  and  $O_{\text{acceptor}}$  groups due to electrostatic shielding.



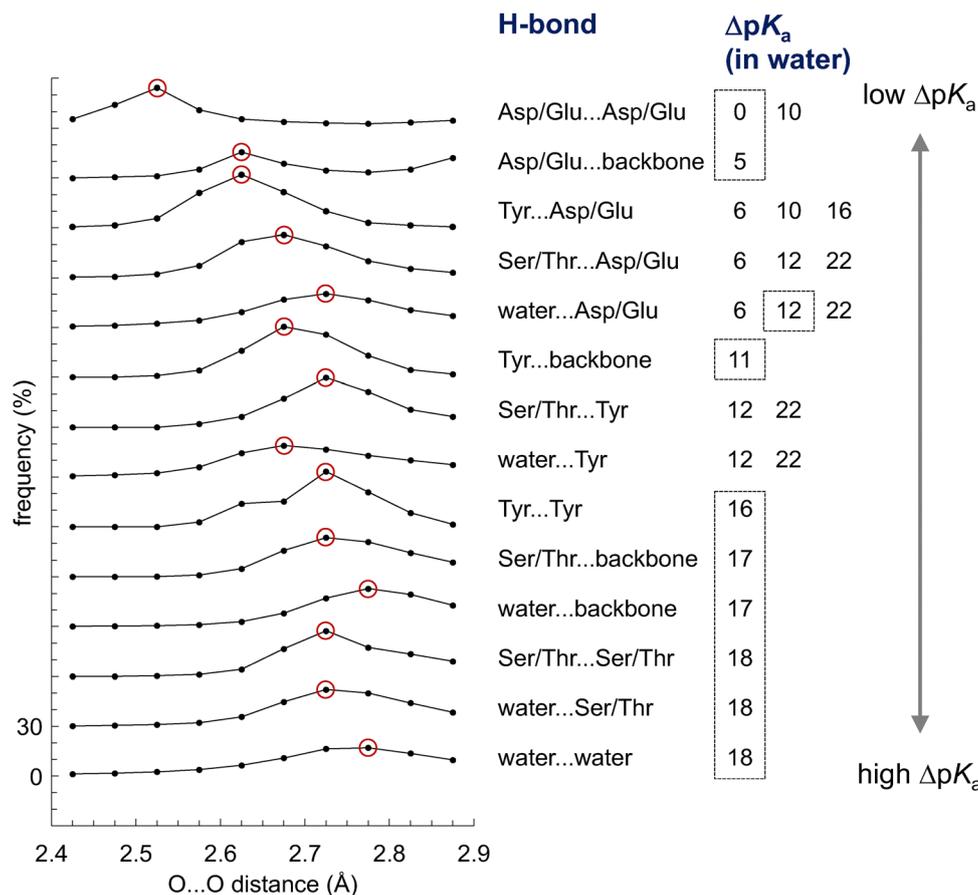


Fig. 3 Distributions of O...O distances for H-bonds in proteins. The bin width was set to 0.05 Å. Red circles indicate bins with the highest frequencies.  $\Delta pK_a$  values for the H-bond types included in each H-bond pair are shown on the right.  $\Delta pK_a$  values of the predominant H-bond types are enclosed by dotted boxes.

This electrostatic interaction is the major stabilizing factor of an H-bond.<sup>71</sup> On the other hand, the strength of the  $O_{\text{donor}}-H \cdots O_{\text{acceptor}}$  electrostatic interaction plays a minor role in determining the O...O distance. When comparing H-bonds with similar O...O distances, the electrostatic interaction energy is higher for H-bonds with total charges of  $-1$  than for charge-neutral H-bonds (Fig. S5b, ESI†). This indicates that a strong  $O_{\text{donor}}-H \cdots O_{\text{acceptor}}$  electrostatic interaction does not necessarily result in a shorter H-bond. This is because the stabilization provided by the electrostatic interaction is largely compensated by the destabilization due to exchange repulsion at short O...O distances (Fig. S5c, ESI†).<sup>62,64</sup> In contrast, a lower  $\Delta pK_a$  leads to an enhanced electron redistribution from the H-bond acceptor to the donor, which decreases the O...O distance (Fig. 2). Therefore, O...O distances are predominantly determined by the  $\Delta pK_a$  of the H-bond without being affected by whether the H-bond is buried in the protein interior or exposed to bulk water.

#### N...O and N...N distances

H-bonds between N and O atoms (N...O) and between N atoms (N...N) are also crucial in the H-bond network of proteins. From the same dataset of 906 protein structures, distances for

144 464 N...O atom pairs and 247 N...N atom pairs with distances  $< 3.2$  Å were obtained (see the discussion in ESI† for details).

The average N...O distance in proteins is highly correlated with  $\Delta pK_a$  in water for each H-bond type (Fig. 5a,  $R^2 = 0.71$ ). Similarly, the average N...N distance in proteins is highly correlated with  $\Delta pK_a$  in water (Fig. 5b,  $R^2 = 0.93$ ). These are consistent with the high correlations between N...O or N...N distances and  $\Delta pK_a$  observed for small-compound H-bonds.<sup>17</sup> The strength of the electrostatic interaction between the donor and acceptor groups differs among H-bonds formed between (i) charge-neutral donor and acceptor (*e.g.*, backbone-N-H...O=C-backbone), (ii) a positively charged donor and a negatively charged acceptor (salt-bridges, *e.g.*,  $Lys-NH_3^+ \cdots ^-OOC-Asp$ ), (iii) a positively charged donor and a charge-neutral acceptor (*e.g.*,  $Lys-NH_3^+ \cdots OH_2$ ), and (iv) a charge-neutral donor and a negatively charged acceptor (backbone-N-H... $^-OOC-Asp$ ). The correlations between average N...O or N...N distances and  $\Delta pK_a$  are largely unaffected by these differences (Fig. 5a and b), indicating that the  $\Delta pK_a$  of the H-bond, rather than the electrostatic interaction strength between the donor and acceptor groups, is the primary determinant of N...O and N...N distances, as well as O...O distances. Indeed, average N...O and N...N distances for buried



in proteins

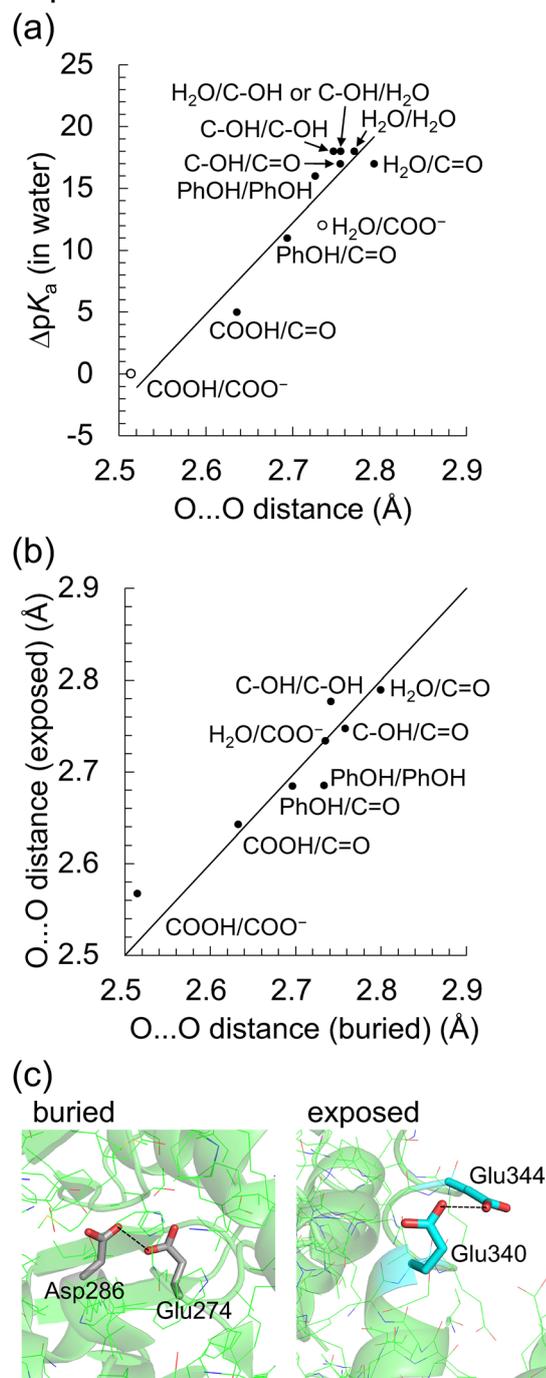


Fig. 4 Average O...O distances for each H-bond type in proteins. (a) Correlation between average O...O distances in proteins and  $\Delta pK_a$  in water ( $R^2 = 0.91$ ). Solid and open circles indicate H-bonds with total charges of 0 and  $-1$ , respectively. (b) Average O...O distances for buried and exposed H-bonds (RMSD = 0.03 Å). The black diagonal line indicates perfect correspondence (*i.e.*, identity line). (a) and (b) Labels indicate donor/acceptor groups. (c) Representative structures of buried (PDB ID: 4KQP<sup>68</sup>) and exposed (PDB ID: 3CLM, unpublished) H-bonds (dotted lines).

H-bonds (with low electrostatic shielding) and solvent-exposed H-bonds (with high electrostatic shielding) are nearly identical (Fig. 5c and d).

### Electrostatic and structural effects

H-bond distances with the same donor and acceptor groups are distributed around a value primarily determined by  $\Delta pK_a$  in water, regardless of whether the H-bond is buried in the protein interior or exposed to bulk water (Fig. 4 and 5). Deviations in H-bond distances from their average values reflect the influence of the protein environment. To investigate the influence of the protein environment in detail, we analyzed H-bonds between water and deprotonated Asp/Glu side-chains using a QM/MM approach. This  $[HOH \cdots ^-OOC]$  H-bond, with medium  $\Delta pK_a$  ( $\sim 12$ ) and average O...O distance (2.73 Å) values, is abundant in proteins (Table 3). Because  $[HOH \cdots ^-OOC]$  H-bonds are frequently found near the active site (retinal Schiff base) of microbial rhodopsins, we analyzed 12  $[HOH \cdots ^-OOC]$  H-bonds in  $H^+$ -pump (ground-state<sup>40</sup> and  $N'$ -state<sup>41</sup> BR),  $Cl^-$ -pump ( $pHR$ <sup>42</sup>), and  $Na^+$ -pump (KR2<sup>43</sup> and  $ErNaR$ <sup>44</sup>) rhodopsins.

$O_{\text{donor}}-H$  stretching vibrational frequencies of H-bonds in proteins are predominantly determined by the  $\Delta pK_a$  of the H-bond.<sup>72</sup> Here,  $\Delta pK_a$  values were estimated from the calculated  $O_{\text{donor}}-D$  stretching vibrational frequencies ( $\nu_{O-D}$ ) using the following equation:<sup>72</sup>

$$\Delta pK_a = 0.019 \nu_{O-D} [\text{cm}^{-1}] - 32 \quad (2)$$

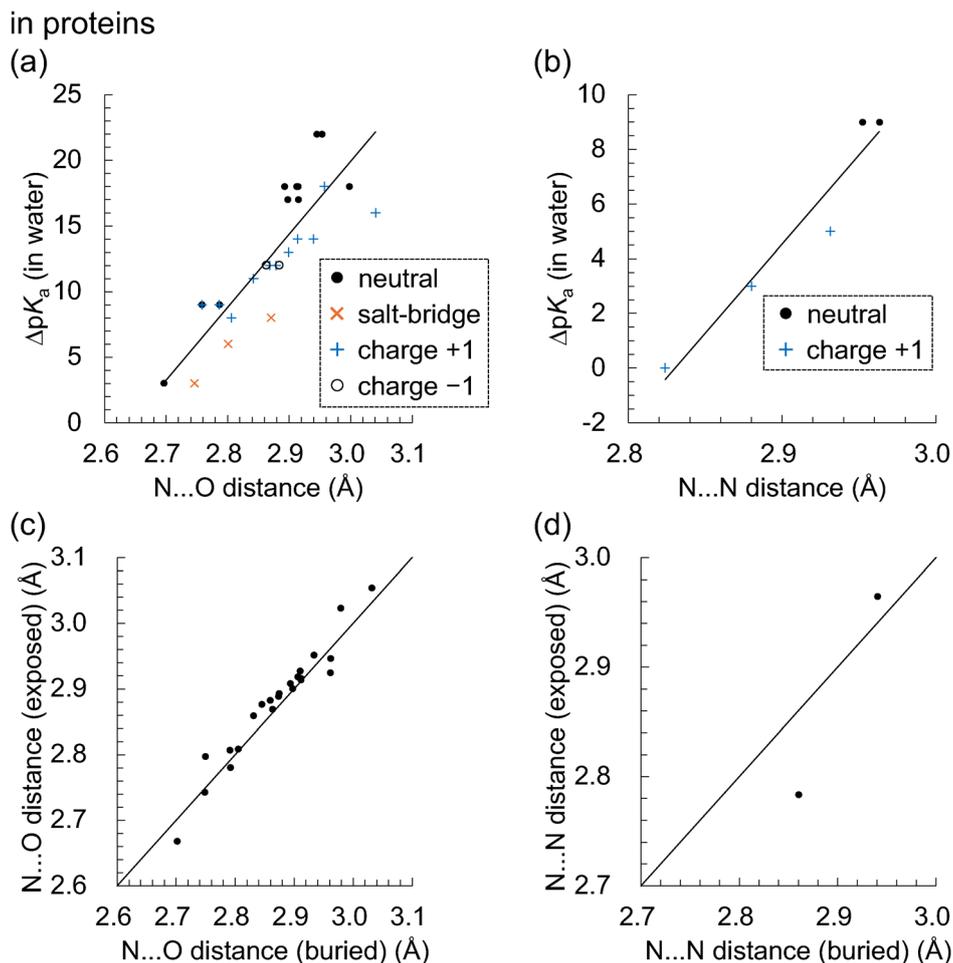
Note that this  $\Delta pK_a$  value, obtained from the  $O_{\text{donor}}-D$  stretching vibrational frequency, represents the  $\Delta pK_a$  value of the H-bond in the protein environment, not the  $\Delta pK_a$  value of  $\sim 12$  in water. O...O distances were obtained from the QM/MM-optimized structures. O...O distances and  $\Delta pK_a$  values were compared for the 12  $[HOH \cdots ^-OOC]$  H-bonds (Fig. 6a and Table S2, ESI<sup>†</sup>).

The O...O distance is correlated with  $\Delta pK_a$  for 12 H-bonds (Fig. 6a,  $R^2 = 0.79$ ). This relatively high correlation suggests that deviations in O...O distances from their average values primarily arise from  $\Delta pK_a$  shifts induced by the protein electrostatic environment.

Although  $O_{\text{donor}}-H$  stretching vibrational frequencies are mostly determined by the  $\Delta pK_a$  of the H-bond,<sup>72</sup> O...O distances are influenced more significantly by factors other than the  $\Delta pK_a$  of the H-bond. The correlation between  $\bar{r}_{O \cdots O}$  and  $\Delta pK_a$  for each H-bond type in proteins (eqn (1), the solid line in Fig. 6a) likely represents the relationship between the O...O distance and  $\Delta pK_a$ . While O...O distances and  $\Delta pK_a$  values for many H-bonds align with this  $\bar{r}_{O \cdots O}$  versus  $\Delta pK_a$  relationship, some H-bonds significantly deviate from this trend (Fig. 6a). These deviations arise from factors other than the  $\Delta pK_a$  of the H-bond. Among these, we focus on the  $[H_2O-406 \cdots Asp212]$  H-bond in the  $N'$ -state BR with a short O...O distance of 2.57 Å (Fig. 6b), and the  $[H_2O-402 \cdots Asp212]$  H-bond in the ground-state BR with a long O...O distance of 2.92 Å (Fig. 6c).

The O...O distance of the  $[H_2O-406 \cdots Asp212]$  H-bond in the  $N'$ -state BR is 2.57 Å, which is 0.16 Å shorter than the average O...O distance for  $[HOH \cdots ^-OOC]$  H-bonds in proteins (2.73 Å) (Fig. 6b). The calculated  $\Delta pK_a$  value of this H-bond is 13, nearly identical to the  $\Delta pK_a$  value of the  $[HOH \cdots ^-OOC]$  H-bond in the absence of the protein environment ( $\sim 12$ ) (Fig. 6a). Therefore,





**Fig. 5** Average N...O and N...N distances for each H-bond type in proteins. (a) Correlation between average N...O distances in proteins and  $\Delta pK_a$  in water ( $R^2 = 0.71$ ). (b) Correlation between average N...N distances in proteins and  $\Delta pK_a$  in water ( $R^2 = 0.93$ ). (a) and (b) Black closed circles, orange crosses, blue plus signs, and black open circles indicate H-bonds formed between (i) charge-neutral donor and acceptor, (ii) positively charged donor and negatively charged acceptor (salt-bridges), (iii) positively charged donor and charge-neutral acceptor, and (iv) charge-neutral donor and negatively charged acceptor, respectively. For H-bond pairs involving His side-chains, H-bond types with the lowest  $\Delta pK_a$  values are assumed (see the discussion in ESI†). (c) Average N...O distances for buried and exposed H-bonds (RMSD = 0.02 Å). (d) Average N...N distances for buried and exposed H-bonds (RMSD = 0.06 Å). (c) and (d) The black diagonal lines indicate perfect correspondence (*i.e.*, identity line). Data points for Trp...Ser/Thr, Trp...Tyr, Arg...His, Lys...His, and Trp...His pairs were excluded due to an insufficient number of exposed H-bonds (<10).

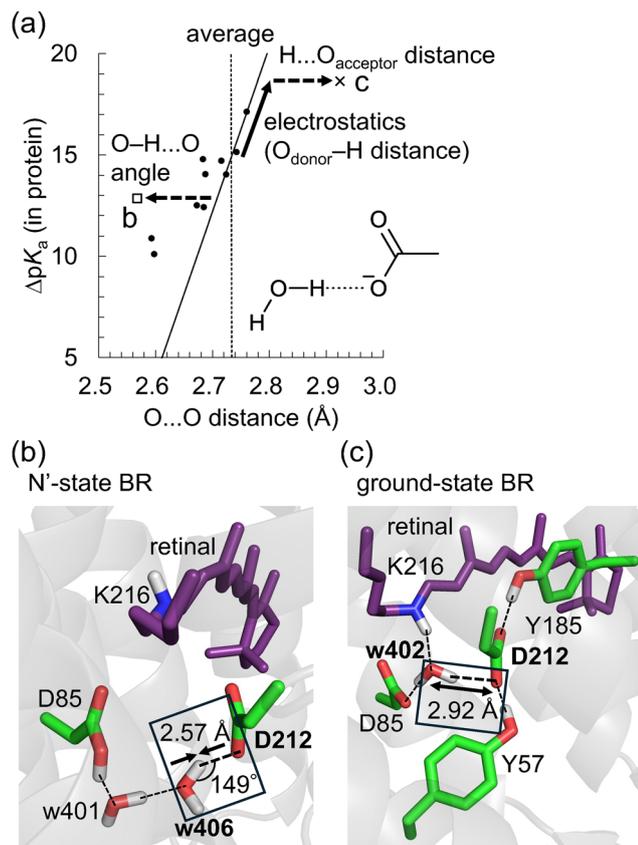
the short O...O distance of this H-bond is not caused by a decreased  $\Delta pK_a$  due to the protein electrostatic environment, but rather by structural constraints in the protein environment.

In the absence of the protein environment, a lower  $\Delta pK_a$  tends to result in a higher binding enthalpy,<sup>11</sup> a shorter O...O distance, and a larger  $O_{\text{donor}}\text{-H}\cdots O_{\text{acceptor}}$  angle approaching  $180^\circ$ .<sup>16,17</sup> On the other hand, the  $O_{\text{donor}}\text{-H}\cdots O_{\text{acceptor}}$  angle of this short H-bond is  $149^\circ$  (Fig. 6b). Short H-bonds with decreased  $O_{\text{donor}}\text{-H}\cdots O_{\text{acceptor}}$  angles, despite high  $\Delta pK_a$  values, are frequently observed in intramolecular H-bonds, where structural constraints are significant.<sup>17</sup> Such short H-bonds with small  $O_{\text{donor}}\text{-H}\cdots O_{\text{acceptor}}$  angles result from rigid anchoring in the protein matrix, which causes deviations from the equilibrium O...O distance. In the case of the [H<sub>2</sub>O-406...Asp212] H-bond in the N'-state BR, the formation of the H-bond network involving protonated Asp85, H<sub>2</sub>O-401, H<sub>2</sub>O-406, and Asp212 decreases the O...O distance (Fig. 6b). Indeed,

when the structure is QM/MM-optimized in the absence of H<sub>2</sub>O-401, the O...O distance of the [H<sub>2</sub>O-406...Asp212] H-bond increases to 2.79 Å (Fig. S7, ESI†).

The O...O distance of the [H<sub>2</sub>O-402...Asp212] H-bond in the ground-state BR is 2.92 Å, which is 0.19 Å longer than the average O...O distance for [HOH...OOC] H-bonds in proteins (2.73 Å) (Fig. 6c). The calculated  $\Delta pK_a$  value of this H-bond is 19, indicating that  $\Delta pK_a$  is increased from ~12 due to the protein electrostatic environment (Fig. 6a). However, a  $\Delta pK_a$  of 19 corresponds to an O...O distance of 2.78 Å based on eqn (1), which is shorter than 2.92 Å. Therefore, the long O...O distance of 2.92 Å cannot be solely attributed to the increased  $\Delta pK_a$ . Instead, it results from both (i) the increased  $\Delta pK_a$  due to the protein environment, and (ii) structural constraints in the protein environment (Fig. 6a). The increase in  $\Delta pK_a$  is due to the decrease in  $pK_a$  of Asp212, caused by H-bond donations from Tyr57 and Tyr185 to Asp212<sup>73</sup> (Fig. 6c). The structural





**Fig. 6** O...O distances for [HOH...OOC] H-bonds in microbial rhodopsins. (a) O...O distances and  $\Delta pK_a$  for 12 [HOH...OOC] H-bonds. The open square, the cross, and circles indicate O...O distances and  $\Delta pK_a$  values of the [H<sub>2</sub>O-406...Asp212] H-bond in the N'-state BR, [H<sub>2</sub>O-402...Asp212] H-bond in the ground-state BR, and the other 10 H-bonds, respectively. The solid line indicates the correlation between  $\overline{O\cdots O}$  and  $\Delta pK_a$  (eqn (1)). The vertical dotted line indicates the average O...O distance for the [HOH...OOC] H-bonds in proteins (Table 3). The solid arrow indicates the shift in O...O distance due to the  $\Delta pK_a$  shift. The dotted arrows indicate the shift in O...O distance due to structural constraints in proteins. (b) [H<sub>2</sub>O (w)-406...Asp212] H-bond in the N'-state BR. (c) [H<sub>2</sub>O (w)-402...Asp212] H-bond in the ground-state BR. (b) and (c) Dotted lines indicate H-bonds. Arrows indicate deviations in O...O distances from the average value in proteins.

constraints likely arise from the H-bond formations between H<sub>2</sub>O-402 and Asp85/Lys216, and between Asp212 and Tyr57/Tyr185, which may increase the H...O<sub>acceptor</sub> distance (Fig. 6c).

To summarize, H-bond distances are influenced by the following two factors of the protein environment: (i) the protein electrostatic environment that shifts the  $\Delta pK_a$  of the H-bond. This  $\Delta pK_a$  shift alters the covalent character of the H...O<sub>acceptor</sub> bond (Fig. 2), thereby altering the O...O distances. (ii) Structural constraints imposed by protein folding. These constraints alter the O<sub>donor</sub>-H...O<sub>acceptor</sub> angle (Fig. 6b) or the H...O<sub>acceptor</sub> distance (Fig. 6c) without affecting  $\Delta pK_a$  (or the corresponding O<sub>donor</sub>-H distance), thereby altering the O...O distance. Such constraints destabilize the H-bond, which is likely compensated by stabilization through other interactions and H-bonds (e.g., H-bonds between H<sub>2</sub>O-402 and Asp85/Lys216, and between Asp212 and Tyr57/Tyr185 in the ground-state BR, Fig. 6c).

## Conclusions

In H-bonds, while the O<sub>donor</sub>-H distance remains around  $\sim 1$  Å, the H...O<sub>acceptor</sub> distance is significantly longer and serves as the primary limiting factor for the O...O distance. Thus, minimizing the H...O<sub>acceptor</sub> distance decreases the overall H-bond distance. The present study demonstrates that electron redistribution counteracts proton migration from the H-bond donor to the acceptor upon H-bond formation. This electron redistribution increases as  $\Delta pK_a$  decreases, introducing a covalent character to the H...O<sub>acceptor</sub> “bond”, thereby significantly decreasing the overall O...O distance (Fig. 2). In contrast, the strength of the electrostatic interaction between the donor and acceptor groups plays a minor role in determining H-bond distances, as indicated by the small dependence of H-bond distances on the total charges of the donor and acceptor groups (Fig. 1, 4 and 5). Therefore, the protein electrostatic environment influences H-bond distances primarily by altering the  $\Delta pK_a$  of the H-bond (Fig. 6a), whereas differences in the dielectric properties of the protein environment do not affect H-bond distances (Fig. 4 and 5).

Beyond this primary determinant of H-bond distance, structural constraints in the protein environment impose a secondary influence. O<sub>donor</sub>-H...O<sub>acceptor</sub> angles (Fig. 6b) and H...O<sub>acceptor</sub> distances (Fig. 6c) are often restricted by rigid anchoring in the protein matrix, causing deviations in the O...O distance from the equilibrium value determined by  $\Delta pK_a$ . While the effect of structural constraints on H-bond distances is weaker than the covalent-bond-like electronic effect, it enables the formation of unique H-bond geometries that are inaccessible in bulk solvent (Fig. 6a). These protein-specific H-bond geometries may play a crucial role in shaping the functional properties of proteins.

## Data availability

The data supporting the findings of this study are available in the main article and the ESI.†

## Conflicts of interest

There are no conflicts of interest to declare.

## Acknowledgements

This research was supported by JSPS KAKENHI (JP22KJ1109 to M. T.; JP23H04963 and JP24K01986 to K. S.; JP23H02444 to H. I.) and the Interdisciplinary Computational Science Program in CCS, University of Tsukuba (K. S.).

## References

- 1 L. Pauling, *The nature of the chemical bond and the structure of molecules and crystals: an introduction to modern structural chemistry*, Cornell University Press, Ithaca, NY, 1960.



- 2 H. M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weissig, I. N. Shindyalov and P. E. Bourne, *Nucleic Acids Res.*, 2000, **28**, 235–242.
- 3 M. M. Flocco and S. L. Mowbray, *J. Mol. Biol.*, 1995, **254**, 96–105.
- 4 S. Rajagopal and S. Vishveshwara, *FEBS J.*, 2005, **272**, 1819–1832.
- 5 G. Wohlfahrt, *Proteins*, 2005, **58**, 396–406.
- 6 A. Langkilde, S. M. Kristensen, L. Lo Leggio, A. Mølgaard, J. H. Jensen, A. R. Houk, J.-C. N. Poulsen, S. Kauppinen and S. Larsen, *Acta Crystallogr., Sect. D: Biol. Crystallogr.*, 2008, **64**, 851–863.
- 7 J. Lin, E. Pozharski and M. A. Wilson, *Biochemistry*, 2017, **56**, 391–402.
- 8 S. Zhou and L. Wang, *Chem. Sci.*, 2019, **10**, 7734–7745.
- 9 H. W. Qi and H. J. Kulik, *J. Chem. Inf. Model.*, 2019, **59**, 2199–2211.
- 10 P. Huyskens and T. Zeegers-Huyskens, *J. Chem. Phys.*, 1964, **61**, 81–86.
- 11 P. Gilli, L. Pretto, V. Bertolasi and G. Gilli, *Acc. Chem. Res.*, 2009, **42**, 33–44.
- 12 B. Brycki and M. Szafran, *J. Chem. Soc., Perkin Trans. 2*, 1984, 223–226.
- 13 P. Huyskens, L. Sobczyk and I. Majerz, *J. Mol. Struct.*, 2002, **615**, 61–72.
- 14 B. Brycki, B. Brzezinski, G. Zundel and T. Keil, *Magn. Reson. Chem.*, 1992, **30**, 507–510.
- 15 J. B. Tobin, S. A. Whitt, C. S. Cassidy and P. A. Frey, *Biochemistry*, 1995, **34**, 6919–6924.
- 16 P. A. Sigala, E. A. Ruben, C. W. Liu, P. M. Piccoli, E. G. Hohenstein, T. J. Martinez, A. J. Schultz and D. Herschlag, *J. Am. Chem. Soc.*, 2015, **137**, 5730–5740.
- 17 D. Herschlag and M. M. Pinney, *Biochemistry*, 2018, **57**, 3338–3352.
- 18 M. M. Pinney, A. Natarajan, F. Yabukarski, D. M. Sanchez, F. Liu, R. B. Liang, T. Doukov, J. P. Schwans, T. J. Martinez and D. Herschlag, *J. Am. Chem. Soc.*, 2018, **140**, 9827–9843.
- 19 J. R. Rumble, *CRC handbook of chemistry and physics*, CRC Press, Boca Raton, FL, 103rd edn, 2022.
- 20 M. B. Smith and J. March, *March's advanced organic chemistry*, John Wiley & Sons, Inc., Hoboken, NJ, 5th edn, 2001.
- 21 M. Śmiechowski, E. Gojlo and J. Stangret, *J. Phys. Chem. B*, 2011, **115**, 4834–4842.
- 22 M. A. Brown, F. Vila, M. Sterrer, S. Thürmer, B. Winter, M. Ammann, J. J. Rehr and J. A. van Bokhoven, *J. Phys. Chem. Lett.*, 2012, **3**, 1754–1759.
- 23 K. Kitaura, E. Ikeo, T. Asada, T. Nakano and M. Uebayasi, *Chem. Phys. Lett.*, 1999, **313**, 701–706.
- 24 D. G. Fedorov and K. Kitaura, *J. Chem. Phys.*, 2004, **121**, 2483–2490.
- 25 D. G. Fedorov and K. Kitaura, *J. Comput. Chem.*, 2007, **28**, 222–237.
- 26 D. G. Fedorov, K. Kitaura, H. Li, J. H. Jensen and M. S. Gordon, *J. Comput. Chem.*, 2006, **27**, 976–985.
- 27 M. W. Schmidt, K. K. Baldridge, J. A. Boatz, S. T. Elbert, M. S. Gordon, J. H. Jensen, S. Koseki, N. Matsunaga, K. A. Nguyen, S. Su, T. L. Windus, M. Dupuis and J. A. Montgomery, *J. Comput. Chem.*, 1993, **14**, 1347–1363.
- 28 J. P. Foster and F. Weinhold, *J. Am. Chem. Soc.*, 1980, **102**, 7211–7218.
- 29 F. Weinhold and C. Landis, *Valency and bonding. A natural bond orbital donor–acceptor perspective*, Cambridge University Press, Cambridge, 2005.
- 30 J. S. Arey, P. C. Aeberhard, I.-C. Lin and U. Rothlisberger, *J. Phys. Chem. B*, 2009, **113**, 4726–4732.
- 31 R. Misra, D. K. Maity and S. P. Bhattacharyya, *Chem. Phys.*, 2012, **402**, 96–104.
- 32 T. Yanai, D. P. Tew and N. C. Handy, *Chem. Phys. Lett.*, 2004, **393**, 51–57.
- 33 E. D. Glendening, J. K. Badenhoop, A. E. Reed, J. E. Carpenter, J. A. Bohmann, C. M. Morales and F. Weinhold, *NBO 5.0*, Theoretical Chemistry Institute, University of Wisconsin, Madison, WI, 2001.
- 34 A. D. Bochevarov, E. Harder, T. F. Hughes, J. R. Greenwood, D. A. Braden, D. M. Philipp, D. Rinaldo, M. D. Halls, J. Zhang and R. A. Friesner, *Int. J. Quantum Chem.*, 2013, **113**, 2110–2142.
- 35 T. Hamelryck and B. Manderick, *Bioinformatics*, 2003, **19**, 2308–2310.
- 36 P. J. A. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A. Dalke, I. Friedberg, T. Hamelryck, F. Kauff, B. Wilczynski and M. J. L. de Hoon, *Bioinformatics*, 2009, **25**, 1422–1423.
- 37 B. Rost and C. Sander, *Proteins*, 1994, **20**, 216–226.
- 38 W. Kabsch and C. Sander, *Biopolymers*, 1983, **22**, 2577–2637.
- 39 W. G. Touw, C. Baakman, J. Black, T. A. H. te Beek, E. Krieger, R. P. Joosten and G. Vriend, *Nucleic Acids Res.*, 2015, **43**, D364–D368.
- 40 N. Hasegawa, H. Jonotsuka, K. Miki and K. Takeda, *Sci. Rep.*, 2018, **8**, 13123.
- 41 B. Schobert, L. S. Brown and J. K. Lanyi, *J. Mol. Biol.*, 2003, **330**, 553–570.
- 42 T. Kouyama, S. Kanada, Y. Takeguchi, A. Narusawa, M. Murakami and K. Ihara, *J. Mol. Biol.*, 2010, **396**, 564–579.
- 43 K. Kovalev, R. Astashkin, I. Gushchin, P. Orekhov, D. Volkov, E. Zinovev, E. Marin, M. Rulev, A. Alekseev, A. Royant, P. Carpentier, S. Vaganova, D. Zabelskii, C. Baeken, I. Sergeev, T. Balandin, G. Bourenkov, X. Carpena, R. Boer, N. Maliar, V. Borschchevskiy, G. Buldt, E. Bamberg and V. Gordeliy, *Nat. Commun.*, 2020, **11**, 2137.
- 44 E. Podoliak, G. H. U. Lamm, E. Marin, A. V. Schellbach, D. A. Fedotov, A. Stetsenko, M. Asido, N. Maliar, G. Bourenkov, T. Balandin, C. Baeken, R. Astashkin, T. R. Schneider, A. Bateman, J. Wachtveitl, I. Schapiro, V. Busskamp, A. Guskov, V. Gordeliy, A. Alekseev and K. Kovalev, *Nat. Commun.*, 2024, **15**, 3119.
- 45 B. R. Brooks, R. E. Bruccoleri, B. D. Olafson, D. J. States, S. Swaminathan and M. Karplus, *J. Comput. Chem.*, 1983, **4**, 187–217.
- 46 A. D. MacKerell, D. Bashford, M. Bellott, R. L. Dunbrack, J. D. Evanseck, M. J. Field, S. Fischer, J. Gao, H. Guo, S. Ha, D. Joseph-McCarthy, L. Kuchnir, K. Kuczera, F. T. Lau, C. Mattos, S. Michnick, T. Ngo, D. T. Nguyen, B. Prodhom,



- W. E. Reiher, B. Roux, M. Schlenkrich, J. C. Smith, R. Stote, J. Straub, M. Watanabe, J. Wiorkiewicz-Kuczera, D. Yin and M. Karplus, *J. Phys. Chem. B*, 1998, **102**, 3586–3616.
- 47 D. Bashford and M. Karplus, *Biochemistry*, 1990, **29**, 10219–10225.
- 48 Y. Nozaki and C. Tanford, *Methods Enzymol.*, 1967, **11**, 715–734.
- 49 M. Tanokura, *Biochim. Biophys. Acta, Protein Struct. Mol. Enzymol.*, 1983, **742**, 576–585.
- 50 M. Tanokura, *Biochim. Biophys. Acta, Protein Struct. Mol. Enzymol.*, 1983, **742**, 586–596.
- 51 M. Tanokura, *J. Biochem.*, 1983, **94**, 51–62.
- 52 B. Rabenstein and E.-W. Knapp, *Biophys. J.*, 2001, **80**, 1141–1150.
- 53 D. M. Philipp and R. A. Friesner, *J. Comput. Chem.*, 1999, **20**, 1468–1494.
- 54 R. B. Murphy, D. M. Philipp and R. A. Friesner, *J. Comput. Chem.*, 2000, **21**, 1442–1457.
- 55 W. L. Jorgensen, D. S. Maxwell and J. Tirado-Rives, *J. Am. Chem. Soc.*, 1996, **118**, 11225–11236.
- 56 A. P. Scott and L. Radom, *J. Phys. Chem.*, 1996, **100**, 16502–16513.
- 57 T. Steiner, *Angew. Chem., Int. Ed.*, 2002, **41**, 48–76.
- 58 L. D'Ascenzo and P. Auffinger, *Acta Crystallogr., Sect. B: Struct. Sci., Cryst. Eng. Mater.*, 2015, **71**, 164–175.
- 59 C. R. Groom, I. J. Bruno, M. P. Lightfoot and S. C. Ward, *Acta Crystallogr., Sect. B: Struct. Sci., Cryst. Eng. Mater.*, 2016, **72**, 171–179.
- 60 L. Pauling, *J. Am. Chem. Soc.*, 1947, **69**, 542–553.
- 61 T. Steiner and W. Saenger, *Acta Crystallogr., Sect. B: Struct. Sci., Cryst. Eng. Mater.*, 1994, **50**, 348–357.
- 62 P. Gilli, V. Bertolasi, V. Ferretti and G. Gilli, *J. Am. Chem. Soc.*, 1994, **116**, 909–915.
- 63 H.-H. Limbach, P. M. Tolstoy, N. Pérez-Hernández, J. Guo, I. G. Shenderovich and G. S. Denisov, *Isr. J. Chem.*, 2009, **49**, 199–216.
- 64 S. J. Grabowski, *Chem. Rev.*, 2011, **111**, 2597–2625.
- 65 W. W. Cleland, *Biochemistry*, 1992, **31**, 317–319.
- 66 C. N. Schutz and A. Warshel, *Proteins*, 2004, **55**, 711–723.
- 67 W. W. Cleland, *Arch. Biochem. Biophys.*, 2000, **382**, 1–5.
- 68 F. Fulyani, G. K. Schuurman-Wolters, A. V. Zagar, A. Guskov, D.-J. Slotboom and B. Poolman, *Structure*, 2013, **21**, 1879–1888.
- 69 G. M. Ullmann and E.-W. Knapp, *Eur. Biophys. J.*, 1999, **28**, 533–551.
- 70 C. N. Schutz and A. Warshel, *Proteins*, 2001, **44**, 400–417.
- 71 K. Morokuma, *Acc. Chem. Res.*, 1977, **10**, 294–300.
- 72 M. Tsujimura, K. Saito and H. Ishikita, *Biophys. J.*, 2023, **122**, 4336–4347.
- 73 K. Saito, H. Kandori and H. Ishikita, *J. Biol. Chem.*, 2012, **287**, 34009–34018.

