

Water regulates the residence time of Benzamidine in Trypsin

Narjes Ansari

Italian Institute of Technology

Valerio Rizzi

Italian Institute of Technology <https://orcid.org/0000-0001-5126-8996>

Michele Parrinello (✉ Michele.parrinello@iit.it)

Italian Institute of technology

Article

Keywords:

Posted Date: April 20th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1546274/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Communications on September 16th, 2022. See the published version at <https://doi.org/10.1038/s41467-022-33104-3>.

Water regulates the residence time of Benzamidine in Trypsin

Narjes Ansari¹, Valerio Rizzi¹, and Michele Parrinello^{1,*}

¹Italian Institute of Technology, Via E. Melen 83, 16152, Genova, Italy

*michele.parrinello@iit.it

April 11, 2022

1 Abstract

2 We simulate with state-of-the-art enhanced sampling techniques the binding of Benzamidine to Trypsin which is a much studied
3 and paradigmatic ligand-protein system. We use machine learning methods and in particular Time-lagged Independent Component
4 Analysis to determine efficient collective coordinates. These coordinates are used to perform On-the-fly Probability Enhanced
5 Sampling simulations, which we adapt to calculate also the ligand residence time. Our results, both static and dynamic, are in good
6 agreement with experiments. We underline the role of water in the unbinding process and find that the presence of a water molecule
7 located at the bottom of the binding pocket allows via a network of hydrogen bonds the ligand to be released into the solution. On
8 a finer scale, even when unbinding is allowed, another water molecule further modulates the exit time.

9 Introduction

10 The process of ligand-protein unbinding is crucial in biophysics and its full understanding would enhance
11 not only our knowledge but also benefit the design of new drugs. In particular, two quantities are of great
12 relevance, the ligand binding free energy and the inverse of the residence time k_{off} [1]. Being able to
13 compute reliably these two quantities would be of great help. Here we describe a strategy that makes
14 it possible to calculate accurately both quantities. We demonstrate this assertion with a state-of-the-art
15 simulation of the Trypsin-Benzamidine system [2, 3, 4, 5, 6] and we find that the unbinding process
16 can be more complex than what one could have anticipated. Although Trypsin-Benzamidine is one of
17 the simplest cases of protein-ligand systems, high-resolution crystallographic experiments have recently
18 demonstrated the presence of an extensive water structure in the binding cavity [7]. Our study reveals
19 that water has not only a structural role but a dynamical one, since it regulates the unbinding process via
20 a complex rearrangement of hydrogen bonds (HBs). In particular, the presence of a water molecule at a
21 specific position in the binding cavity allows the unbinding process to take place. As the ligand begins to
22 leave its binding pose, the number of water molecules in the binding cavity increases, leading to a finer
23 regulation in which water determines two possible escape pathways with a k_{off} differing by one order of
24 magnitude.

25 Several technical advances have made this study possible. Ligand unbinding is a rare event that takes
26 place on a timescale of milliseconds and thus its study requires the use of an enhanced sampling method.
27 Here, we profit from the flexibility and efficiency of the On-the-fly Probability Enhanced Sampling (OPES)
28 method [8], that is the latest evolution of Metadynamics [9, 10]. Like umbrella sampling [11] and other
29 enhanced sampling methods [12], OPES is based on the use of a set of collective variables (CVs) that are
30 functions $s(R)$ of the atomic coordinates R and describe the slow modes of the system. A good choice
31 of $s(R)$ is essential to obtain converged results in an affordable time. To determine good CVs we use

32 two machine-learning-based tools that we recently developed, Deep Local Discriminant Analysis (Deep-
33 LDA) [13] and Deep Time-lagged Independent Component Analysis (Deep-TICA) [14]. In building these
34 CVs, we shall pay great attention to the role of water and make use of the experience gained in Ref. 15.

35 Standard OPES is very efficient in calculating static properties such as binding free energies, but, in
36 doing so, it alters the natural dynamics of the system so that a sensitive quantity like k_{off} cannot be easily
37 extracted. Nevertheless, it has been pointed out that if an enhanced sampling method can be engineered
38 such that no bias is added to the transition region, the value of k_{off} can still be computed [16, 17, 18, 19, 20].
39 Here we show that by an appropriate setting of the input parameters, OPES can be made to satisfy this
40 condition. We call this approach OPES flooding (OPES_f) since it is inspired by the flooding approach [19]
41 and use it to calculate the ligand residence time.

42 Results

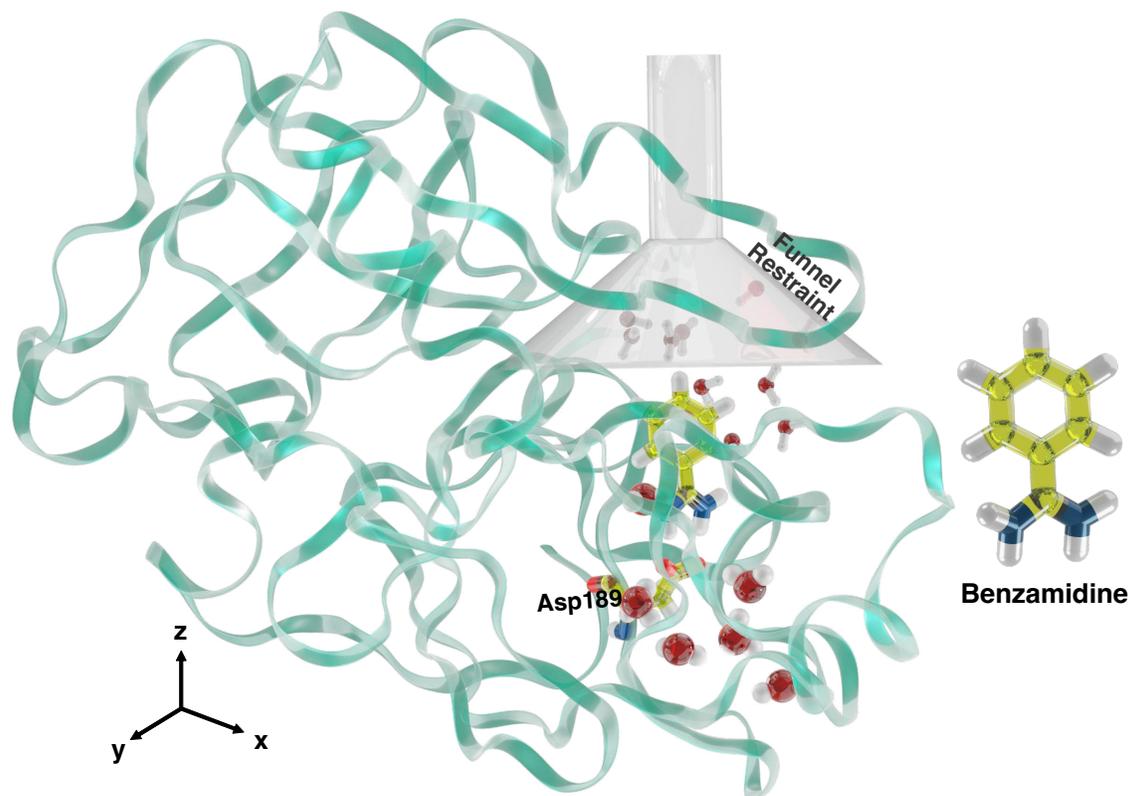


Figure 1: **The Trypsin-Benzamidine system.** A cartoon representation of Trypsin structure with the ligand Benzamidine and the Funnel restraint geometry.

43 Enhanced sampling technique for estimating free energies: OPES

44 To accelerate the occurrence of binding and unbinding events, we use OPES [8]. This method allows the
45 system to overcome kinetic barriers by transforming the original CV probability distribution $P(s)$ into a
46 smoother and thus simpler to sample one $P^{\text{tg}}(s)$. In the variant of OPES that we use here, we transform the
47 s probability distribution $P(s) = \frac{\int \delta(s-s(R))e^{-\beta V(R)} dR}{Z}$, where $\beta = 1/k_B T$ is the inverse temperature, $V(R)$

48 the interaction potential and $Z = \int e^{-\beta V(R)} dR$ the partition function, into the smoother well-tempered
 49 distribution $P^{\text{tg}}(s) \propto [P(s)]^{1/\gamma}$ [10], where the parameter $\gamma > 1$ regulates the broadening of the target
 50 distribution. In standard metadynamics this transformation is achieved by iteratively building a bias
 51 potential $V(s)$, while in OPES one instead reconstructs the probability distribution $P(s)$ on-the-fly. At
 52 iteration step n , the probability distribution $P_n(s)$ is estimated as

$$P_n(s) = \frac{\sum_k^n w_k G(s, s_k)}{\sum_k^n w_k} \quad (1)$$

53 where $G(s, s_k)$ are multivariate Gaussian kernels evaluated at every step k , the weights $w_k = e^{\beta V_{k-1}(s_k)}$ are
 54 computed from the bias at step $k - 1$. In turn, the bias is

$$V_n(s) = \left(1 - \frac{1}{\gamma}\right) \frac{1}{\beta} \log \left(\frac{P_n(s)}{Z_n} + \epsilon \right) \quad (2)$$

55 where Z_n is a factor that measures the configuration space thus far explored and ϵ is a very important
 56 parameter that controls the maximum bias that can be deposited in the system.

57 Machine Learning-based CVs: Deep-LDA and Deep-TICA

58 As in our previous work [15], we use the machine learning-based Deep-LDA method [13] to design
 59 effective CVs to be used in conjunction with OPES. The method is based on the time-honored Linear
 60 Discriminant Analysis technique [21] employed in classification applications. Our aim is to build a CV
 61 that is able to distinguish between two sets of data. In our case, data come from unbiased simulations of
 62 the system in the bound (B) and unbound (U) states.

63 In standard LDA, one optimizes Fisher's ratio $\frac{w^T S_b w}{w^T S_w w}$ to obtain the linear combination $s(R) = w^T \mathbf{d}(R)$
 64 of descriptors $\mathbf{d}(R)$ that best separates the two states. Fisher's ratio is written in terms of the scatter
 65 matrix $S_b = (\mu_B - \mu_U)(\mu_B - \mu_U)^T$ and the within matrix $S_w = S_B + S_U$, where μ_B, μ_U indicate the
 66 average descriptors values and S_B, S_U the descriptors variance matrices in the two states. The vector w
 67 that maximizes this ratio is the direction that optimally discriminates the states and provides the best-
 68 separated projection of the data in the one-dimensional s space [22].

69 In Deep-LDA, the set of N_d descriptors \mathbf{d} is fed into a neural network (NN) that is trained by applying
 70 the LDA criterion to the last hidden layer \mathbf{h} of the network. In analogy with LDA, a projection of the N_h
 71 components of the last hidden layer produces the Deep-LDA CV $s = w^T \mathbf{h}$. This CV, being by construction
 72 a non-linear combination of the original input descriptors \mathbf{d} , is more expressive than a simple linear
 73 combination and has been successfully applied to solve a number of problems [15, 23, 24]. As s tends to
 74 produce sharp distributions, we apply the cubic transformation $s_w = s + s^3$ [25, 15] to make the CV more
 75 easily applicable to enhanced simulations.

76 While often successful in driving a system forth and back between different states, a Deep-LDA CV
 77 does not encode any information on the transition state. This information could have been obtained if
 78 one had access to a long dynamical trajectory in which many state-to-state transitions did occur. In this
 79 case, a way of building a good CV, would have been to use the Time-lagged Independent Component
 80 Analysis [26, 27], in which one looks for the most slowly decorrelating modes. These slow modes can
 81 be found using the variational principle [28]. If the modes are expressed as a linear combination of
 82 descriptors, the variational principle leads to a generalized eigenvalue equation. As in Deep-LDA, one
 83 can apply TICA not only to a linear combination of descriptors but also to the last hidden layer of a
 84 NN. This greatly improves the variational flexibility of the solution and thus its quality. In its original
 85 formulation, this approach was meant to be applied to unbiased trajectories. However, McCarty and
 86 Parrinello [29] using a linear approach and later Bonati et al. [14] using a non-linear one (Deep-TICA)
 87 have shown how to extract useful CVs from biased simulations. Here, we are going to apply Deep-TICA
 88 to the set of N_d descriptors on a converged OPES trajectory where the Deep-LDA CV was biased.

Enhanced sampling technique for estimating residence times: OPES flooding

As discussed in the introduction, to calculate rates from a biased simulation no bias must be deposited in the transition region. In such a case, the physical residence time t is related to the physical simulation time t_{MD} by [16, 17]

$$t = \langle e^{\beta V(s)} \rangle_V t_{\text{MD}} \quad (3)$$

where the acceleration factor $\langle e^{\beta V(s)} \rangle_V$ is computed as an average along the simulation. To ensure that the condition under which Eq. 3 is satisfied, we introduce a variation of OPES that we call OPES_f that is inspired by the variational flooding method [19]. The condition for OPES_f to allow calculating physical rates are easily satisfied. The maximum amount of bias deposited can be controlled by the choice of the ϵ parameter in Eq. 2, making sure that it is lower than the free energy barrier. Furthermore via the parameter EXCLUDED_REGION, one can prevent OPES_f from depositing bias in a preassigned region of configuration space. Since, from an initial free energy surface (FES) estimate one can roughly estimate both the height and the location of the transition, the setup of a subsequent OPES_f rate calculation is relatively simple and straightforward.

Designing water CVs

In our OPES simulations we shall use as CVs the distance z from the binding site and a CV that encodes water behavior. Water is known to play a non-trivial role in ligand binding [30, 31, 32, 33, 34, 35, 36, 37, 38], therefore we want to use the power of machine learning techniques to capture and encode its behavior in a CV, generalizing the strategy that we have proposed for a smaller host-guest system in Ref. 15. The choice of an effective set of descriptors \mathbf{d} is critical as it must capture the solvation in different situations that are relevant to the binding-unbinding process, i.e. the ligand itself, the binding position, the binding pocket and the binding path.

Some of these descriptors can be identified following physical intuition. For instance, to describe ligand’s solvation we focus on its charged tail and use the coordination number of water around the Carbon atom of the amidine group ($\{\text{G}\}$) (see Fig. 2c). Regarding the binding position, we analogously evaluate the coordination number between water and the Carbon atom of the carboxylate group in residue Asp189 ($\{\text{H}\}$). These choices are similar in spirit to other solvation variables that have been devised in the past [39, 40, 41, 42, 15].

However, we follow a different approach to describe water behavior around the binding pocket and along the binding path. The S1 binding pocket of apo Trypsin form (see Fig. 1) is known from high resolution experiments [7, 43] to be characterized by the presence of a set of deeply buried water molecules. These water molecules have a long residence time and form what has been called in Ref. 7 the reservoir. To encode their role in the descriptor set, we propose a general method that can be applied to describe trapped water molecules in biological systems. This method consists of four steps: 1) selecting the α -Carbon atoms of the relevant part of the protein that encloses water (label 1 of Fig. 2a), 2) building the convex hull surface with the coordinates from step 1 (label 2 of Fig. 2a), 3) collecting the position of the water molecules that lie within the convex hull for longer than a pre-defined lifetime, and 4) clustering the collected data to determine areas of high water density and their centers (label 4 of Fig. 2a). We call those points hydration spots ($\{\text{V}_i\}$) and use them as centers around which we calculate water coordination. The method is implemented in a Python script [44].

In Trypsin, we restrict the analysis to the region around the S1 binding pose where we build the convex hull. That area encloses both the deeply buried water molecules and the binding path of the ligand (see Fig. 2b). From a number of unbiased trajectories (both in states B and U) we collect the positions of the water molecules that have a residence time inside the convex hull longer than 100 ps. Then, using the clustering method introduced in step 4, we identify 16 Vs (see Fig. 2e) position. The number of clustering centers is chosen so that the relative distance between the Vs lies in a range of 2-3 Å. We find that there is a correspondence between the $\text{V}_5 - \text{V}_{12}$ centers and the position of the reservoir water molecules reported in Ref. 7. Furthermore, V_3 lies at the position of the long-lived water molecule that stabilizes binding

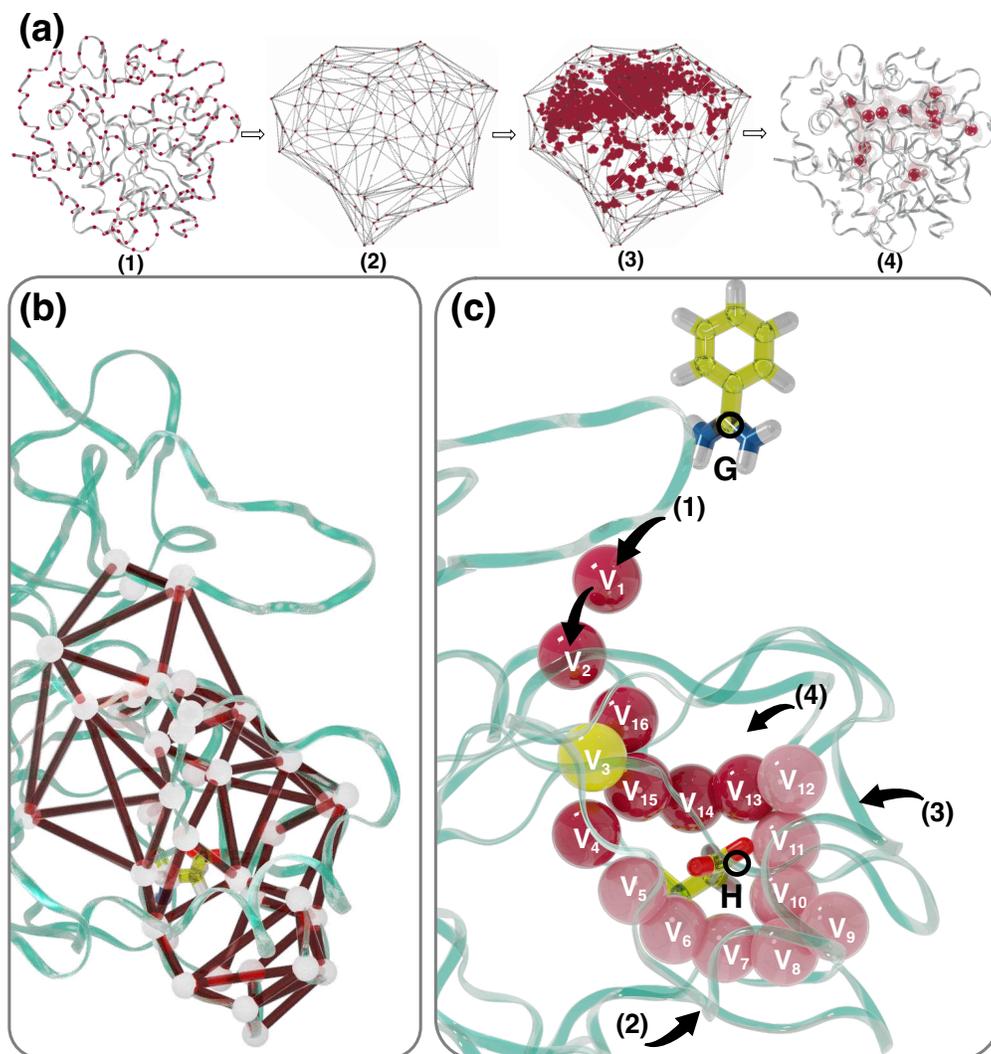


Figure 2: **Identification of long-lived water molecules.** (a) Graphical representation of the four steps strategy used to identify the long-lived hydration spots. (1) α -Carbon atoms of the residues located on a protein's outer surface, represented by red dots. (2) Convex hull surface built from the α -Carbon atoms shown in (1). (3) Distribution of the long-lived water molecules inside the surface. (4) Centers of the water distribution from step (3) obtained with a clustering algorithm. (b) Convex hull surface built around the binding pose of Trypsin. (c) Position of the 16 hydration spots V_i , shown by spheres. The hydration spots $V_5 - V_{12}$ coincide with the position of the reservoir water molecules of Ref. 7 and the dark red spheres lie on the binding path. The yellow sphere on V_3 shows the position of the key W_1 water molecule [7] that forms a hydrogen bond with the ligand. Black arrows indicate the four possible paths of entrance/exit for water molecules in the binding pocket.

136 and that in Ref. 7 is called W_1 . We use as descriptors all the water coordination on ($\{G\}$, $\{H\}$, $\{V_i\}$) to
 137 generate Deep-LDA s_w and Deep-TICA s_t water CVs.

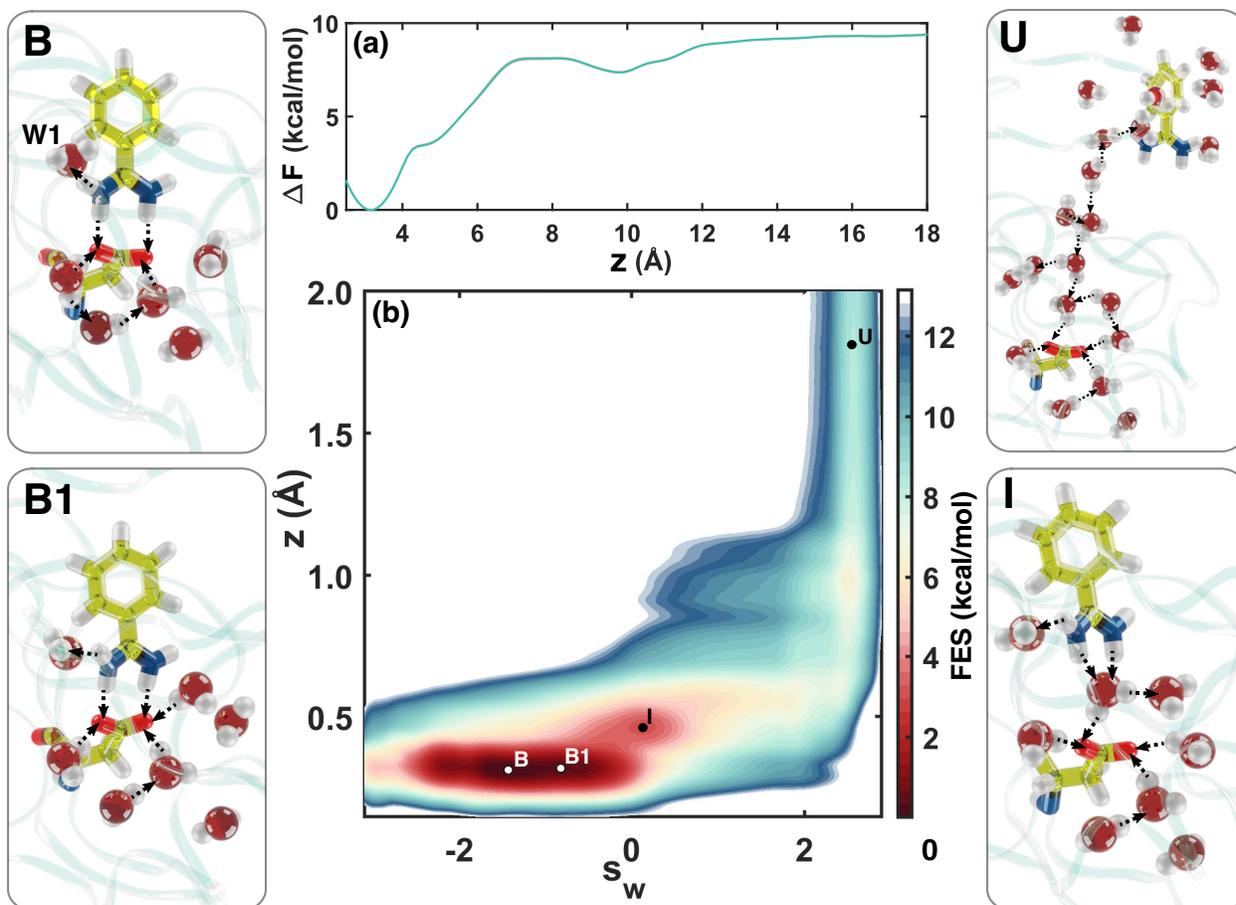


Figure 3: **Mechanistic interpretation of the binding process.** (a) 1D free energy surface of Trypsin-Benzamidine reconstructed using reweighting along z with the statistical uncertainty within the line width. (b) 2D free energy landscape along z and s_w . In the panels, representative configurations of the relevant states are shown with a focus on the solvation pattern around the ligand and the binding pose.

Table 1: Our simulation estimates of the free energy, enthalpy and entropy of binding at $T=300$ K and corresponding experimental data [45, 46].

	ΔF	ΔU	$-T\Delta S$
This work	6.36 ± 0.07	4.12 ± 1.56	2.23 ± 1.56
EXP.	6.36	4.52	1.84

138 **Static properties: binding free energy, enthalpy and entropy**

139 To calculate the binding free energy, we perform an OPES simulation using 32 walkers, biasing z and
 140 the Deep-LDA water CV s_w , for a total simulation time of $3.2 \mu s$. More details about the simulation are
 141 provided in the Supporting Information (SI). Convergence is achieved as several binding and unbinding
 142 events occur in a quasi-static regime of the bias (see Fig. S-1). In Fig. 3a, we show the FES projection on z ,
 143 calculated as a block average with an error bar that is within the line width of the plot. We calculate the
 144 free energy difference between the B and U states using the funnel correction in Eq. 4 obtaining a value of
 145 6.36 ± 0.07 kcal/mol, in an embarrassing and certainly fortuitous agreement with the experimental value

146 of 6.36 kcal/mol [45, 46].

147 The high level of accuracy that the combination of OPES and good quality CVs, makes it possible to
148 estimate separately enthalpy ΔU and entropy ΔS of binding. This is achieved by converging binding free
149 energy calculations in a range of temperatures and using the relationship $\Delta F = \Delta U - T\Delta S$. In Tab. 1 we
150 report our estimate for these thermodynamic quantities and observe that they are also in good agreement
151 with experiments. Further details about these simulations are provided in the SI.

152 In Fig. 3b, we show the two-dimensional FES of binding projected on z and s_w , along with a number of
153 representative snapshots of the different states and their typical water arrangement. State B is the global
154 minimum of the FES and corresponds to the protein-ligand crystallographic structure (e.g. PDB 3at1 [47]).
155 In line with experiments [7], the W_1 water molecule connects the ligand to Trypsin residues Tyr228, and
156 Ser190, forming a total of 3 HBs. Furthermore, in the reservoir region below residue Asp189 we observe
157 on average the presence of 5 water molecules, in agreement with the X-ray structure [7]. State B₁ is a
158 less stable binding pose where the ligand is in the same configuration as B, but the reservoir contains on
159 average one more water molecule. In state I, the reservoir region contains another extra water molecule
160 that tends to bridge the amidine group of the ligand and the binding site of Asp189, thus weakening
161 binding. In the fully dissociated state U, the binding cavity returns to its apo form in which the 5 water
162 of the reservoir are in the experimental structure of the holo form. The space previously occupied by the
163 ligand is replaced by about 4-5 water molecules. The fact that in the apo form the reservoir structure of
164 the water is preserved underlines once more the relevance of this structure (see Fig. 3).

165 Note that, as the ligand progresses towards the U state, the number of water molecules in the binding
166 site increases, underlining the role of water throughout the whole unbinding process. In spite of the limit
167 of being trained on data coming exclusively from B and U, the Deep-LDA CV in combination with OPES
168 is able to capture this non-trivial water behavior and converge the FES by virtue of the ability of OPES to
169 deal with non-optimal CVs [48]. Nevertheless, s_w is not able to resolve the diverse water arrangements in
170 states such as B and B₁, as we shall see below. In order to capture the finer details of the role of water in
171 the binding pose, the intermediate states and the path to unbinding, we use Deep-TICA to determine a
172 new water CV that we call s_t . Since Deep-TICA is trained on trajectories coming from biased simulations,
173 the resulting CV is informed on the water modes that Deep-LDA is not able to resolve.

174 Using the data generated in the Deep-LDA-driven simulation, we project in Fig. 4 (a,b) the FES along
175 different pairs of CVs: z, s_t and s_w, s_t , respectively. We find that s_t resolves the original B state into states
176 B and B', with state B being about 12 kJ/mol more stable than B'. From these two states, two different
177 unbinding pathways depart. We denote the states belonging to the less stable branch with a prime. Fig. 4
178 (d,e) shows typical configurations of the ligand and the surrounding water molecules in state B and B'.
179 The number of water molecules in the reservoir is the same, but the water arrangement is different. In B,
180 water molecules are part of an extended HB network that includes residue Asp189, while in state B' this
181 network is less structured.

182 One striking difference between B and B' is that in state B there is one water molecule trapped deeply
183 in the binding cavity, in an intermediate position between residues Tyr182, Lys219 and Gln219 which
184 corresponds to descriptor V₉ (see the semi transparent sphere in Fig. 4d). The FES projection along V₉
185 and s_t in Fig. 4c, confirms that V₉ discriminates well between B and B'. Further analysis reveals that state B
186 presents a slightly larger volume of the reservoir region than B' (see SI), which, in turn, possibly facilitates
187 the formation of an extended HB network. A study of the relative importance of the water descriptors on
188 the NN CVs can be also found in the SI.

189 The presence of a water molecule in V₉ discriminates well between B and B', as can be seen if we
190 project the FES along V₉ and s_t . Thus this molecule stabilizes the cavity HB network that is such an
191 important feature of this protein.

192 **Dynamic properties: ligand residence time**

193 We compute the ligand-unbinding rate by employing the OPES flooding technique on simulations that
194 start from state B. In this approach a choice of CVs that captures well the complexity of the path from
195 state B to state U is essential. The Deep-LDA coordinate built only on the knowledge of the B and U states

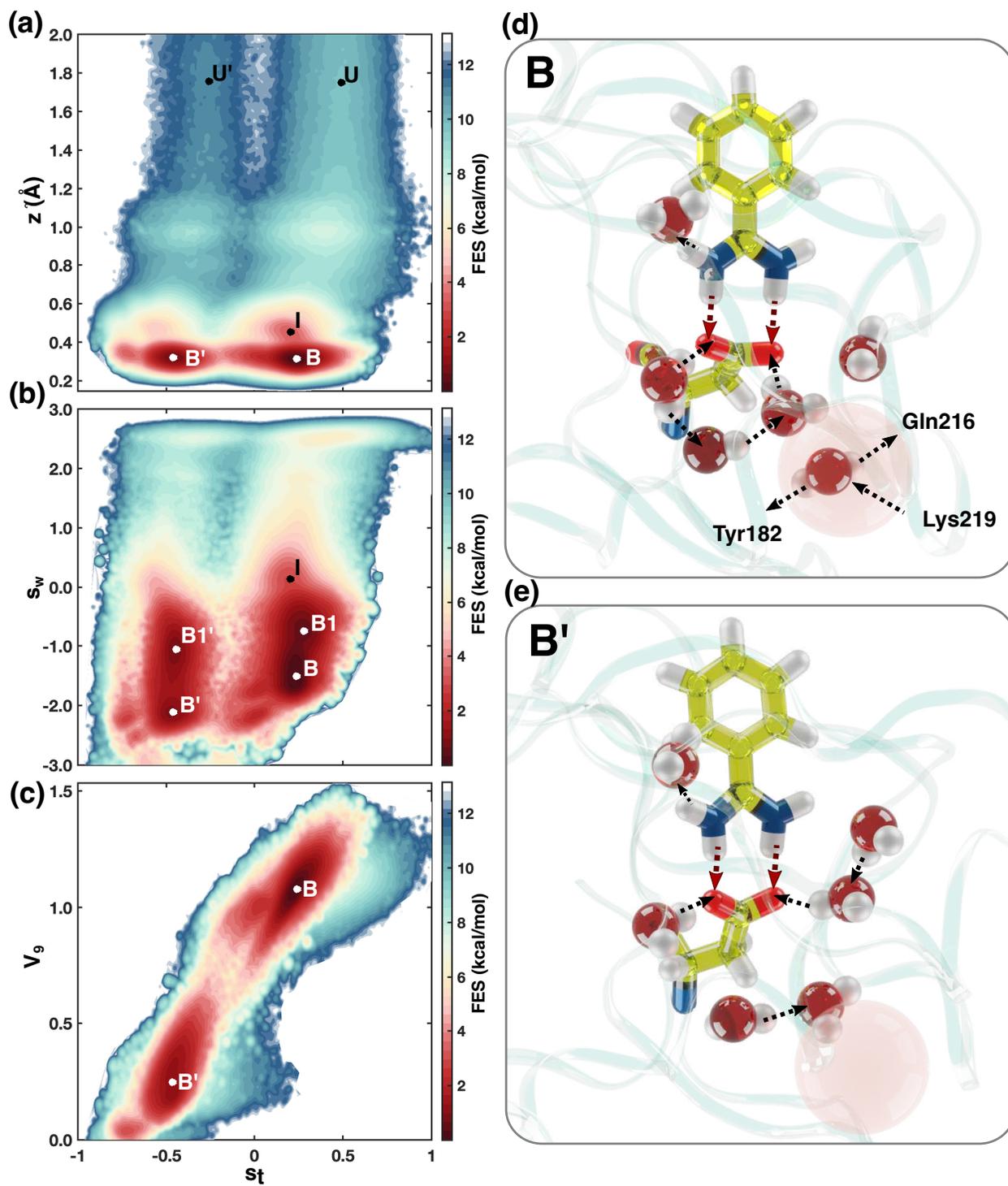


Figure 4: 2D FES along s_t and (a) z , (b) s_w , (c) V_9 . In (d) and (e), representative configurations of the ligand and its solvation pattern in states B and B', respectively. The location of V_9 is highlighted with a semi-transparent pink sphere.

196 is not adequate to the purpose. In fact, being able to fill all the different metastable bound states is crucial
 197 for promoting transitions to the U state without remaining stuck in intermediate states. For this reason
 198 besides z we use Deep-TICA CV s_t . By construction, s_t is able to extract the slow modes of the system
 199 and as a consequence to drive it out of deep basins towards the transition state. The use of OPES_f for the
 200 calculation of rates requires that we define an excluded region to avoid depositing bias in the transition
 201 state region. An analysis of the FES in Fig. 4a suggests to prevent depositing bias in the region $z > 6$ Å.
 202 We run a total of 55 ligand unbinding simulations and, by using Eq. 3, we determine for each simulation
 203 a physical ligand residence time t .

204 The distribution of transition times of a rare event dominated by a single barrier is expected to be
 205 Poissonian $\frac{1}{\tau}e^{-t/\tau}$ where τ is the characteristic time of the associated homogeneous process [49]. We fit
 206 all our data to such a model and find that the quality of the fit is poor (see Fig. S-7), as indicated by its low
 207 p-value. It is known that complex biological processes present multiple timescales and are not expected to
 208 necessarily follow such a simple model [50, 1, 51]. This led us to further analyze the unbinding trajectories
 209 and to identify the presence of two possible unbinding mechanisms: a faster and a slower one. The key
 210 difference between the two is related in the water network arrangement in the binding pocket (see Fig. 5).

211 All unbinding events start with the water reservoir acquiring an extra molecule and the system reach-
 212 ing state B1 (see Fig. 5b). Then two possibilities occur. In the faster case, the presence of a water molecule
 213 in the vicinity of W1 weakens the bond between the ligand and W1 (see Fig. 5c), which leads to another
 214 water molecule to bind to residue Asp189 (see Fig. 5d) and finally brings about the ligand-unbinding
 215 event. In the slower mechanism, W1 is stable and, as the reservoir increases its water content (see Fig. 5e),
 216 eventually one water molecule bridges the ligand and Asp189 (see Fig. 5f), driving the system towards
 217 unbinding. As reported in Tab. 2, both mechanisms fit well the homogeneous model and present a discrep-
 218 ancancy in the resulting τ of about one order of magnitude. In particular, the slower mechanism τ that will
 219 dominate the residence time has a $k_{\text{off}} = 1/\tau = 687 \text{ s}^{-1}$ that is in close agreement with the experimental
 220 value of $k_{\text{off}} = 600 \pm 300 \text{ s}^{-1}$ [52].

Table 2: **Ligand residence time.** τ is the characteristic time from an exponential fit [49], $k_{\text{off}} = 1/\tau$ and
 p-value measures the quality of the fit. μ and σ are the average and the standard deviation of the data,
 respectively. We show results from all our data and from data split into the two observed mechanisms.
 The experimental results are taken from Ref. 52.

	τ (10^{-3} s)	k_{off} (10^2 s^{-1})	p-value	$\mu \pm \sigma$ (10^{-3} s)
All data	0.64	15.6	0.06	1.10 ± 1.49
Faster path	0.18	55.4	0.62	0.23 ± 0.29
Slower path	1.45	6.87	0.63	1.58 ± 1.59
EXP.	-	6.00 ± 3.00	-	-

221 Discussion

222 Our work shows clearly that water plays a major role in the activation of Trypsin and modulates the release
 223 process of Benzamidine. Through the use of water-focused machine-learned CVs, we estimate static and
 224 dynamic properties in excellent agreement with experiments and we analyze the results with a level of
 225 resolution that lets us identify the importance of each individual water molecules. We observe how the
 226 presence of water in key positions affects the binding stability and determines unbinding mechanisms
 227 with significantly different ligand residence times. In conclusion, one might say that the set of reservoir
 228 water molecules, present both in the apo and the holo form, are an essential part of the enzyme, its
 229 structure and its activity. We believe that this is not just specific to the Trypsin-Benzamidine system and
 230 we argue that water would play such a non-trivial role at least in all the cases where the ligand has to
 231 negotiate its way out of a water-rich enzymatic cavity.

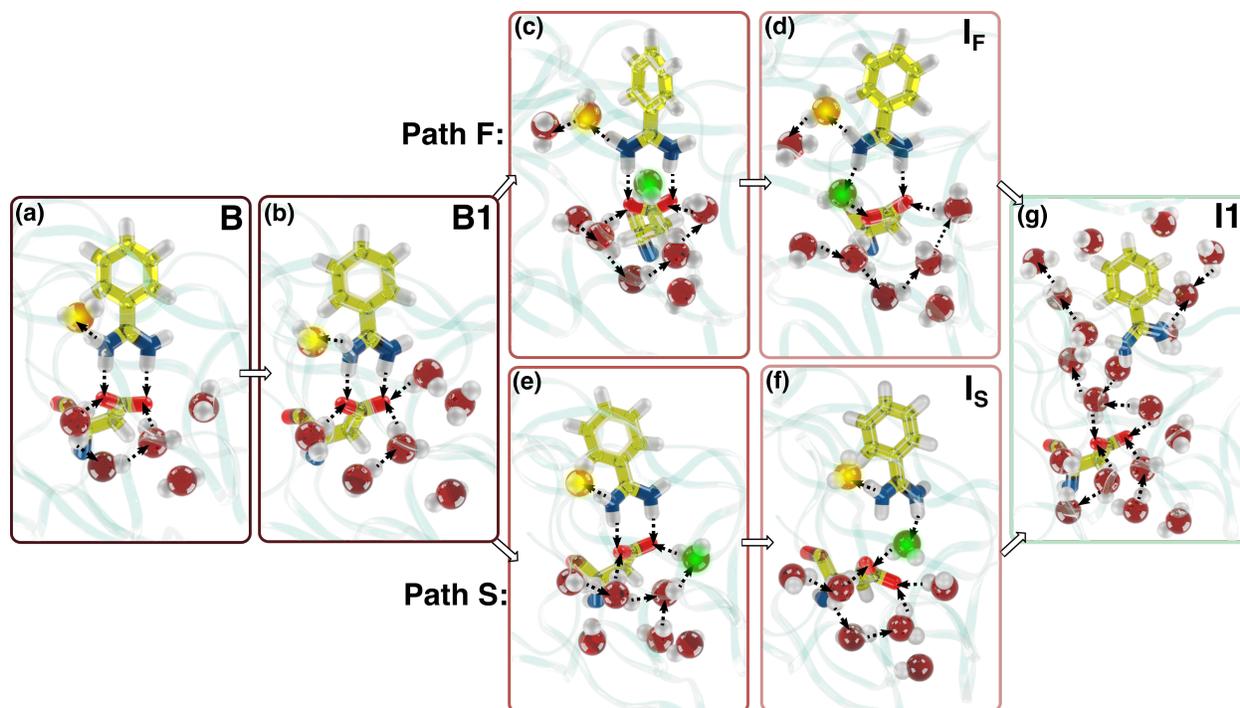


Figure 5: **Two ligand unbinding mechanisms.** Representative configurations of states (a) B and (b) B1. (c) and (e) the pre-intermediate step of fast and slow mechanisms, respectively. (d) and (f) the intermediate step of the faster (I_F) and the slower (I_S) unbinding mechanisms. (g) Intermediate I_1 that presents a number of water molecules between ligand and the binding pose. The yellow sphere highlights the location of the W_1 water molecule. The green sphere indicates the position of the water molecule that bridges the ligand and the Asp189 residue. The protein is shown as a semi-transparent ribbon.

232 Methods

233 We perform all simulations with the molecular dynamics code GROMACS 2020.4 [53] patched with
 234 PLUMED 2.8 [54] and the Pytorch library 1.4 [55, 13]. We use TIP3P water, while the protein is described
 235 by the Amber-14SB force field and the ligand interactions are taken from the Amber GAFF library [56]. We
 236 employ Ewald summation [57] for long-range electrostatic interactions with a cutoff of 10 Å for both the
 237 Coulomb and the van der Waals interactions. All simulations are run in the NPT ensemble with a timestep
 238 of 2 fs. We use the Parrinello-Rahman barostat [58] and the stochastic velocity rescaling thermostat [59],
 239 both with a coupling constant of 1 ps.

240 We train a Deep-LDA CV by running unbiased simulations on states B and U for 60 ns and evaluating
 241 the water coordination on a descriptor set \mathbf{d} made by the 18 components ($\{G\}$, $\{H\}$, $\{V_i\}$). For training,
 242 we take as state B the initial configuration used in Ref. 60. We use a NN made of 4 layers with 2.5×10^{-5}
 243 learning rate. To train the Deep-TICA CV, we take the converged OPES trajectory used for calculating the
 244 binding free energy and feed the descriptors set \mathbf{d} to a NN made of 4 layers with 1.0×10^{-4} learning rate
 245 and 0.07 lag time. Further details can be found in the SI.

246 One of the difficulties in ligand-unbinding problem arises from the fact that once the ligand leaves the
 247 protein, it has to explore a large conformational space. To tackle this issue, we use the Funnel restraint
 248 proposed in Ref. 61. In this method, the space available to the ligand in the unbound state is limited
 249 by confining it to a cylindrical volume above the binding site (see Fig. 1) which introduces an entropic
 250 restraint in the U state. To calculate the absolute free energy of binding, one needs to apply a correction

251 to the apparent free energy $W(z)$ that comes from the simulation:

$$\Delta F = -\frac{1}{\beta} \log \left(C^0 \pi R_{\text{cyl}}^2 \int_{\text{B}} dz \exp(-\beta(W(z) - W_{\text{U}})) \right) \quad (4)$$

252 where $C^0 = 1/1660 \text{ \AA}^{-3}$ is the standard concentration, z is the distance between the center of mass of the
253 ligand and the binding site along the funnel's axis, $W(z)$ is the free energy along the funnel axis and W_{U}
254 is its reference value in state U.

255 Data Availability

256 All the inputs and instructions to reproduce the results presented in this manuscript are soon to be
257 uploaded to the PLUMED-NEST repository. A tutorial on Deep-LDA training can be found at this [link](#)
258 while a tutorial on Deep-TICA is at this [link](#).

259 References

- 260 [1] Pan, A. C., Borhani, D. W., Dror, R. O. & Shaw, D. E. Molecular determinants of drug-receptor
261 binding kinetics. *Drug Discovery Today* **18**, 667–673 (2013). URL [http://dx.doi.org/10.1016/j.](http://dx.doi.org/10.1016/j.drudis.2013.02.007)
262 [drudis.2013.02.007](http://dx.doi.org/10.1016/j.drudis.2013.02.007)<https://linkinghub.elsevier.com/retrieve/pii/S1359644613000627>.
- 263 [2] Gervasio, F. L., Laio, A. & Parrinello, M. Flexible Docking in Solution Using Metadynamics. *Journal*
264 *of the American Chemical Society* **127**, 2600–2607 (2005). URL [https://pubs.acs.org/doi/10.1021/](https://pubs.acs.org/doi/10.1021/ja0445950)
265 [ja0445950](https://pubs.acs.org/doi/10.1021/ja0445950).
- 266 [3] Tiwary, P., Limongelli, V., Salvalaglio, M. & Parrinello, M. Kinetics of protein–ligand unbinding:
267 Predicting pathways, rates, and rate-limiting steps. *Proceedings of the National Academy of Sciences* **112**
268 (2015).
- 269 [4] Wolf, S., Lickert, B., Bray, S. & Stock, G. Multisecond ligand dissociation dynamics from atomistic
270 simulations. *Nature Communications* **11**, 2918 (2020).
- 271 [5] Brotzakis, Z. F., Limongelli, V. & Parrinello, M. Accelerating the calculation of protein–ligand binding
272 free energy and residence times using dynamically optimized collective variables. *Journal of Chemical*
273 *Theory and Computation* **15**, 743–750 (2019).
- 274 [6] Ray, D., Stone, S. E. & Andricioaei, I. Markovian weighted ensemble milestoning (m-wem): Long-time
275 kinetics from short trajectories. *Journal of Chemical Theory and Computation* **18**, 79–95 (2022).
- 276 [7] Schiebel, J. *et al.* Intriguing role of water in protein-ligand binding studied by neutron crystallography
277 on trypsin complexes. *Nature Communications* **9**, 3559 (2018). URL [http://dx.doi.org/10.1038/](http://dx.doi.org/10.1038/s41467-018-05769-2)
278 [s41467-018-05769-2](http://dx.doi.org/10.1038/s41467-018-05769-2)<http://www.nature.com/articles/s41467-018-05769-2>.
- 279 [8] Invernizzi, M. & Parrinello, M. Rethinking Metadynamics: From Bias Potentials to Probability Distri-
280 butions. *The Journal of Physical Chemistry Letters* **11**, 2731–2736 (2020). URL [https://link.aps.org/](https://link.aps.org/doi/10.1103/PhysRevLett.115.070601)
281 [doi/10.1103/PhysRevLett.115.070601](https://link.aps.org/doi/10.1103/PhysRevLett.115.070601).
- 282 [9] Laio, A. & Parrinello, M. Escaping free-energy minima. *Proceedings of the National Academy of Sciences*
283 **99**, 12562–12566 (2002). URL [http://www.pnas.org/cgi/doi/10.1073/pnas.](http://www.pnas.org/cgi/doi/10.1073/pnas.202427399)
[202427399](http://www.pnas.org/cgi/doi/10.1073/pnas.202427399).
- 284 [10] Barducci, A., Bussi, G. & Parrinello, M. Well-Tempered Metadynamics: A Smoothly Converging and
285 Tunable Free-Energy Method. *Physical Review Letters* **100**, 020603 (2008). URL [https://link.aps.](https://link.aps.org/doi/10.1103/PhysRevLett.100.020603)
286 [org/doi/10.1103/PhysRevLett.100.020603](https://link.aps.org/doi/10.1103/PhysRevLett.100.020603).
- 287 [11] Torrie, G. & Valleau, J. Nonphysical sampling distributions in monte carlo free-energy estimation:
288 Umbrella sampling. *Journal of Computational Physics* **23**, 187–199 (1977).
- 289 [12] Valsson, O., Tiwary, P. & Parrinello, M. Enhancing Important Fluctuations: Rare Events and Metady-
290 namics from a Conceptual Viewpoint. *Annual Review of Physical Chemistry* **67**, 159–184 (2016). URL
291 <http://www.annualreviews.org/doi/10.1146/annurev-physchem-040215-112229>.
- 292 [13] Bonati, L., Rizzi, V. & Parrinello, M. Data-Driven Collective Variables for Enhanced Sampling.
293 *The Journal of Physical Chemistry Letters* **11**, 2998–3004 (2020). URL [http://arxiv.org/abs/2002.](http://arxiv.org/abs/2002.06562)
294 [06562](http://arxiv.org/abs/2002.06562)<https://pubs.acs.org/doi/10.1021/acs.jpcllett.0c00535>.

- 295 [14] Bonati, L., Piccini, G. & Parrinello, M. Deep learning the slow modes for rare events sampling. *Proceedings of the National Academy of Sciences* **118** (2021). URL <http://arxiv.org/abs/2107.03943><http://www.pnas.org/lookup/doi/10.1073/pnas.2113533118><https://pnas.org/doi/full/10.1073/pnas.2113533118>.
- 299 [15] Rizzi, V., Bonati, L., Ansari, N. & Parrinello, M. The role of water in host-guest interaction. *Nature Communications* **12**, 93 (2021). URL <http://dx.doi.org/10.1038/s41467-020-20310-0><http://www.nature.com/articles/s41467-020-20310-0>.
- 302 [16] Grubmüller, H. Predicting slow structural transitions in macromolecular systems: Conformational flooding. *Physical Review E* **52**, 2893–2906 (1995). URL <https://link.aps.org/doi/10.1103/PhysRevE.52.2893>.
- 305 [17] Voter, A. F. A method for accelerating the molecular dynamics simulation of infrequent events. *The Journal of Chemical Physics* **106**, 4665–4677 (1997). URL <http://aip.scitation.org/doi/10.1063/1.473503>.
- 308 [18] Tiwary, P. & Parrinello, M. From Metadynamics to Dynamics. *Physical Review Letters* **111**, 230602 (2013). URL <https://link.aps.org/doi/10.1103/PhysRevLett.111.230602>.
- 310 [19] McCarty, J., Valsson, O., Tiwary, P. & Parrinello, M. Variationally Optimized Free-Energy Flooding for Rate Calculation. *Physical Review Letters* **115**, 070601 (2015). URL <https://link.aps.org/doi/10.1103/PhysRevLett.115.070601>.
- 313 [20] Debnath, J. & Parrinello, M. Gaussian Mixture-Based Enhanced Sampling for Statics and Dynamics. *The Journal of Physical Chemistry Letters* **11**, 5076–5080 (2020). URL <https://pubs.acs.org/doi/10.1021/acs.jpcllett.0c01125>.
- 316 [21] Welling, M. Fisher Linear Discriminant Analysis. Tech. Rep., Dep. Comput. Sci. Univ. Toronto (2005).
- 317 [22] Mendels, D., Piccini, G. & Parrinello, M. Collective Variables from Local Fluctuations. *The Journal of Physical Chemistry Letters* **9**, 2776–2781 (2018). URL <http://pubs.acs.org/doi/10.1021/acs.jpcllett.8b00733>.
- 320 [23] Karmakar, T., Invernizzi, M., Rizzi, V. & Parrinello, M. Collective variables for the study of crystallisation. *Molecular Physics* **119**, e1893848 (2021). URL <http://arxiv.org/abs/2101.03150><https://www.tandfonline.com/doi/full/10.1080/00268976.2021.1893848>.
- 323 [24] Ansari, N., Rizzi, V., Carloni, P. & Parrinello, M. Water-Triggered, Irreversible Conformational Change of SARS-CoV-2 Main Protease on Passing from the Solid State to Aqueous Solution. *Journal of the American Chemical Society* **143**, 12930–12934 (2021). URL <http://biorxiv.org/content/early/2021/05/21/2021.05.21.445090.abstract><https://pubs.acs.org/doi/10.1021/jacs.1c05301>.
- 327 [25] Bjelobrk, Z. *et al.* Naphthalene crystal shape prediction from molecular dynamics simulations. *CryStEngComm* **21**, 3280–3288 (2019).
- 329 [26] Naritomi, Y. & Fuchigami, S. Slow dynamics in protein fluctuations revealed by time-structure based independent component analysis: The case of domain motions. *The Journal of Chemical Physics* **134**, 065101 (2011). URL <http://aip.scitation.org/doi/10.1063/1.3554380>.
- 332 [27] Pérez-Hernández, G., Paul, F., Giorgino, T., De Fabritiis, G. & Noé, F. Identification of slow molecular order parameters for Markov model construction. *The Journal of Chemical Physics* **139**, 015102 (2013). URL <http://aip.scitation.org/doi/10.1063/1.4811489>.
- 335 [28] Wu, H. *et al.* Variational Koopman models: Slow collective variables and molecular kinetics from short off-equilibrium simulations. *The Journal of Chemical Physics* **146**, 154104 (2017). URL <http://aip.scitation.org/doi/10.1063/1.4979344>.

- 338 [29] McCarty, J. & Parrinello, M. A variational conformational dynamics approach to the selection of
339 collective variables in metadynamics. *The Journal of Chemical Physics* **147**, 204109 (2017). URL <http://arxiv.org/abs/1703.08777><http://aip.scitation.org/doi/10.1063/1.4998598>.
340
- 341 [30] Ladbury, J. E. Just add water! The effect of water on the specificity of protein-ligand binding sites
342 and its potential application to drug design. *Chemistry & Biology* **3**, 973–980 (1996). URL <https://linkinghub.elsevier.com/retrieve/pii/S1074552196901647>.
343
- 344 [31] Homans, S. Water, water everywhere — except where it matters? *Drug Discovery Today* **12**, 534–539
345 (2007). URL <https://linkinghub.elsevier.com/retrieve/pii/S135964460700222X>.
- 346 [32] Ewell, J., Gibb, B. C. & Rick, S. W. Water inside a hydrophobic cavitand molecule. *Journal of Physical*
347 *Chemistry B* **112**, 10272–10279 (2008).
- 348 [33] Mahmoud, A. H., Masters, M. R., Yang, Y. & Lill, M. A. Elucidating the multiple roles of hydration
349 for accurate protein-ligand binding prediction via deep learning. *Communications Chemistry* **3**, 19
350 (2020). URL <http://dx.doi.org/10.1038/s42004-020-0261-x><http://www.nature.com/articles/s42004-020-0261-x>.
351
- 352 [34] Hummer, G. Under water’s influence. *Nature Chemistry* **2**, 906–907 (2010).
- 353 [35] Mitusińska, K., Raczyńska, A., Bzówka, M., Bagrowska, W. & Góra, A. Applications of water
354 molecules for analysis of macromolecule properties. *Computational and Structural Biotechnology Journal*
355 **18**, 355–365 (2020). URL <https://linkinghub.elsevier.com/retrieve/pii/S2001037019304556>.
- 356 [36] Maurer, M. & Oostenbrink, C. Water in protein hydration and ligand recognition. *Journal of Molecular*
357 *Recognition* **32** (2019).
- 358 [37] Bodnarchuk, M. S. Water, water, everywhere... it’s time to stop and think. *Drug Discovery Today* **21**,
359 1139–1146 (2016).
- 360 [38] Bellissent-Funel, M.-C. *et al.* Water determines the structure and dynamics of proteins. *Chemical*
361 *Reviews* **116**, 7673–7697 (2016).
- 362 [39] Pietrucci, F., Marinelli, F., Carloni, P. & Laio, A. Substrate Binding Mechanism of HIV-1 Protease
363 from Explicit-Solvent Atomistic Simulations. *Journal of the American Chemical Society* **131**, 11811–11818
364 (2009). URL <https://pubs.acs.org/doi/10.1021/ja903045y>.
- 365 [40] Tiwary, P., Mondal, J. & Berne, B. J. How and when does an anticancer drug leave its binding site?
366 *Science Advances* **3**, e1700014 (2017). URL <https://advances.sciencemag.org/lookup/doi/10.1126/sciadv.1700014>.
367
- 368 [41] Pérez-Conesa, S., Piaggi, P. M. & Parrinello, M. A local fingerprint for hydrophobicity and hy-
369 drophilicity: From methane to peptides. *The Journal of Chemical Physics* **150**, 204103 (2019). URL
370 <http://dx.doi.org/10.1063/1.5088418><http://aip.scitation.org/doi/10.1063/1.5088418>.
- 371 [42] Brotzakis, Z. F., Limongelli, V. & Parrinello, M. Accelerating the Calculation of Protein–Ligand Bind-
372 ing Free Energy and Residence Times Using Dynamically Optimized Collective Variables. *Journal of*
373 *Chemical Theory and Computation* **15**, 743–750 (2019). URL <https://pubs.acs.org/doi/10.1021/acs.jctc.8b00934>.
374
- 375 [43] Hüfner-Wulsdorf, T. & Klebe, G. Role of Water Molecules in Protein–Ligand Dissociation and Se-
376 lectivity Discrimination: Analysis of the Mechanisms and Kinetics of Biomolecular Solvation Us-
377 ing Molecular Dynamics. *Journal of Chemical Information and Modeling* **60**, 1818–1832 (2020). URL
378 <https://pubs.acs.org/doi/abs/10.1021/acs.jcim.0c00156>.

- 379 [44] Ansari, N. & Raucci, U. Hydration spot. https://github.com/narjesansari/hydration_spot
380 (2022).
- 381 [45] Talhout, R., Villa, A., Mark, A. E. & Engberts, J. B. F. N. Understanding binding affinity: A combined
382 isothermal titration calorimetry/molecular dynamics study of the binding of a series of hydrophobi-
383 cally modified benzamidine chloride inhibitors to trypsin. *Journal of the American Chemical Society*
384 **125**, 10570–10579 (2003).
- 385 [46] Talhout, R. & Engberts, J. B. F. N. Probing the effect of the amidinium group and the phenyl ring
386 on the thermodynamics of binding of benzamidine chloride to trypsin. *Organic and Biomolecular*
387 *Chemistry* **2**, 3071 (2004).
- 388 [47] Yamane, J. *et al.* In-crystal affinity ranking of fragment hit compounds reveals a relationship with
389 their inhibitory activities. *Journal of Applied Crystallography* **44**, 798–804 (2011).
- 390 [48] Invernizzi, M. & Parrinello, M. Exploration vs Convergence Speed in Adaptive-bias Enhanced Sam-
391 pling. *arXiv* (2022). URL <http://arxiv.org/abs/2201.09950>.
- 392 [49] Salvalaglio, M., Tiwary, P. & Parrinello, M. Assessing the Reliability of the Dynamics Reconstructed
393 from Metadynamics. *Journal of Chemical Theory and Computation* **10**, 1420–1425 (2014). URL <http://pubs.acs.org/doi/10.1021/ct500040r>.
- 395 [50] Lindorff-Larsen, K., Piana, S., Dror, R. O. & Shaw, D. E. How Fast-Folding Proteins Fold. *Science* **334**,
396 517–520 (2011). URL <http://www.ncbi.nlm.nih.gov/pubmed/22034434><https://www.sciencemag.org/lookup/doi/10.1126/science.1208351>.
- 398 [51] Peña Ccoa, W. J. & Hocky, G. M. Assessing models of force-dependent unbinding rates via infrequent
399 metadynamics. *The Journal of Chemical Physics* **156**, 125102 (2022). URL <http://arxiv.org/abs/2112.03249><https://aip.scitation.org/doi/10.1063/5.0081078>.
- 401 [52] Guillain, F. & Thusius, D. Use of proflavine as an indicator in temperature-jump studies of the
402 binding of a competitive inhibitor to trypsin. *Journal of the American Chemical Society* **92**, 5534–5536
403 (1970).
- 404 [53] Abraham, M. J. *et al.* GROMACS: High performance molecular simulations through multi-level par-
405 allelism from laptops to supercomputers. *SoftwareX* **1-2**, 19–25 (2015). URL <https://linkinghub.elsevier.com/retrieve/pii/S2352711015000059>.
- 407 [54] Tribello, G. A., Bonomi, M., Branduardi, D., Camilloni, C. & Bussi, G. PLUMED 2: New feathers
408 for an old bird. *Computer Physics Communications* **185**, 604–613 (2014). URL <http://linkinghub.elsevier.com/retrieve/pii/S0010465513003196>.
- 410 [55] Paszke, A. *et al.* Automatic differentiation in PyTorch. *Adv. Neural Inf. Process. Syst.* **32**, 8024–8035
411 (2019).
- 412 [56] Maier, J. A. *et al.* ff14sb: Improving the accuracy of protein side chain and backbone parameters from
413 ff99sb. *Journal of Chemical Theory and Computation* **11**, 3696–3713 (2015).
- 414 [57] Essmann, U. *et al.* A smooth particle mesh ewald method. *The Journal of Chemical Physics* **103**, 8577–
415 8593 (1995). URL <http://aip.scitation.org/doi/10.1063/1.470117>.
- 416 [58] Parrinello, M. & Rahman, A. Polymorphic transitions in single crystals: A new molecular dynamics
417 method. *Journal of Applied Physics* **52**, 7182–7190 (1981).
- 418 [59] Bussi, G., Donadio, D. & Parrinello, M. Canonical sampling through velocity rescaling. *The Journal of*
419 *Chemical Physics* **126**, 014101 (2007).

- 420 [60] Brotzakis, Z. F. Vac metad applications. [https://gitlab.e-cam2020.eu/brotzakis/vac_metad_](https://gitlab.e-cam2020.eu/brotzakis/vac_metad_applications/-/tree/master/TRYPsin/GROMACS_INPUte)
421 [applications/-/tree/master/TRYPsin/GROMACS_INPUte](https://gitlab.e-cam2020.eu/brotzakis/vac_metad_applications/-/tree/master/TRYPsin/GROMACS_INPUte) (2019).
- 422 [61] Limongelli, V. *et al.* Sampling protein motion and solvent effect during ligand binding. *Proceedings of*
423 *the National Academy of Sciences of the United States of America* **109**, 1467–1472 (2012).

424 **Acknowledgements**

425 Part of this work was performed while the authors were associated with ETH Zürich and Università della
426 Svizzera italiana, Lugano.

427 The authors thank Faidon Brotzakis for providing the simulations initial input and they are grateful to
428 Michele Invernizzi, Luigi Bonati, Jayashrita Debnath and Umberto Raucci for many useful conversations.
429 The authors thank Michele Invernizzi for the implementation of OPES flooding. N.A. thanks Umberto
430 Raucci for the support in developing the convex hull script. M.P. thanks Maria Ramos and Paolo Carloni
431 for enlightening discussions on protein regulation.

432 **Author contributions statement**

433 N.A. performed the simulations. N.A., V.R. and M.P. discussed the results and reviewed the manuscript.

434 **Competing financial interests**

435 The authors declare no competing financial interests.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SI.pdf](#)