



An online repository of solvation thermodynamic and structural maps of SARS-CoV-2 targets

Brian Olson^{3,4} · Anthony Cruz^{1,2} · Lieyang Chen^{1,3} · Mossa Ghattas^{1,2} · Yeonji Ji^{1,3} · Kunhui Huang^{1,3} · Steven Ayoub Jr⁶ · Tyler Luchko⁷ · Daniel J. McKay⁵ · Tom Kurtzman^{1,2,3}

Received: 4 June 2020 / Accepted: 29 August 2020 / Published online: 12 September 2020
© Springer Nature Switzerland AG 2020

Abstract

SARS-CoV-2 recently jumped species and rapidly spread via human-to-human transmission to cause a global outbreak of COVID-19. The lack of effective vaccine combined with the severity of the disease necessitates attempts to develop small molecule drugs to combat the virus. COVID19_GIST_HSA is a freely available online repository to provide solvation thermodynamic maps of COVID-19-related protein small molecule drug targets. Grid inhomogeneous solvation theory maps were generated using AmberTools cpptraj-GIST, 3D reference interaction site model maps were created with AmberTools rism3d.snglpnt and hydration site analysis maps were created using SSTMap code. The resultant data can be applied to drug design efforts: scoring solvent displacement for docking, rational lead modification, prioritization of ligand- and protein- based pharmacophore elements, and creation of water-based pharmacophores. Herein, we demonstrate the use of the solvation thermodynamic mapping data. It is hoped that this freely provided data will aid in small molecule drug discovery efforts to defeat SARS-CoV-2.

Keywords COVID-19 · Solvation thermodynamics · Virtual screening · Rational lead modification · Drug discovery

Introduction

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) recently emerged and spread to cause a pandemic of coronavirus disease 2019 (COVID-19). Given the failure to contain the initial outbreak, the global failure to restrain the pandemic, and the absence of an effective vaccine, we may need to identify existing drugs or develop new drugs to interrupt COVID-19 at a critical juncture.

A number of targets may be of interest for the development of small molecule therapeutics for COVID-19: main protease (M^{pro} , 3CL^{pro}), helicase (Nsp13), endoribonuclease (Nsp15), and 2'-*O*-methyltransferase (Nsp10/16) are known viral protein drug targets for SARS-CoV-2. Small molecule drugs may target the substrate binding site of M^{pro} , the ADP binding site of Nsp13, the active site of Nsp15, or the *S*-adenosylmethionine (SAM) binding site of Nsp16.

Water is essential to the description of interactions between drugs and their biomolecular targets because solvation is a key contributor to molecular recognition and binding. Energies, entropies, and structural features of water molecules can be used to identify waters that may produce favorable or unfavorable contributions to the free energy of

✉ Tom Kurtzman
thomas.kurtzman@lehman.cuny.edu

¹ Lehman College Department of Chemistry, 205 W Bedford Park Blvd, Bronx, NY 10468, USA

² Ph.D. Program in Chemistry, The Graduate Center of the City University of New York, 365 5th Avenue, New York, NY 10016, USA

³ Ph.D. Program in Biochemistry, The Graduate Center of the City University of New York, 365 5th Avenue, New York, NY 10016, USA

⁴ Department of Biology and Chemistry, County College of Morris, 214 Center Grove Rd, Randolph, NJ 07869, USA

⁵ Ventus Therapeutics, Frederick-Banting, Montreal, QC H9S 2A1, Canada

⁶ Department of Chemistry and Biochemistry, California State University, Northridge, 18111 Nordhoff Street, Northridge, CA 91330, USA

⁷ Department of Physics and Astronomy, Center for Biological Physics, California State University, Northridge, 18111 Nordhoff Street, Northridge, CA 91330, USA

binding upon displacement and therefore aid in the identification of ligand interactions that may or may not be desirable. Water networks and tightly bound structural waters can affect ligand–receptor binding affinities. Information on water structure and thermodynamics may be useful to screen virtual compound databases, to identify new lead drug candidates, and inform rational lead modification to improve affinity and specificity for its target [1–4]. Ignoring water molecules in binding sites may reduce the chance that a drug design project will be successful.

Solvation thermodynamic mapping (STM) is widely used in academic studies of drug–protein interactions and has been widely integrated into the workflow of drug discovery and rational design efforts at major pharmaceutical companies. The utility of STM spans a number of areas in early-stage drug development efforts including virtual screening [1, 2], formation or improvement of pharmacophores [2, 5], docking [1, 6], and rational lead modification [3, 4]. While the utility of STM is apparent, there are significant obstacles to widespread use. Of particular concern is that many existing software packages for characterizing water properties are commercial and, hence, not available to all and/or they require computational expertise in molecular dynamics, computer modeling, and statistical mechanics in order to apply. This set of skills often does not exist in wet chemistry labs whose research is dedicated to discovering and optimizing new pharmaceutical compounds.

The goal of this publication is to remove these obstacles and make publicly available solvation thermodynamic and structural maps of SARS-CoV-2 targets as a resource to the academic and industrial drug design community to aid in their pursuit of identifying small molecule treatments for COVID-19. In order to aid in screening and modification of drugs, we offer a free public repository of solvation thermodynamic maps of significant small molecule COVID-19 drug targets. Here we present solvation maps of seven targets that are likely viable for small molecule modulation. All GIST maps, active site 3D-RISM maps, and simulation data are publically available on the KurtzmanLab github (github): github.com/KurtzmanLab/COVID19_GIST_HSA. Full system 3D-RISM maps can be accessed at <https://scholarworks.csun.edu/handle/10211.3/217209>.

Methods

Protein preparation

Protein monomer structures were prepared using the Protein Preparation Wizard [7] in Maestro [8] with default settings. ACE and NMA groups were used to cap the protein termini. Active sites were visually inspected and compared to ligand-bound structures to ensure that protonation states

and conformations were consistent with known ligand–protein interactions. Proteins were left as prepared by Maestro except for 6YB7 for which the protonation state of H163 was changed from being protonated in the delta position (HID) to the epsilon position (HIE) to produce an expected ligand protein interaction. No changes were made for other proteins. Energy minimization for hydrogen atoms was then performed in Maestro. The resulting protonation states can be found in Table 1.

A second set of structure models for SARS-CoV-2 M^{pro} (PDB IDs 6YB7 and 6W63) were manually prepared by one of the authors (McKay). All histidine side chains were assigned as either HIE or HID given the local environment. Protonation states for histidine residues found near the active site can be found in Table 1. All asparagine and glutamine side chains were examined and found to be in reasonable rotameric states. For these systems, the PARM@FROSST [9] small molecule extension to ff14SB [10] and AM1-BCC [11] charges were used for the ligands. All prepared structures can be found in the github repository.

Molecular dynamics simulations

Molecular dynamics simulations were performed in GPU accelerated AMBER 16 [12] using the ff14SB [10] force field and the optimal point charge (OPC) model [13] of water. Ligand force field parameters were assigned with the general AMBER force field, GAFF [14] using the Antechamber package [15] in AmberTools. Antechamber assigns charges, missing bonds, angles, dihedral angles and Lennard–Jones parameters for each atom. Ligand charges were assigned using AM1-BCC [11].

For systems with a co-crystallized ligand, the ligand was removed from the protein, and then the protein was solvated in a box of OPC water molecules with dimensions that ensured there were at least 10 Å between any atom of the protein and the box edge. Sodium or chlorine counterions were added accordingly to neutralize the system. Each system was then energetically minimized in a two-step process. The first minimization step was performed with 1500 steps of steepest descent with all protein atoms restrained

Table 1 Protonation states of M^{pro}

PDB ID	H41	H163	H164	H172
6LU7	HIE	HIE	HID	HIE
6M03	HID	HID	HIE	HIE
6W63	HIE	HIE	HIE	HIE
6Y84	HID	HID	HIE	HIE
6YB7	HIE	HIE	HIE	HIE

Protonation states for histidine residues near the active site of main protease. These protonation states were shared for all simulations

harmonically using a force constant of $100 \text{ kcal/mol}\cdot\text{\AA}^2$. For the second minimization step, only main chain heavy atoms were restrained. Following minimization, the system was heated to 300 K in a 240 ps NVT simulation with the main chain heavy atoms restrained; the temperature was regulated by Langevin thermostat with collision frequency of 1 ps. This was followed by a 20 ns NPT simulation with the atom restraints declining from 100 to $2.5 \text{ kcal/mol}\cdot\text{\AA}^2$ in the first 10 ns. In the production phase, the temperature was regulated via a Langevin thermostat set to 300 K with a collision frequency of 2 ps. The constant pressure (1 atm) was maintained by isotropic position scaling with a relaxation time of 0.5 ps.

The SARS-CoV-2 main protease is expected to be functionally active as a dimer due to its structural homology with SARS-CoV-1 [16–18]. To investigate possible differences between the monomer and dimer we simulated both the monomeric crystal structure and the reported dimeric biological assemblies for 6W63 and 6YB7. We present data for several different sets of heavy atom restraints to model different amounts of protein flexibility. In one set (denoted rigid in the repository), restraints of $2.5 \text{ kcal/mol}\cdot\text{\AA}^2$ were applied to all heavy atoms. In the second set (denoted SCflex) no restraints were applied to the side chain heavy atoms and restraints of $2.5 \text{ kcal/mol}\cdot\text{\AA}^2$ were applied to the backbone heavy atoms. In the third set (denoted McKay structures) for both monomers and dimeric biological assemblies, carbon atoms were restrained ($2.0 \text{ kcal/mol}\cdot\text{\AA}^2$) for six residues (L87, V91, A206, A211, F294, R298) located distal to the active site. These six-residue restraints were applied to each monomer during the production simulations. All MD simulation files for all simulations are included in the github repository.

GIST

GIST maps were created using the GPU port [19] of AmberTools cpptraj-GIST [20]. Analyses were performed on the complete 50 ns production trajectory for each system (25,000 configurational snapshots). For each system, maps were created in a cubical region with 30\AA length sides centered on the geometric mean position of the co-crystallized ligand for the PDB (see Fig. 1). The resolution of the grid was 0.5\AA (0.125\AA^3 per voxel). For structures with no co-crystallized ligand for the PDB entry, a homologous protein with a co-crystallized ligand was structurally aligned to the PDB structure and the geometric center of that ligand was used to define the GIST analysis region. In the case of 6JYT, the region was defined for HSA by a partial set of the residues found in the active site (K288, S289, D374, E375, R567). For the GIST analysis of 6JYT, the geometric center of ADP from a structurally aligned 2XZL was used as the

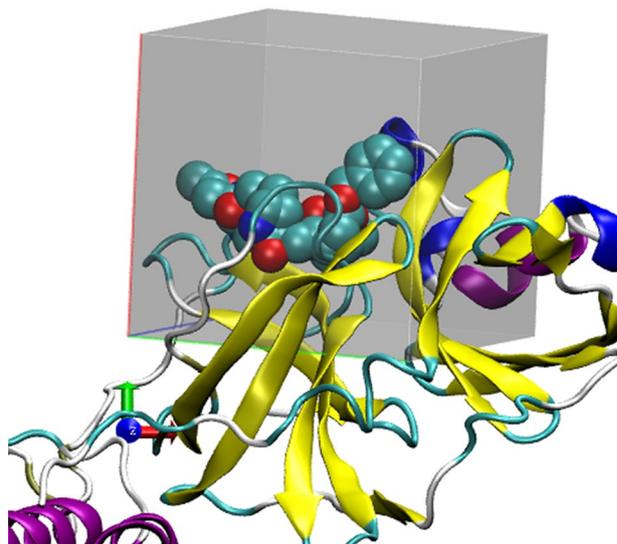


Fig. 1 The co-crystallized structure of M^{Pro} (cartoon) with ligand N3 from 6LU7. The GIST analysis was performed in the cubical region shaded in gray

center of the box. The ligands used for defining the GIST region for each structure can be found in the repository.

Hydration site analysis

Hydration site analysis, HSA [21] was performed using the publicly available SSTMap code [22] with the default settings except for the region analysis which was set to within 10\AA of the ligand ($-d 10$). For each system, the analysis was per the first 20 ns (10,000 frames) of the MD production run for each protein.

Briefly, the method analyzes all the water positions from an MD trajectory and identifies high-density 1\AA radius spherical regions called hydration sites. In each hydration site, average quantities of the water molecules found in the hydration site are calculated and provide estimates for the local IST thermodynamic quantities. A number of measures that describe the local solvent structure and characterize the hydrogen-bonding environment of the water in each hydration site are also calculated. These measures can be used to characterize the enhancement or disruption of local water structure, describe the local enclosure, and describe the average hydrogen bonding interactions that water has in each hydration site with both its water neighbors and protein. Full details of the calculations are specified in a previous publication and the code is available on the github.

We also use newly developed code to determine the most probable orientations for water molecules in each hydration site. To do this, the orientations of all water molecules in each hydration site are clustered using a quaternion distance metric and the centroid orientation of each high-density

cluster (generally at least 10% of the population) is recorded. The code and complete details of the method are in the github.

3D-RISM

3D-RISM maps were created for all rigid structures in the repository using `rism3d.snglpnt` [23] from AmberTools 19 [24]. In each case, the initial protein conformations of the production-phase molecular dynamics simulations were used after all waters had been removed. Water and solvation thermodynamics and number density distributions [25] were calculated for each protein structure with a residual tolerance of 10^{-6} using the partial series expansion of order-3 (PSE-3) closure [26]. The solvation box was extended to include Lennard–Jones interactions $< 10^{-7}$ kcal/mol, which corresponded to a minimum distance between the protein and box edge of 47 Å. A 0.5 Å grid spacing was used with the centering = 3 option, which rounds the positions of the grid points to the nearest 0.5 Å. Oxygen and hydrogen solvent site contributions to the thermodynamics maps were combined using the molecular reconstruction technique [27].

Bulk water solvent-susceptibility input for 3D-RISM was generated with `rism1d` using the coincident SPC/E (cSPC/E) water model [23], the PSE-3 closure and dielectrically consistent RISM (DRISM) theory [28]. The temperature, density and dielectric constant were set to 298.15 K, 55.41 M and 78.45 respectively. A grid of 16,384 points was used with a grid spacing of 0.025 Å.

Repository contents

Structures 1–5 (SARS-CoV-2 structures)

Main protease (M^{pro} , $3CL^{\text{pro}}$): 6LU7 [29] (2.16 Å), 6YB7 (1.25 Å), 6M03 (2.00 Å), 6Y84 (1.39 Å), 6W63 (2.10 Å). Target the substrate binding site of M^{pro} .

Structure 6 (SARS-CoV-1 structure)

Helicase (Nsp13): 6JYT [30] (2.80 Å). Target (1) the ADP binding site but discourage (2) the nucleic acids binding site. No SARS-CoV-2 structure exists for this protein.

Structure 7 (SARS-CoV-2 structure)

Nsp16 (2'-*O*-methyltransferase, nsp 10/16): 6W4H (1.80 Å). Target the *S*-adenosylmethionine (SAM) binding site.

All files with prepared structures, topologies files, and molecular dynamics input and restart files are provided as well as solvation structural and thermodynamic maps described below.

Solvation thermodynamic maps

Inhomogeneous solvation theory, IST [31–33] provides the statistical mechanical framework for the solvation thermodynamic quantities from explicit solvent molecular dynamics simulations. Here, we use two methods: grid-based inhomogeneous solvation theory, GIST [34, 35] and HSA [21] to localize the IST thermodynamic quantities onto a three-dimensional grid and onto high density 1 Å radius spherical “hydration sites”, respectively. These localization approaches both process snapshots of the system configurations generated in molecular dynamics simulations to estimate local IST thermodynamic quantities including local energies, entropies, and number densities.

Grid based solvation maps

The repository contains grid-based solvation maps of calculated IST and 3D-RISM entropies, energies, and densities in Data Explorer (dx) format. The dx format enables visualization in standard graphics packages such as VMD and Pymol. For each target, energetic maps are provided for water's interactions with the protein, with other water molecules, and the total interactions of the water in each voxel with the system as a whole.

GIST provides entropy maps for the total entropy as well maps that separately include the translational and orientational contributions to the total entropy. Maps are provided for all of the entropy and energy quantities for both normalized (per water quantities) and density (per voxel) quantities. A complete list of quantities can be found in Table 2. Detailed descriptions of these quantities can be found in our prior work [20, 35].

3D-RISM provides maps of the total entropy and the solvation free energy, as well as maps solvation energy. In all cases, only densities maps are provided. A list of quantities is provided in Table 3 and descriptions can be found in Nguyen et al. [27].

3D-RISM maps differ from GIST maps in that they are not limited to the ligand binding site and do not rely on solvent sampling, which may be incomplete. However, the orientations of water molecules are lost in this process and the density distributions may differ from those of explicit solvent, primarily in the height and breadth of the maxima and minima. As the files for the full thermodynamic and number density distributions are quite large, files available in the github repository have been truncated to the binding site and match the dimensions of the GIST maps.

Table 2 Key GIST quantities

Quantity	Description	Units
$^{[a]}TS_{\text{six}}$	Total entropy density	kcal/mol/Å ³
$^{[a]}TS_{\text{trans}}$	Translational entropy density	kcal/mol/Å ³
$^{[a]}TS_{\text{orient}}$	Oriental entropy density	kcal/mol/Å ³
$^{[a]}E_{\text{ww}}$	Water–water energy density	kcal/mol/Å ³
$^{[a]}E_{\text{sw}}$	Solute–water energy density	kcal/mol/Å ³
Neighbor count	Mean number of water neighbors ^[b]	Molecules

[a] Corresponding normalized quantities also reported

[b] Neighbors are defined as two water molecules with an O–O distance of 3.5 Å or less

Table 3 Key 3D-RISM quantities

Quantity	Description	Units
ΔG	Solvation free energy	kcal/mol/Å ³
–TS	Total entropy density	kcal/mol/Å ³
E_{tot}	Total energy density	kcal/mol/Å ³
E_{ww}	Water–water energy density	kcal/mol/Å ³
E_{sw}	Solute–water energy density	kcal/mol/Å ³

Hydration site solvation maps

For each target, the positions and calculated thermodynamic and structural quantities for the water in each hydration site are summarized in a space delimited spreadsheet file.

The same energetic quantities as calculated for GIST (above) are calculated for each hydration site and reported in per water (normalized) units. Additionally, the HSA data includes a breakdown of the total energy into contributions from Lennard–Jones, electrostatic, and first solvation shell water–water interactions.

SSTMap also calculates a number of quantities that are aimed at characterizing the local environment surrounding each hydration site. These are aimed at better describing local water structure and the interactions of the water in the hydration site with the protein surface.

Quantities that provide a measure of local water structure include the average number of first shell neighbors each water has in its first solvation shell, the fraction of these neighbors to which the hydration site water is hydrogen bonded, and the average energy of interaction with each neighboring water. When compared to bulk water values, these quantities provide measures of whether the local water structure is enhanced or frustrated [36].

Additional quantities that characterize the interaction of the water in each hydration site with the protein include: (1) an enclosure parameter that describes how much of the region around the hydration site is protein and how much

is water, (2) the average number of hydrogen bond donor and acceptor interactions that water molecules found in the hydration site have with the protein surface, and (3) lists of the protein residues that donate and accept hydrogen bonds to the water in the hydration site.

A list of thermodynamic and structural quantities can be found in Table 4. A text delimited spreadsheet file summarizing all calculated water properties is found in the HSA directory for each protein.

In addition, to facilitate visualization, each HSA directory includes PDB files that feature (1) the hydration site centers, (2) water molecules located at the center of each hydration site that have the most probable orientation, and (3) water molecules located at the center of each hydration site that include all probable orientation clusters.

Potential applications

Solvation thermodynamic mapping has been used in a variety of applications aimed at aiding the discovery and design of new pharmaceutical compounds. In docking, scoring terms have been added to explicitly account for solvent displacement upon ligand binding and the modified docking scoring functions have been used to help improve AUC, pose prediction, and identify novel binding ligands [1, 6, 37]. Solvation maps have also been used to create pharmacophores [2] as well as provide criteria to prioritize the selection of pharmacophore sites [5]. Both water thermodynamics and water interactions with protein surfaces have been used to direct lead modification [4, 38].

Here, we describe by example several potential applications for the GIST and HSA solvation maps provided in this repository. 3D-RISM solvation maps can be used as a complement or alternative to GIST and HSA results as they treat regions where water molecules may not be able to exchange during normal sampling, allow the hydration of the entire protein to be explored, and are based on a

Table 4 HSA structural quantities

Quantity	Description	Units
N_{nbr}	Average # first shell neighbors	None
$N_{\text{ww}}^{\text{HB}}$	Average # water–water hydrogen bonds	None
$N_{\text{sw}}^{\text{HB}}$	# Solute–water hydrogen bonds	kcal/mol
$E_{\text{nbr}}^{\text{ww}}$	Average water–water interaction energy by neighbor	kcal/nbr
$N_{\text{ww}}^{\text{HB,don}}$	# Water–water hydrogen bonds donated	None
$N_{\text{ww}}^{\text{HB,acc}}$	# Water–water hydrogen bonds accepted	None
$N_{\text{sw}}^{\text{HB,don}}$	# Solute–water hydrogen bonds donated	None
$N_{\text{sw}}^{\text{HB,acc}}$	# Solute–water hydrogen bonds accepted	None
$f_{\text{ww}}^{\text{HB}}$	Fraction of hydrogen-bonded neighbors	None

different theoretical framework. Certainly, users should have greater confidence for regions where all three methods agree. As the 3D-RISM maps extend beyond the active site, they may be used to investigate the dimer interface as well as investigate potential allosteric modulation distal from the substrate binding pocket.

Rational lead modification

The properties of water in and around the binding site may be used to direct the design of chemical modifications to a lead compound or fragment. The physical principles of this are that the displacement of thermodynamically unfavorable surface water upon the binding of a ligand will lead to favorable contributions to the free energy as the water is displaced to the more thermodynamically favorable environment of bulk biological water.

Here, we illustrate how solvation structural and thermodynamic solvation mapping in this repository can be used to provide insight into which modifications may lead to boosts in binding affinity.

The binding site of M^{pro} features a large number of energetically unfavorable hydration sites (see Fig. 2). Prior work [39, 40] suggests that the displacement of water from these hydration sites may be correlated with differences in binding affinities between congeneric pairs of ligands. Most of the hydration sites identified in Fig. 1 are displaced by N3. However, the two leftmost sites are not. We will focus on the upper left site, hydration site 7 (HS7), as it has an exceptionally unfavorable thermodynamic profile.

HS7 occupies a small cleft on the surface of the protein, which is formed by seven different residues (N28, G143, N119, T26, Y118, and C145). The water in this cleft is resolved the crystal structures of 6LU7, 6W63, 6Y84, and 6YB7. However, this water is not reported in 6M03. The water is highly enclosed by the protein (81.7%) having slightly less than 1 (0.96) water neighbor, on average, in its first solvation shell. Despite the hydration site being highly occupied (84.5% occupancy), the water is exceptionally unfavorable energetically (+2.6 kcal/mol) and entropically (−TS of 4.45 kcal/mol) by IST estimates. Its low entropy result is based on the water's highly restricted translational and orientational motion. The water's high enclosure in the protein cleft and its formation of two hydrogen bonds with the protein surface severely restrict the water's translational freedom leading to a translational entropic penalty of 2.11 kcal by IST estimates. The two hydrogen bonds it forms with the protein surface as well as forming a hydrogen bond 82% of the time with its water neighbor located above the cleft (HS56), further restrict its orientational freedom resulting in an entropic penalty of 2.33 kcal/mole.

Despite being on a hydrophilic surface (forming on average 2.00 hydrogen bonds with the protein), the water in HS7

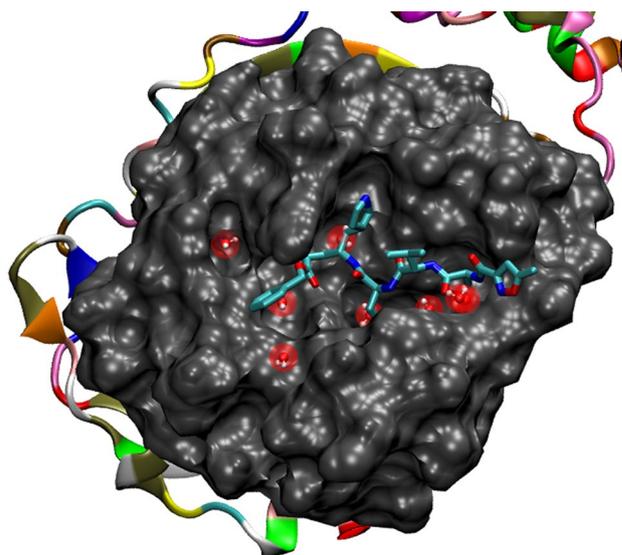


Fig. 2 N3 bound to M^{pro} (PDB ID 6LU7). Hydration sites that are located within 7.5 angstroms of N3 and have highly unfavorable energy ($\Delta E > 0.5$ kcal/mol with respect to neat water) are shown as transparent red spheres. The most probable water orientation for each hydration site is represented by a water molecule at the center of each sphere. The protein surface proximal (within 11 Å) to N3 is shown in gray

cannot form a full complement of hydrogen bonds, instead forming only 2.85 geometric hydrogen bonds on average compared to a bulk OPC water which would form 3.62. This deficiency of more than three quarters of a hydrogen bond, on average, is a significant contribution to the unfavorable energetic profile (+2.6 kcal/mole overall) of HS7.

Both the unfavorable IST energy and entropy suggest that displacing this HS7 water could lead to gains in binding affinity. In order to displace this water, an optimal chemical group must replace interactions that the water makes with the protein without disrupting the hydrogen bond network that the water is making with its neighbors. As the water in HS7 is located in a cleft, any chemical group would also need to displace its water neighbor (h-bonded water in Fig. 3). The optimal chemical group would need to both donate a hydrogen bond to the backbone carbonyl of G143 and accept a hydrogen bond from the backbone amine of N119. A hydroxy group seems ideal for this.

All of the numerical data in the above analysis is located in the HSA summary spreadsheet for 6LU7 (6LU7_apo_flex_hsa_summary.csv). All the data for the visualizations is likewise located in the repository.

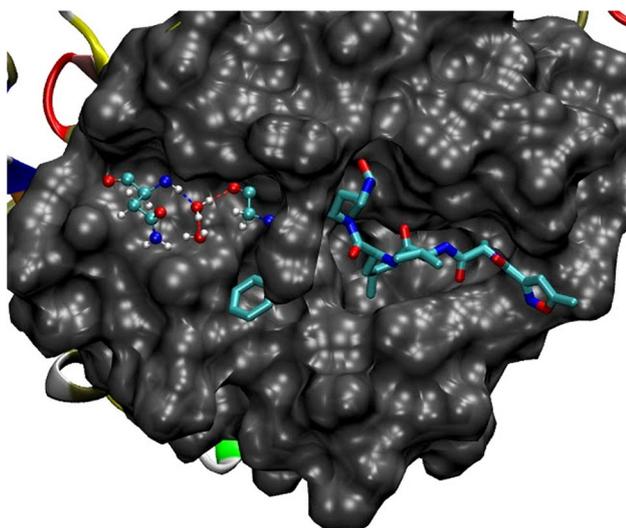


Fig. 3 The most probable orientation of the water in HS7 donates a hydrogen bond (red dashed line) to the backbone carbonyl of Gly143, accepts a hydrogen bond (blue dashed line) from the backbone NH of Asn19, and donates a hydrogen bond to HS56 above the cleft wherein lies HS7

Scoring solvation displacement in docking

Four studies outline how solvation thermodynamic mapping can be used to aid in the discovery of new leads in docking. The first two of these studies are based on our prior work on Factor Xa [39, 40], in which a displaced solvent functional used high energy and high density voxels as functional inputs to correlate with experimental measurements of differences in binding free energies between congeneric pairs of ligands [39, 40]. The third docking study [6] by Uehara and Tanaka instead used a displaced solvent functional with free energetic maps created by GIST as input whereas the fourth study [1] by Balias et al. used the displacement of voxels with high energy densities as input. The third study showed improvements in pose prediction and enrichment and the fourth showed only nominal measurable improvements to docking enrichment and pose prediction, though the method was successfully used to prospectively identify new tightly binding compounds, including the tightest binding compound to cytochrome *c* peroxidase. A map showing related unfavorable and favorable energy density regions for M^{pro} is shown in Fig. 4.

The GIST maps in this repository provide the data to create the maps used in all three of the GIST-based studies. Necessary modifications of the provided GIST dx maps (e.g. creating a free energy density map from the energy and entropy density maps) can be easily created using the GIST post-processing (GISTPP) code provided on the github.

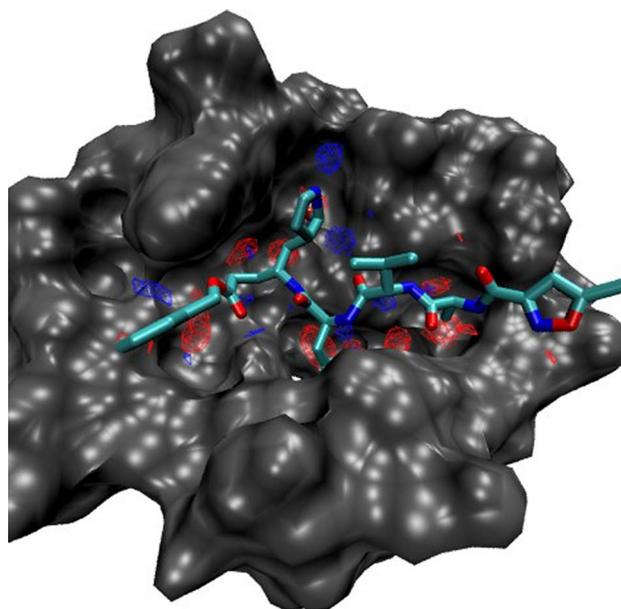


Fig. 4 Unfavorable and favorable solvation energy density map of M^{pro} . Regions of unfavorable energy density ($E_{dens} > 0.1$ kcal/mole/ \AA^3) and favorable energy density ($E_{dens} < -0.1$ kcal/mole/ \AA^3) are shown in red or blue wireframe, respectively. The predicted score for a docked ligand would be penalized for displacing water from the favorable blue regions or given an affinity boost for displacing water from the red regions

Pharmacophore creation

Solvation mapping can be used to generate water-based pharmacophore hypotheses [2] and to prioritize ligand- or protein-based pharmacophore sites [5]. Here we combine several interesting hydration sites with ligand-based pharmacophore elements.

Three pharmacophore sites were constructed using ligand–protein interactions based on analyses of co-crystallized ligands found inside the binding sites of SARS-CoV-2 M^{pro} structures (PDB ID 6W63, 6LU7, 6Y2F, 6Y2G, and 6M2N). These ligand-based sites appear as dotted spheres in Figs. 5 and 6.

The leftmost ligand-acceptor site (Figs. 5, 6) lies inside the oxyanion hole. All five of the co-crystallized ligands accept a hydrogen bond from the backbone amino group of G143 while three of five (6Y2F, 6Y2G, and 6M2N) also accept a hydrogen bond from C146. The pharmacophore site shown in Fig. 6 shows both of these interactions. The middle ligand-based site donates a hydrogen bond to the backbone carbonyl of H164. Ligands from 6W63, 6Y2G and 6Y2F make this contact. The rightmost ligand site, inside the S1 subsite, accepts a hydrogen bond from the backbone amino group of E166. Four of the five (all except 6M2N) co-crystallized ligands accept a hydrogen bond from this group. Each ligand-based site is proximal to a hydration site

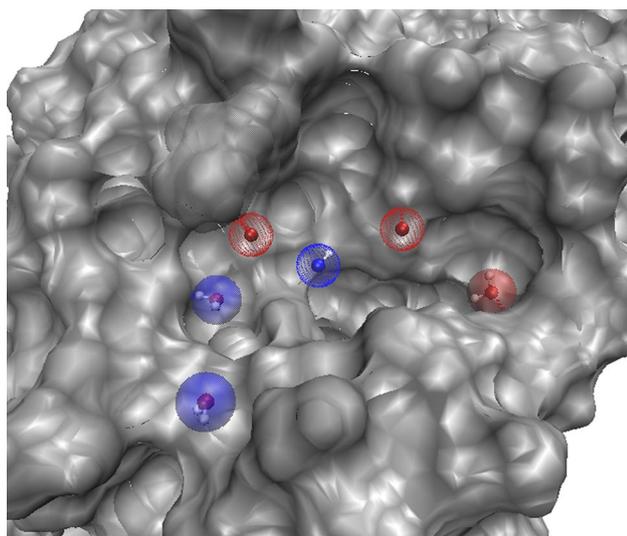


Fig. 5 Hybrid ligand- and water-based pharmacophore within the binding site of M^{Pro} (PDB ID 6LU7). The ligand-based sites are shown as dotted spheres and the water-based sites are shaded spheres. Ligand-based sites have an NH group for donors or an oxygen for acceptors. The most probable water orientation is found at the center of each water-based pharmacophore site. Acceptor sites are red and donor sites are blue spheres

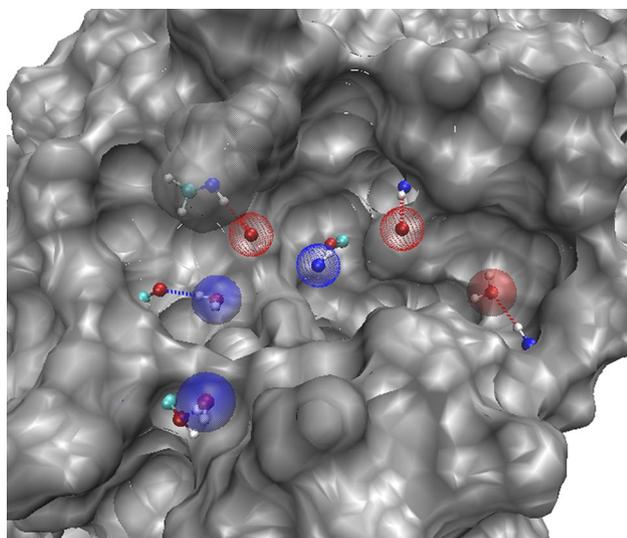


Fig. 6 The same hybrid pharmacophore hypothesis as shown in Fig. 5, except the interactions with chemical groups on the surface are shown explicitly. Blue dashed lines show the pharmacophore sites donation of hydrogen bonds and red dashed lines show acceptance

and GIST high-density group of voxels but none have any significant thermodynamic signal for use in prioritization. These ligand-based sites were chosen by the fact that they were well conserved across the limited number of structures available with co-crystallized ligands.

We used hydration site analysis to add three additional ligand-based pharmacophore sites. These sites are shown in shaded spheres in Figs. 5 and 6. While water-based pharmacophore sites can be chosen using other criteria (as outlined in Jung et al. [2]), here we simply chose water-based sites that are energetically unfavorable and categorized them based on their donor/acceptor interactions with the protein surface.

The first site (on the far right of Figs. 5, 6) is from HS9, which primarily accepts a hydrogen bond from the backbone amino group of residue T190 and has an unfavorable energy of -11.28 kcal/mole (almost 1 kcal above the bulk energy of -12.26 kcal/mole). The second site HS52 (middle left in Figs. 5, 6) has an energy of -10.46 kcal/mole (1.8 kcal above bulk energy) and donates a hydrogen bond to T26. The third site, HS56 has an energy of -11.69 kcal/mole (0.57 kcal less favorable than bulk) and donates a hydrogen bond to T25.

Together, the conserved ligand sites and the water-based sites create a pharmacophore hypothesis that can be used to screen virtual compound databases.

While we arbitrarily chose three conserved sites from the ligand and three proximal hydration sites to construct the hypothesis outlined here, this approach allows a drug designer flexibility to choose ligand and water sites on virtually any solvent exposed surface of the protein, allowing different regions of the active site or potential allosteric sites to be targeted.

How to access data

All hydration site, 3D-RISM and GIST data is available on github with a readme.md that details directory structure and the descriptive file naming convention. Briefly, each PDB structure has its own subdirectory named after its PDB ID. Each PDB ID subdirectory has further subdirectories for simulations with apo or complexed structures and different protein restraints. Additional subdirectories for each of these include the hydration site, 3D-RISM and GIST analyses, as well as the prepared protein input files and Amber MD restart files in case longer simulations are desired. All of the above can be found on the github (github.com/KurtzmanLab/COVID19_GIST_HSA).

Acknowledgements We'd like to thank Daniel Roe for extensive help in integrating GIST into Amber Tools. José Duca, Camilo Velez-Vega, Callum Dickson, Andre Golosov, and the Novartis CADD Team for provocative discussions on protein structure and preparation.

Funding This research was supported by the USA National Institutes of Health R01GM100946, National Science Foundation Grant No. 1566638 and the Research Corporation for Science Advancement (RCSA) Cottrell Scholar Award #23967.

Data availability All data are publicly available on github.

Code availability All water analysis code used to produce this data is open-source with extensive documentation and has been made publicly available for download. Four sets of code were used for the water analysis in the repository: SSTMap, GIST-cpptraj, and GISTPP. SSTMap was used for hydration site analyses, GIST-cpptraj was used for the GIST analyses, GISTPP was used to make numerical manipulations to the GIST dx files, rism3d.snglpt was used for the 3D-RISM analysis. Usage tutorials and documentation can be found on the SSTMap project page (SSTMap.org) and on the AMBER website. GIST-cpptraj code is available on the Amber-MD github (<https://github.com/Amber-MD>). All other code is available on the github (<https://github.com/KurtzmanLab>).

Compliance with ethical standards

Conflict of interest We declare no conflict of interest.

Ethical approval This publication has followed the guidelines of the Committee on Publication Ethics.

References

- Balius TE, Fischer M, Stein RM, Adler TB, Nguyen CN, Cruz A, Gilson MK, Kurtzman T, Shoichet BK (2017) Testing inhomogeneous solvation theory in structure-based ligand discovery. *Proc Natl Acad Sci USA* 114:E6839–E6846. <https://doi.org/10.1073/pnas.1703287114>
- Jung SW, Kim M, Ramsey S, Kurtzman T, Cho AE (2018) Water pharmacophore: designing ligands using molecular dynamics simulations with water. *Sci Rep* 8:10400. <https://doi.org/10.1038/s41598-018-28546-z>
- Harriman G, Greenwood J, Bhat S, Huang X, Wang R, Paul D, Tong L, Saha AK, Westlin WF, Kapeller R, Harwood HJ (2016) Acetyl-CoA carboxylase inhibition by ND-630 reduces hepatic steatosis, improves insulin sensitivity, and modulates dyslipidemia in rats. *Proc Natl Acad Sci USA* 113:E1796–E1805. <https://doi.org/10.1073/pnas.1520686113>
- Collin M-P, Lobell M, Hübsch W, Brohm D, Schirok H, Jautelat R, Lustig K, Bömer U, Vöhringer V, Héroult M, Grünewald S, Hess-Stumpp H (2018) Discovery of Rogaratinib (BAY 1163877): a pan-FGFR inhibitor. *ChemMedChem* 13:437–445. <https://doi.org/10.1002/cmdc.201700718>
- Hu B, Lill MA (2012) Protein pharmacophore selection using hydration-site analysis. *J Chem Inf Model* 52:1046–1060. <https://doi.org/10.1021/ci200620h>
- Uehara S, Tanaka S (2016) AutoDock-GIST: incorporating thermodynamics of active-site water into scoring function for accurate protein-ligand docking. *Molecules* 21:1604. <https://doi.org/10.3390/molecules21111604>
- Madhavi Sastry G, Adzhigirey M, Day T, Annabhimoju R, Sherman W (2013) Protein and ligand preparation: parameters, protocols, and influence on virtual screening enrichments. *J Comput Aided Mol Des* 27:221–234. <https://doi.org/10.1007/s1082-013-9644-8>
- Schrödinger LLC (2017) Schrödinger Release 2017-3. Schrödinger, LLC, New York
- Bayly C, McKay D, Truchon J (2010) An informal AMBER small molecule force field: Parm@ Frosst
- Maier JA, Martinez C, Kasavajhala K, Wickstrom L, Hauser KE, Simmerling C (2015) ff14SB: improving the accuracy of protein side chain and backbone parameters from ff99SB. *J Chem Theory Comput* 11:3696–3713. <https://doi.org/10.1021/acs.jctc.5b00255>
- Jakalian A, Jack DB, Bayly CI (2002) Fast, efficient generation of high quality atomic charges. AM1-BCC Model: II. Parameterization and validation. *J Comput Chem*. <https://doi.org/10.1002/jcc.10128>
- Case DA, Betz RM, Cerutti DS, Cheatham TE III, Darden TA, Duke RE, Giese TJ, Gohlke H, Goetz AW, Homeyer N, Izadi S, Janowski P, Kaus J, Kovalenko A, Lee TS, LeGrand S, Li P, Lin C, Luchko T, Luo R, Madej B, Mermelstein D, Merz KM, Monard G, Nguyen H, Nguyen HT, Omelyan I, Onufriev A, Roe DR, Roitberg A, Sagui C, Simmerling CL, Botello-Smith WM, Swails J, Walker RC, Wang J, Wolf RM, Wu X, Xiao L, Kollman PA (2016) Amber 16. University of California, San Francisco
- Izadi S, Anandakrishnan R, Onufriev AV (2014) Building water models: a different approach. *J Phys Chem Lett* 5:3863–3871. <https://doi.org/10.1021/jz501780a>
- Wang J, Wolf RM, Caldwell JW, Kollman PA, Case DA (2004) Development and testing of a general Amber force field. *J Comput Chem*. <https://doi.org/10.1002/jcc.20035>
- Wang J, Wang W, Kollman PA, Case DA (2006) Automatic atom type and bond type perception in molecular mechanical calculations. *J Mol Graph Model* 25:247–260. <https://doi.org/10.1016/j.jmgm.2005.12.005>
- Anand K, Palm GJ, Mesters JR, Siddell SG, Ziebuhr J, Hilgenfeld R (2002) Structure of coronavirus main proteinase reveals combination of a chymotrypsin fold with an extra alpha-helical domain. *EMBO J* 21:3213–3224. <https://doi.org/10.1093/emboj/cdf327>
- Shi J, Song J (2006) The catalysis of the SARS 3C-like protease is under extensive regulation by its extra domain. *FEBS J* 273:1035–1045. <https://doi.org/10.1111/j.1742-4658.2006.05130.x>
- Zhang L, Lin D, Sun X, Curth U, Drosten C, Sauerhering L, Becker S, Rox K, Hilgenfeld R (2020) Crystal structure of SARS-CoV-2 main protease provides a basis for design of improved α -ketoamide inhibitors. *Science*. <https://doi.org/10.1126/science.abb3405>
- Kraml J, Kamenik AS, Waibl F, Schauerl M, Liedl KR (2019) Solvation free energy as a measure of hydrophobicity: application to serine protease binding interfaces. *J Chem Theory Comput* 15:5872–5882. <https://doi.org/10.1021/acs.jctc.9b00742>
- Ramsey S, Nguyen C, Salomon-Ferrer R, Walker RC, Gilson MK, Kurtzman T (2016) Solvation thermodynamic mapping of molecular surfaces in AmberTools: GIST. *J Comput Chem* 37:2029–2037. <https://doi.org/10.1002/jcc.24417>
- Young T, Abel R, Kim B, Berne BJ, Friesner RA (2007) Motifs for molecular recognition exploiting hydrophobic enclosure in protein–ligand binding. *Proc Natl Acad Sci USA* 104:808–813. <https://doi.org/10.1073/pnas.0610202104>
- Haider K, Cruz A, Ramsey S, Gilson MK, Kurtzman T (2018) Solvation structure and thermodynamic mapping (SSTMap): an open-source, flexible package for the analysis of water in molecular dynamics trajectories. *J Chem Theory Comput* 14:418–425. <https://doi.org/10.1021/acs.jctc.7b00592>
- Luchko T, Gusarov S, Roe DR, Simmerling C, Case DA, Tuszynski J, Kovalenko A (2010) Three-dimensional molecular theory of solvation coupled with molecular dynamics in Amber. *J Chem Theory Comput* 6:607–624. <https://doi.org/10.1021/ci900460m>
- Case DA, Ben-Shalom IY, Brozell SR, Cerutti DS, Cheatham TEI, Cruzeiro VWD, Duke RE, Darden TA, Ghoreishi D, Giambasu G, Giese T, Gilson MK, Gohlke H, Goetz AW, Greene D, Harris R, Homeyer N, Huang Y, Izadi S, Kovalenko A, Krasny R, Kurtzman T, Lee TS, LeGrand S, Li P, Lin C, Liu J, Luchko T, Luo R, Man V, Mermelstein DJ, Merz KM, Miao Y, Monard G, Nguyen C, Nguyen H, Onufriev A, Pan F, Qi R, Roe DR, Roitberg A, Sagui

- C, Schott-Verdugo S, Shen J, Simmerling CL, Smith J, Swails J, Walker RC, Wang J, Wei H, Wilson L, Wolf RM, Wu X, Xiao L, Xiong Y, York DM, Kollman PA (2019) Amber 2019. University of California, San Francisco
25. Johnson J, Case DA, Yamazaki T, Gusarov S, Kovalenko A, Luchko T (2016) Small molecule hydration energy and entropy from 3D-RISM. *J Phys Condens Matter* 28:344002. <https://doi.org/10.1088/0953-8984/28/34/344002>
 26. Kast SM, Kloss T (2008) Closed-form expressions of the chemical potential for integral equation closures with certain bridge functions. *J Chem Phys* 129:236101. <https://doi.org/10.1063/1.3041709>
 27. Nguyen C, Yamazaki T, Kovalenko A, Case DA, Gilson MK, Kurtzman T, Luchko T (2019) A molecular reconstruction approach to site-based 3D-RISM and comparison to GIST hydration thermodynamic maps in an enzyme active site. *PLoS ONE* 14:e0219473. <https://doi.org/10.1371/journal.pone.0219473>
 28. Perkyns J, Pettitt BM (1992) A site–site theory for finite concentration saline solutions. *J Chem Phys* 97:7656–7666
 29. Jin Z, Du X, Xu Y, Deng Y, Liu M, Zhao Y, Zhang B, Li X, Zhang L, Peng C, Duan Y, Yu J, Wang L, Yang K, Liu F, Jiang R, Yang X, You T, Liu X, Yang X, Bai F, Liu H, Liu X, Guddat LW, Xu W, Xiao G, Qin C, Shi Z, Jiang H, Rao Z, Yang H (2020) Structure of Mpro from SARS-CoV-2 and discovery of its inhibitors. *Nature* 582:289–293. <https://doi.org/10.1038/s41586-020-2223-y>
 30. Jia Z, Yan L, Ren Z, Wu L, Wang J, Guo J, Zheng L, Ming Z, Zhang L, Lou Z, Rao Z (2019) Delicate structural coordination of the Severe Acute Respiratory Syndrome coronavirus Nsp13 upon ATP hydrolysis. *Nucleic Acids Res* 47:6538–6550. <https://doi.org/10.1093/nar/gkz409>
 31. Morita T, Hiroike K (1961) A new approach to the theory of classical fluids. III: general treatment of classical systems. *Prog Theor Phys* 25:537–578. <https://doi.org/10.1143/PTP.25.537>
 32. Lazaridis T (1998) Inhomogeneous fluid approach to solvation thermodynamics. 1. Theory. *J Phys Chem B* 102:3531–3541. <https://doi.org/10.1021/jp9723574>
 33. Lazaridis T (1998) Inhomogeneous fluid approach to solvation thermodynamics. 2. Applications to simple fluids. *J Phys Chem B* 102:3542–3550. <https://doi.org/10.1021/jp972358w>
 34. Nguyen CN, Young TK, Gilson MK (2012) Grid inhomogeneous solvation theory: hydration structure and thermodynamics of the miniature receptor cucurbit[7]uril. *J Chem Phys* 137:044101–44117
 35. Nguyen C, Gilson MK, Young T (2011) Structure and thermodynamics of molecular hydration via grid inhomogeneous solvation theory. arXiv:11084876
 36. Haider K, Wickstrom L, Ramsey S, Gilson MK, Kurtzman T (2016) Enthalpic breakdown of water structure on protein active-site surfaces. *J Phys Chem B* 120:8743–8756. <https://doi.org/10.1021/acs.jpcc.6b01094>
 37. Murphy RB, Repasky MP, Greenwood JR, Tubert-Brohman I, Jerome S, Annabhimoju R, Boyles NA, Schmitz CD, Abel R, Farid R, Friesner RA (2016) WScore: a flexible and accurate treatment of explicit water molecules in ligand–receptor docking. *J Med Chem* 59:4364–4384. <https://doi.org/10.1021/acs.jmedchem.6b00131>
 38. Abel R, Mondal S, Masse C, Greenwood J, Harriman G, Ashwell MA, Bhat S, Wester R, Frye L, Kapeller R, Friesner RA (2017) Accelerating drug discovery through tight integration of expert molecular design and predictive scoring. *Curr Opin Struct Biol* 43:38–44. <https://doi.org/10.1016/j.sbi.2016.10.007>
 39. Abel R, Young T, Farid R, Berne BJ, Friesner RA (2008) Role of the active-site solvent in the thermodynamics of Factor Xa ligand binding. *J Am Chem Soc* 130:2817–2831. <https://doi.org/10.1021/ja0771033>
 40. Nguyen CN, Cruz A, Gilson MK, Kurtzman T (2014) Thermodynamics of water in an enzyme active site: grid-based hydration analysis of coagulation Factor Xa. *J Chem Theory Comput* 10:2769–2780. <https://doi.org/10.1021/ct401110x>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.