

Supercomputer-Based Ensemble Docking Drug Discovery Pipeline with Application to Covid-19

A. Acharya¹, R. Agarwal²⁻⁴, M. Baker⁵, J. Baudry⁶, D. Bhowmik⁷, S. Boehm⁵, K. G. Byler⁶, L. Coates⁸, S.Y. Chen⁹, C.J. Cooper^{2,4}, O. Demerdash¹⁰, I. Daidone¹¹, J.D. Eblen^{2,3}, S. Ellingson¹³, S. Forli¹⁴, J. Glaser¹⁵, J. C. Gumbart¹, J. Gunnels¹⁶, O. Hernandez⁵, S. Irle^{7,17,18}, J. Larkin¹⁹, T.J. Lawrence¹⁰, S. LeGrand¹⁹, S.-H. Liu^{2,3}, J.C. Mitchell¹⁰, G. Park⁹, J.M. Parks²⁻⁴, A. Pavlova¹, L. Petridis^{2,3}, D. Poole¹⁹, L. Pouchard⁹, A. Ramanathan²⁰, D. Rogers¹⁵, D. Santos-Martins¹⁴, A. Scheinberg²¹, A. Sedova¹⁰, S. Shen²⁻⁴, J.C. Smith^{*2,3}, M.D. Smith^{2,3}, C. Soto⁹, A. Tsaris¹⁵, M. Thavappiragasam¹⁰, A.F. Tillack¹⁴, J.V. Vermaas¹⁵, V.Q. Vuong^{7,17,18}, J. Yin¹⁵, S. Yoo⁹, M. Zahran²², L. Zanetti-Polzi²³

¹School of Physics, Georgia Institute of Technology, Atlanta, GA 30332

²UT/ORNL Center for Molecular Biophysics, Oak Ridge National Laboratory, TN, 37830

³The University of Tennessee, Knoxville. Department of Biochemistry & Cellular and Molecular Biology, 309 Ken and Blaire Mossman Bldg. 1311 Cumberland Avenue Knoxville, TN, 37996

⁴Graduate School of Genome Science and Technology, University of Tennessee, Knoxville, TN, 37996

⁵Computer Science and Mathematics Division, Oak Ridge National Lab, Oak Ridge, TN 37830

⁶The University of Alabama in Huntsville, Department of Biological Sciences. 301 Sparkman Drive, Huntsville, AL 35899

⁷Computational Sciences and Engineering Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831

⁸Neutron Scattering Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831

⁹Computational Science Initiative, Brookhaven National Laboratory, Upton, NY 11973

¹⁰Biosciences Division, Oak Ridge National Lab, Oak Ridge, TN 37830

¹¹Department of Physical and Chemical Sciences, University of L'Aquila, I-67010 L'Aquila, Italy

¹³University of Kentucky, Division of Biomedical Informatics, College of Medicine, UK Medical Center MN 150, Lexington KY, 40536

¹⁴Scripps Research, La Jolla, CA, 92037

¹⁵National Center for Computational Sciences, Oak Ridge National Laboratory, Oak Ridge, TN 37830

¹⁶HPC Engineering, Amazon Web Services, Seattle, WA 98121

¹⁷Chemical Sciences Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831

¹⁸Bredesen Center for Interdisciplinary Research and Graduate Education, University of Tennessee, Knoxville, TN 37996

¹⁹NVIDIA Corporation, Santa Clara, CA 95051

²⁰Data Science and Learning Division, Argonne National Lab, Lemont, IL 60439

²¹Jubilee Development, Cambridge MA 02139

²²Department of Biological Sciences, New York City College of Technology, The City University of New York (CUNY), Brooklyn, NY 11201

²³CNR Institute of Nanoscience, I-41125 Modena, Italy

This manuscript has been authored by UT-Battelle, LLC under Contract No. DE-AC05-00OR22725 with the U.S. Department of Energy. The United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United States Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes. The Department of Energy will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

Abstract

We present a supercomputer-driven pipeline for in-silico drug discovery using enhanced sampling molecular dynamics (MD) and ensemble docking. Ensemble docking makes use of MD results by docking compound databases into representative protein binding-site conformations, thus taking into account the dynamic properties of the binding sites. We also describe preliminary results obtained for 23 systems involving eight proteins of the proteome of SARS-CoV-2. The MD involves temperature replica exchange enhanced sampling, making use of massively-parallel supercomputing to quickly sample the configurational space of protein drug targets. Using the Summit supercomputer at the Oak Ridge National Laboratory, more than 1 ms of enhanced sampling MD can be generated per day. We have ensemble docked repurposing databases to ten configurations of each of the 23 SARS-CoV-2 systems using AutoDock Vina. We also demonstrate that using Autodock-GPU on Summit, it is possible to perform exhaustive docking of one billion compounds in under 24 hours. Finally, we discuss preliminary results and planned improvements to the pipeline, including the use of quantum mechanical (QM), machine learning, and AI methods to cluster MD trajectories and rescore docking poses.

Introduction

A typical drug takes 10-15 years and \$1-1.5B to develop, due mainly to the need for extensive preclinical testing and high failure rates in clinical trials.¹ Accelerating the drug discovery process is therefore of major concern, and this need has now been placed in even sharper relief by the Covid-19 pandemic. When a target is known and a specific assay developed, experimental high-throughput screening of compounds synthesized using combinatorial chemistry can be performed, and this is a mainstay of drug discovery. However, this approach has had a checkered past, with notable failures, for example, in antimicrobial efforts.² A logical alternative is to mimic what happens in nature, using ‘structure-based’ drug discovery. This can involve structural biology experiments, such as crystallography and cryo-EM, combined with a variety of assays. Moreover, the availability of many experimental protein structures combined with massive increases in computational power and methodological advances have led to a resurgence of computational studies in which trial compounds are docked into binding sites in three-dimensional models of the protein targets and then ranked according to their strength of binding. Computational docking has been particularly useful in early stages of molecular discovery in order to identify initial hits to be prioritized for experimental validation.

Early docking studies were performed with static target crystal structures and rigid ligands. These were quite successful in some cases, such as in the discovery of antivirals for HIV and influenza^{3, 4}. Unfortunately, though, at that time, structures for few targets existed, and the process was relatively inefficient: calculations were relatively inaccurate, and computers could dock only ~100 compounds in a reasonable timeframe. Since the 1990s, the power of supercomputers has increased by a factor of a million or so. Rigid docking of over a billion compounds has been performed in a few days. Thus, virtual high-throughput screening has outperformed equivalent experimental high-throughput screening and has been shown to rapidly identify very tightly binding compounds⁵.

Strictly rigid docking does not often take place in protein: ligand interactions, as both ligands and proteins, undergo thermally-driven internal motions, which lead to fluctuating binding site conformations⁶. Therefore, a particularly important development has been the recognition that incorporating target flexibility into drug discovery protocols can improve the drug discovery process⁷. Ensemble docking uses different conformations of the protein targets of interest, and combinatorially performs the docking of databases of compounds against each of the protein target conformations. This process models the

conformational selection binding mechanism, as opposed to a more limited induced-fit mechanism. The method requires the generation of an ensemble of protein conformations to be used in the docking calculations.

Ensemble docking of small probe molecules for flexible pharmacophore modeling was introduced in 1999. It was shown that consensus pharmacophore models, based on multiple MD structures or on multiple crystallographic structures, were more successful than models based on single conformations in yielding successful predictions of binding⁸. In our own labs, ensemble docking has produced experimentally validated hits against each of the 16 protein targets presented to us over the past few years. Our groups have increasingly used an ensemble approach to perform docking⁹⁻²⁵. In addition, we have shown that the clustering of protein target MD trajectories usually brings a large improvement in the quality of ensemble docking compared to what is obtained using single structure docking²⁶.

Ensemble generation using MD and docking both require significant computational power - to perform MD simulations of sufficient duration and to dock large databases of compounds. This combinatorially large computational time requirement essentially limits this approach to high-performance computing (HPC) architectures for large database screenings, even when only a subset of protein conformations is used in docking, for example, following clustering of the target's MD configurations. HPC involves the use of specialized, large supercomputing systems to perform large calculations that are parallelized over many compute nodes, each consisting of dozens of cores. Traditionally, the use of a high-speed interconnect allows rapid communication between separate compute units and clever parallelization schemes to enable rapid calculations on problems too large to fit on a single compute unit. These schemes have historically involved specialized programs that focus efforts to optimize communication overlap. The use of graphics processing units (GPUs) has helped to accelerate many types of calculations. The Summit supercomputer, housed at the Oak Ridge Leadership Computing Facility (OLCF), is currently the fastest supercomputer in the United States. Summit is an IBM AC922 system consisting of 4608 large nodes, each with six NVIDIA Volta V100 GPUs per node. Each node also contains two POWER9 CPU sockets for a total of 42 cores per node. The GROMACS MD program²⁷⁻²⁹ is able to make use of all aspects of the Summit supercomputer's HPC utilities, including the GPUs and the interconnect, providing for both strong and weak scaling, which dramatically decreases time per MD step and increases the size of the system that can be simulated efficiently. The temperature replica exchange molecular dynamics (T-REMD) routine³⁰⁻³³ which was chosen here for the MD calculations (see below) uses the interconnect not only to allow for parallelization of a single simulated biomolecule, but also to communicate between separate replicas of the system, each carried out at a different temperature, and performs exchanges between replicas to accelerate the conformational sampling of the biomolecule³⁴.

Protein-ligand docking has hitherto not been considered a traditional HPC task, as each docking calculation is short and does not require multiple nodes to complete. In fact, many docking programs can run on a single CPU core. However, the number of cores on a supercomputer or cluster can provide a resource to perform many simultaneous docking calculations that greatly decrease the time-to-solution for screening a large dataset of ligands. Cloud and distributed computing resources also provide this type of completely parallel solution for high-throughput screening^{35, 36}. The use of GPUs has recently been made possible for the widely-used program AutoDock^{37, 38} resulting in the program AutoDock-GPU, which provides up to 50X speedup over AutoDock4 (available at <https://github.com/ccsb-scripts/AutoDock-GPU>)³⁹⁻⁴¹. Thus, the use of leadership HPC facilities for ensemble docking can provide the ability to screen billions of ligands to a full set of conformations generated with HPC-based MD simulations. Quantum mechanical refinement of classical docking ranking based on fragment molecular orbital (FMO) techniques also naturally benefits greatly from massively parallel supercomputer capabilities⁴².

The need for rapid time-to-solution in drug discovery has become accentuated by the Covid-19 pandemic. Combating viral illnesses such as Covid-19 encompasses two main modalities: vaccination for prevention and drug therapy for those already infected. Unfortunately, vaccine development can be a long process due to the need to ensure immunity while also determining effective booster schedules to reinforce immunity and safety⁴³. Regarding the latter, there can be a paradoxical enhancement of infection after vaccine administration wherein viruses hijack the host's immune system and actually utilize antibodies to *facilitate* entry into host cells⁴⁴. Partly for these reasons, there is an impetus to develop in parallel drugs that target the virus directly. These are typically small molecules that bind viral proteins, abolishing their function and, in turn, inhibiting viral replication, although drugs targeting human proteins can also be considered.

Given the urgency of the pandemic, repurposing established antivirals appears to be an expedient and reasonable approach. Among the present candidates, the antiviral remdesivir, a nucleoside analog that acts by interfering with RNA synthesis, is a leading candidate for repurposing, as it is active against coronaviruses related to SARS-CoV-2, including the coronavirus responsible for MERS⁴⁵⁻⁴⁹. Recent publications have indicated that remdesivir shortens recovery from Covid-19 in hospitalized patients.⁵⁰ Other promising candidates, such as dexamethasone, may modulate the host response⁵¹.

When surveying the existing clinical trials for Covid-19, one is struck by the number of trials that are not based on knowledge of the drug interacting with a known target. There are several examples. As one illustration, baloxavir is a specific inhibitor of the cap-snatching endonuclease of influenza virus, which is a member of the PD-(D/E)xK two-metal nuclease superfamily⁵². Coronaviruses have an endonuclease but of a completely different fold (NendoU) and different active site residues⁵³. NendoU also oligomerizes into a hexamer. Although it would seem unlikely that baloxavir would bind to NendoU, it nevertheless is in clinical trials. Similarly, lopinavir and ritonavir are also undergoing testing, even though they target proteases of the unrelated HIV⁴⁸. Although in principle, enzyme active sites with similar chemical functions may bind similar ligands, the steric and physicochemical substrates of drug-protein binding are nuanced. The frequent trial of drugs specific for targets known to be absent in SARS-CoV-2 seems to us to be symptomatic of a lack of precision in combating this pandemic. For further information on other ongoing trials on small molecule drugs, biologics, passive immunization with antibodies, and vaccines, the reader is referred to the comprehensive review by Liu *et al.*⁵⁴. The aforementioned challenges, taken together, demonstrate that there exists an unequivocal need for *de novo* drug discovery campaigns as well as repurposing studies.

There are a number of events in the SARS-CoV-2 viral replication cycle⁵⁵ to target for antiviral therapeutic development; from viral entry to membrane fusion, travel to the host endoplasmic reticulum where translation of the viral genome occurs, to formation of the viral replication complex and formation from host membranes of double-membrane vesicles (DMVs)^{56,57}, the passage of the replicon through the Golgi and the release of the virion from the cell. Each of these steps involves key viral proteins and occurs in a different compartment of the host cell. For example, the binding of the virus to the ACE2 receptor involves the receptor-binding domain (RBD) of the virus S (spike) protein, pre-fusion cleavage involves the binding of host TMPRSS and furin proteases to the S1/S2 dibasic domain⁵⁸ formation of the replication complex and the DMVs involves the non-structural proteins (NSPs), and the N protein is required for packaging of the viral genome into the newly assembled virion⁵⁹. The replication complex is made up of 15 mature NSPs, which are encoded by *orf1ab* and *orf1a* genes as the pp1ab and pp1a polyproteins⁶⁰. Currently, many efforts are targeting the main protease, MPro⁶¹, which is required for cleavage of the large viral polypeptide into its functional proteins, the RNA-dependent RNA polymerase (RdRp)⁶² responsible for the production of new viral RNA, and some efforts target prevention of S cleavage⁶³. In addition, viral proteins also function to impede the host's defense mechanisms: both proteases have been shown to inhibit the human immune response by interacting with immune proteins in SARS-CoV⁶⁴. It is, therefore, important to understand regions on these proteins that act as binding sites for both substrates (as in the case of the proteases) and for protein-protein interaction.

There are 16 NSPs that are involved in SARS-CoV-2 viral replication and the circumvention of the immune system. The NSPs vary in size and complexity, and the function of some of them is not currently known. For instance, NSP3 is a large multi-domain protein⁶⁵ composed of around 1900 amino acids that form multiple domains, making it a particularly difficult target for structural biology studies. To make structural studies more tractable, individual domains of NSP3, such as the papain-like protease domain and the ADP ribose phosphatase domain, can be overexpressed in bacteria and then purified and crystallized for X-ray diffraction studies.

Besides importance to the viral replication cycle⁶⁶, there are other factors to consider when choosing targets for *in-silico* drug discovery. One important concept is the idea of druggability⁶⁷, i.e., the likelihood that a drug will be able to change the behavior of the target in a way that brings a therapeutic response. Traditionally druggability has been based on previous responses to drugs of proteins in similar families or with similar structures⁶⁶. Newer metrics include analyses of the structure of the target's active site and the desolvation that occurs when the drug binds to the target⁶⁸. For antivirals, especially for coronaviruses such as SARS, MERS, and SARS-CoV-2 viruses, certain targets such as the RdRp have been considered to be more druggable than others⁶⁹.

Another consideration for the medium to longer-term is the ability of the pathogen's proteome to mutate and cause drug resistance^{70, 71}. It is possible that less conserved proteins in the proteome may withstand more mutations and thus could confer increased viral resistance to a drug than more conserved proteins⁷²⁻⁷⁴. For more rapidly mutating viruses, combination drugs may be used in an attempt to counteract drug resistance that occurs as the virus mutates in response to drugs^{70, 71}. Incorporating information about mutation frequencies when considering drug targets can help to prevent resistance-- it is possible that ancestrally conserved proteins across taxa may be less likely to mutate and cause resistance. Therefore, sequence analyses can be helpful when deciding which viral proteins should be targeted by *in-silico* efforts or experimental screenings that make use of *in-silico* filtering. An analysis of the current mutations in the SARS-CoV-2 virus and their frequencies can help to elucidate emerging trends that may affect the propensity for the resistance of the various targets and help to focus choices for lead optimization.

In previous work, very early in the pandemic, we combined restrained temperature replica-exchange molecular dynamics (restrained T-REMD) simulations with virtual high-throughput screening in a supercomputer-based ensemble docking campaign to identify well-characterized drugs, metabolites, or natural products that bind to either the S-protein: ACE2 receptor interface or the RBD of the S-protein²⁵. From this ensemble docking campaign, we provided a ranking of the predicted binding affinities of over 8000 drugs, metabolites, and natural products (and their isomers) with regards to the SARS-CoV-2 S-protein and the S-protein: ACE2 receptor complex. The ranked list has been incorporated into experimental testing using a high throughput screen that was implemented in the SARS CoV outbreak, and new compounds will be added as discovered. Two of the top compounds in this screen are now in clinical trials (Trial Numbers: NCT04357990 and NCT04377789).

In this work, we describe our efforts establishing a supercomputing-based pipeline for ensemble docking and preliminary results on its application to discovering therapeutics that target viral proteins of SARS-CoV-2. The pipeline and results presented here represent both the continuation of the above work and our contribution to date to the work of the USA HPC Covid-19 Consortium that was created on March 29th, 2020. We describe the choice of 8 targets and the preparation of protein models from experimental data. We report on T-REMD simulations performed for the targets totaling about half a millisecond of simulation time. We have docked repurposing databases to ten configurations of each protein simulated using the popular docking program Autodock Vina. We also describe efforts deploying Autodock-GPU³⁹⁻⁴¹ at scale on Summit that demonstrate the docking of over a billion compounds in 24 hours with full structural optimization of the ligand. Future developments involving the use of AI and quantum chemistry in rescoring

and clustering are also outlined. The pipeline described here can also be used in future work to target human proteins^{75, 76} known to interact with viral proteins, or in disease-causing responses in Covid-19 and more generally in computational structure-based drug discovery.

Methods

Computational methods in drug discovery narrow a vast chemical search space to a tractable set of compounds suitable for experimental testing. Experimental work can involve a variety of tests, including live virus testing as well as target engagement studies, and will not be considered further here. Rather, we discuss the procedures of structural modeling, MD simulation, and docking.

a. Choice of target proteins and generation of structural models from experimental work

Multiple groups have been using structural biology techniques, including X-ray crystallography, small-angle scattering (SAS), and cryogenic electron microscopy (cryoEM), to investigate the structure of proteins and protein complexes from SARS-CoV-2⁷⁷⁻⁸⁰. However, obtaining a structure from the Protein Data Bank (PDB) or perhaps a revised model from another resource is only the starting point. Often structures are missing residues that were not resolved, and a determination must be made whether and how to model them. Also, as structural models are rapidly being released to aid in the fight against COVID-19, the potential inclusion of a few structural errors is an unfortunate reality. In particular, the identification of metal cations in protein structures requires careful thought and examination of its local coordination environment.

Even with perfectly assigned and complete experimental structures, it may not be enough to consider viral protein targets as chemically invariant structures for modeling and binding calculations. Large differences in pH in various cell compartments as the virus travels through the host cell⁸¹⁻⁸³ can qualitatively change the protein's structure and function. Differences in pH also affect the protonation states of the proteins and the small molecules being tested as drugs, altering drug binding preferences. Finally, the oligomerization states of the target proteins are important to consider as the interactions between protein monomers may influence the shapes of the active sites. Another important factor to consider when performing *in silico* screens using ensemble docking is the ability to construct a useful model of a particular protein for MD simulation. For instance, certain metal-containing regions of a protein may not have an existing classical mechanics model (force field parameters), or existing models may be inadequate. In addition, highly charged, disordered, and ion-dependent biomolecules have been known to have less accurate force fields and may perform poorly in an MD simulation⁸⁴⁻⁸⁷.

Proteins chosen for ensemble docking in this study were those that had a crystal structure available with a reasonable resolution, were amenable to accurate simulation with classical MD force fields, and were also known to be important for viral pathogenicity based on either recent studies or those on SARS-CoV. The 23 systems studied comprise nine protein domains. Two of these, RBD of the S (spike) protein and the N-terminal region of the N (nucleocapsid) protein, are domains in structural proteins found attached/within the virion⁸⁸. The remaining seven domains come from non-structural proteins (NSPs) 3, 5, 9, 10, 15, and 16, which form the replication complex and are involved in a number of key tasks leading to the creation of new virus particles⁸⁹. Two domains come from NSP3⁶⁵, the ADP ribose phosphatase (ADRP, also known as macro- or "X") domain, and the papain-like protease domain (PLPro)⁹⁰. The ADRP seems to be involved in ADP-ribosylation, which is used in cell signaling and thus may act to inhibit the host immune response⁹¹. PLPro cleaves regions of the polyprotein to release non-structural proteins and also is involved in the mechanisms the virus uses to counteract the host immune response, for instance by interaction with host

immune proteins^{89, 92} Nsp 5 is the main protease (MPro), which self-cleaves and also cleaves other regions along the polyprotein, releasing essential proteins to perform their tasks in the assembly of the replication complex, and is also involved in interacting with and preventing the action of host immune proteins^{92, 93}. The exact function of NSP 9 is unclear, but has been found in SARS-CoV to be required for viral replication and has been shown to bind to RNA oligonucleotides⁹⁴, Nsp 15 is an endoribonuclease specific for uridine whose exact function is also unknown, but has been implicated in interfering with host immune response both through direct interaction and by cleaving viral RNA to prevent detection by the host⁹⁵. NSP 16 is thought to be a methyltransferase that requires NSP 10 as a co-factor, and acts to disguise viral mRNA from the host immune response by adding a methylation onto the RNA cap which host cells use to mark RNA as belonging to “self” versus “pathogen.”⁹⁶⁻⁹⁸.

The explosion in research and literature fueled by the Covid-19 pandemic, together with the need for searching through related literature on other coronaviruses, has created a challenge for researchers needing to understand the structural details and cellular contexts of the SARS-CoV-2 proteins. To help navigate this challenging landscape, we have been developing new tools based on natural language processing for enabling a more robust search for specific questions required for our modeling, simulation, and ligand docking work⁹⁹⁻¹⁰¹ featuring targeted filtering and exploiting external resources (e.g., Wikidata, ChEBI, PubChem) to expand our search capability. For example, after we generate a set of related keywords, the service will screen for the terms referring to a chemical substance and fetch the chemical information (e.g., SMILES string) from the PubChem automatically. In addition, using this keyword search enables the ontologies (e.g., Wikidata, ChEBI) to be used to link related chemicals and their properties for document annotations in query results. The main data resource of the system is a collection of scientific papers, which are collated from major publications. The full-text article access and download from the publishers’ archives are performed under the publishers' agreements, and the internal article corpus in our system is updated on a weekly basis.

To provide a diverse survey of the conformational ensembles of the SARS CoV-2 viral proteome, we performed T-REMD simulations of 23 different model systems listed in Table 1. An additional supplementary table (SI Table I) is also provided, which summarizes the PDB entry simulated, complete protonation state choices (where applicable), and the number of replicas used for the T-REMD.

Table 1. Model Systems Simulated

Protein/ System Notes		
S (Spike) Protein Receptor Binding Domain (RBD) / “Apo”	S Protein RBD / Complexed with ACE2	MPro / monomer, CHARMM-GUI default protonation
MPro / dimer, CHARMM-GUI default protonation	MPro / dimer, ‘charged’ protonation variant	MPro monomer/ HIE41 protonation variant
MPro dimer / HIE protonation variant	MPro monomer / HID41 protonation variant	MPro dimer / HID41 protonation variant
NSP15 (endoribonuclease) / hexamer	NSP15 (Endoribonuclease) / monomer	NSP10:NSP16 Complex (Methyltransferase)

NSP10 / monomer	NSP16 / monomer	N (nucleocapsid) N-terminus phosphoprotein / monomer
N (nucleocapsid) N-terminus phosphoprotein / tetramer	N (nucleocapsid) N-terminus phosphoprotein / tetramer complexed with Zn	N (nucleocapsid) N-terminus phosphoprotein / monomer alternate crystal structure
NSP9 / monomer	NSP9 / dimer	NSP3 ADP ribose phosphatase / asymmetric unit
PLPro / monomer 'charged' protonation variant	PLPro / monomer 'neutral' variant	

b. Simulation Model Preparation

To engage in the use of MD for the rapid generation of conformational ensembles for drug discovery one requires that the input for MD be generated in a semiautomated fashion by which the atomic coordinates, obtained from experimental or *in silico* protein structure prediction methods, can be quickly processed into MD input files. To facilitate this semiautomated approach, CHARMM-GUI was used for most model building¹⁰². The general system building method used here involves the direct retrieval of structures from the PDB and processing to model missing residues, assign protonation states, add disulfide bonds (where noted in the PDB annotation), add glycosylation (where resolved in the crystal structures of the S-protein receptor binding domain and ACE2), neutralize the charge of the system (using Na⁺ and Cl⁻ ions), and solvate (with TIP3P water). Many proteins have coordinated ions that serve structural roles, such as the Zn²⁺ cations in Nsp10, or catalytic roles. Thus, the treatment of Zn-complexes in fixed-charged classical MD force-fields is a challenge, and for some systems, it may result in the failure to maintain Zn-protein coordination^{103 104 105}, and when found necessary (as noted below and also summarized in SI Table I) an explicit bond representation was used. All of these considerations mandate an abundance of care when preparing a biologically accurate model for simulation. Below we discuss considerations taken into account when modeling some of the proteins simulated.

S (Spike) Protein

Presently available structures of the S protein have nine gaps totaling approximately 150 residues, in addition to over 20 and over 100 missing residues at the N- and C-termini, respectively. Current models also lack post-translational modifications, including glycosylation and formation of disulfide bonds. The S protein is heavily glycosylated, with roughly 20% of its mass in glycan chains, yet at most, a few mono or disaccharides are present in the structure.

In our preliminary study²⁵, we made use of a homology model of the entire spike with restraints applied such that only the human ACE2-Spike interface was unconstrained. Here, using crystal structures of the ACE2-S protein complex, simulations of the receptor-binding domain of the S Protein (Spike) both in complex with the human ACE2 receptor and on its own (referred to within the text as the “Apo” RBD) were performed. The viral spike receptor-binding domain was chosen to provide insight into the details of the initial viral-host recognition process. Glycosylation resolved from crystallographic imaging was used, and an annotated disulfide bond was also included.

Main protease.

The main protease, MPro, is an attractive target for the development of antiviral drugs. There is compelling evidence that the enzymatically active species is the dimeric assembly of MPro. A dimer is observed in most crystal structures of CoV MPro, as well as in solution at sufficiently high concentrations. In addition, a linear increase in the enzyme activity at increasing concentration suggests catalytic incompetence of the monomer¹⁰⁶. Therefore, the full dimer was considered in the present MD simulations for MPro, using as starting coordinates the apo-homodimer in the crystal structures 6Y2E and 6WQF^{107, 108}.

The crystal structures show that SARS-CoV-2 MPro, similarly to other MPro's^{93, 107, 109, 110}, is composed of three domains: Domains I (residues 8–101) and II (residues 102–184) are arranged in an antiparallel β -barrel structure, and domain III (residues 201–303) contains five α -helices arranged in a globular cluster. Domain III is a specific feature of CoV MPro proteins and was suggested to be essential in the proteolytic activity by keeping domain II and the long loop connecting domains II and III (residues 185–200) in the proper orientation, and/or by orienting the N-terminal residues that are essential for the dimerization¹⁰⁹. Dimerization occurs through interactions between the helical domains of the two monomers and through hydrogen bonding interactions between the N-terminal residues of one monomer and key residues in the other monomer. In particular, the salt bridge between the N-terminal Ser1 of one monomer and Glu166 of the other monomer has been suggested to be essential to maintaining the catalytically competent conformation^{109, 111}. The substrate-binding site is located in a cleft between domains I and II and contains a highly conserved catalytic dyad formed by Cys145 and His41. Comparison among the two apo crystal structures and the crystal structure obtained in the presence of an inhibitor reveals⁹³ only minor structural differences in the position of a few side-chains and no relevant changes in the substrate-binding site, except for the rotation of the side chain of Met165, which is in the proximity of His41.

Although the catalytic mechanism is not fully understood, there is a general agreement in considering that the proteolytic activity of CoV MPros is initiated by activation of the enzyme through a proton transfer reaction in the catalytic dyad, leading to a charge-separated state with a highly reactive thiolate. It has also been suggested that such a proton transfer reaction is induced by the presence of the native substrate^{112, 113}. Therefore, in the present MD simulations of the apoenzyme, Cys145 and His41 were simulated in their neutral state, with His41 protonated at N δ (i.e., HSD). This choice is based on the observation that the His41-N ϵ appears to be the best candidate as proton acceptor from Cys145 because in the crystal structures the His41-N ϵ is closer than the His41-N δ to the Cys145-S and the His41-N δ is already involved in a hydrogen bond to a highly conserved water molecule, which is considered the third element of the catalytic site. A recent QM/MM study¹¹⁴ also supports this proton transfer mode and the role of water in catalysis. However, the ϵ -nitrogen protonation state for His41 (HIE) cannot be decisively ruled out, and MD simulations were performed also considering this alternative, although less probable, protonation state.

The protonation states of two additional His residues, namely His163 and His172, were also highlighted as being crucial for the enzymatic activity of CoV MPro. In particular, the doubly protonated (cationic) state of His163 at pH 6.0 was suggested to modulate relevant conformational variations involving Glu166, Phe140, and His172, leading to a catalytically inactive conformation¹⁰⁹. At higher pH values, both His163 and His172 should be uncharged, and, on the basis of the hydrogen-bonding pattern that can be inferred from the crystal structures, the HSE protonation state was used for both His163 and His172 in the present MD simulations. All other His residues were also simulated in their neutral state, assigning the N δ or N ϵ protonation state on the basis of their chemical environment and hydrogen-bonding patterns. The selected protonation states are as follows: HSD64, HSD80, HSE164, HSE246.

PLPro. For PLPro (PDB entry 6W9C), the Zn cation was coordinated to C189, C192, C224, and C226. Similar to MPro, the protonation states of the His residues in PLPro were not readily available. These were obtained with the use of the PropKa 3.0 server assuming pH ~5, corresponding to its presumed cellular (lysosome) environment¹¹⁵, with 6W9C being assigned to pH 5 based on its physiological role in acidic

environments. TopoGromacs¹¹⁶ was used to convert the system and associated force field to GROMACS format, allowing Zn-coordinating proteins to be incorporated into the overall temperature-replica exchange workflow.

NSP15. Large (His) tags were added during crystallization for NSP15. Prior to simulating the monomeric and hexameric forms, these His tags were removed from NSP15 using MOE2019 and subjected to a “quick prep” with the *prepare protein* module of MOE to resolve potential issues in the resulting structure. The resulting truncated PDBs were then uploaded to CHARMM-GUI for neutralization and solvation.

NSP10. For NSP10, both in its monomeric and a complexed form with one Zn cation liganded to C4370, C4373, C4381, and C4383, while the other bound Zn was liganded to C4327, C4330, H4336, and C4343. For 6VYO, H59 and H145 were liganded.

c. Molecular Mechanics and Molecular Dynamics System Preparation (Force Fields, Counter-Ions, Energy Minimization, and Equilibration)

All simulations were performed using the GROMACS^{27, 117} software suite, and the CHARMM36m force field¹¹⁸, which is the default of choice using the CHARMM-GUI. For all systems, the protein was solvated in water-boxes with edge-distances of 1nm, and only neutralizing Na⁺ and Cl⁻ ions were used. Short-range interactions were treated with a smooth force-switch cutoff of 1.2 nm, and long-range electrostatics were treated using the particle-mesh Ewald (PME) formalism, as implemented in GROMACS¹¹⁹. To facilitate the use of a 2-fs MD timestep, all covalent bonds to hydrogen were restrained with the LINCS algorithm¹²⁰ in all simulations. Following system preparation, all solvated models generated were subject to steepest-descent energy minimization with a stopping condition of either reaching the force-convergence criteria of 1000 kJ mol⁻¹ nm⁻¹ or a maximum of 5000 iterations. Energy minimization was performed primarily to remove potential clashes between the solvent, ions, and the protein (or protein complex) of interest. Post-clash removal minimization, short (250 ps) NPT relaxation simulations (with default positions restrains generated from CHARMM-GUI) were performed to relax the simulation box dimensions for each replica (at different temperatures) independently (*see T-REMD Protocol*). For these relaxation simulations, the Berendsen baro/thermostats¹²¹ (as implemented in GROMACS) and an integration time step of 1 fs were used.

d. T-REMD Protocol

MD simulations provide a means to study the conformational dynamics of proteins. However, frequently MD becomes ‘trapped,’ resulting in the need for many long simulations to effectively sample a protein’s conformational landscape¹²². To overcome this sampling challenge, enhanced sampling techniques can be used. For the present work, temperature replica-exchange molecular dynamics (T-REMD) was employed, whereby multiple copies of a target system are simulated simultaneously with each copy (replica) at a different temperature, with periodic coordinate swapping (performed in such a manner as that preserves detailed balance) between the copies³¹⁻³³. By running at multiple temperatures, with exchanges, the dynamics of the system avoids ‘kinetic traps’ and provides a robust sampling of the protein free-energy landscape, and thus the protein conformational diversity¹²³. T-REMD was chosen for several reasons:

- 1) it guarantees an increase in sampling efficiency over straightforward MD¹²⁴,
- 2) it does not require the assignment of reaction coordinates (or collective-variables) *a priori* to accelerate conformational sampling
- 3) it does not require direct modification to the system Hamiltonian¹²⁵.

T-REMD simulations for each system were performed with the GROMACS simulation suite. A limited temperature range of 310 K to ~350 K was chosen to maintain physiological configuration space. For each protein system, the number of replicas and temperatures for each replica was chosen using the temperature predictor server by Patriksson and van der Spoel¹²⁶ with a target exchange probability set at 0.2 though the actual exchange was found to be ~0.3 for all systems. All simulations were performed for a total of 750 ns per replicate.

After relaxation, production T-REMD simulations were performed with a frame saving rate of 10 ps and an integration time step of 2 fs. Production simulations were performed, similar to the relaxation simulations, in the NPT ensemble. Unlike the relaxation simulations, the V-rescale (Bussi) thermostat¹²⁷ and the Parrinello-Rahman barostat^{128, 129}, were used. Regardless of the temperature window, the target pressure for each replicate was set to 1 bar.

e. Trajectory Analysis

For all systems, the measures of the gyration tensor (from which shape anisotropies are derived), solvent-accessible surface area (SASA), and pairwise simulation frame versus simulation frame RMSD matrices, and RMSD based clusters, were obtained using a combination in-house VMD¹³⁰ scripts, NumPy, and SciPy¹³¹. The RMSD clustering specifically only considered the lowest temperature replica (310 K), and rapidly generated the pairwise RMSD matrices using the QCP algorithm¹³². Clustering was performed using hierarchical clustering with a complete linkage method, as implemented within SciPy. For generality in evaluating structural diversity, clustering was initially performed based on the RMSD of all heavy protein atoms, and where additional diversity of active sites was of interest for subsequent docking, a second round of clustering was performed based on binding site residues and protein-protein interfaces. VMD atom selections for the docking specific clustering are summarized in SI Table II.

Although T-REMD is an efficient simulation method, and the 310K data do correspond to a formal statistical mechanical ensemble generated at this temperature, as with other enhanced sampling methods the risk is always present that the enhancement of the sampling takes the system to regions of configurational space beyond that that would be significantly sampled by the protein physiologically; for example, to partially or wholly unfolded states. We, therefore, take care to identify these and to not perform docking screens on such configurations.

f. Docking

Two different docking databases were used.

- 1) A **smaller database** of potential ligands was built merging together the content of the SWEETLEADS^{133, 134} repurposing database SuperDRUG2^{135, 136}, and the NCI-diversity database¹³⁷, yielding 13,757 unique compounds. This database has been ensemble docked to all systems, as listed in Table II, with noted targeted binding sites. This database was docked using local HPC clusters using Autodock Vina.
- 2) Supercomputing docking runs were performed involving **billion-plus** compound screens of the Enamine database using an accelerated version of Autodock: Autodock-GPU. To date, these runs have been performed on two crystal structures of MPro.

f.1 Smaller database docking

Data and Protocols

Docking to the target structures obtained from the MD simulations (as listed in Table I) was performed on various HPC clusters using Vina MPI¹³⁸ and MOE. Two sets of structures were used in the ensemble docking. In the first series of docking calculations, only the first 100ns of the T-REMD trajectories were used, and the results of the docking simulations were passed on to collaborators for experimental testing. In the second series of docking, as the MD trajectories were expanded beyond their initial first 100ns, the clustering was performed on the entire 750ns trajectories, as described in the results section below.

For the VinaMPI¹³⁸ calculations, the “Exhaustiveness” parameter was kept at its default value of 10. Databases of potential ligands were built merging together the content of the SWEETLEADS^{133, 134}, SuperDRUG2^{135, 136}, and NCI-diversity databases¹³⁷, yielding 13,757 unique compounds.

Using the program MOE, compounds with more than 49 rotatable bonds were deleted from the database, and only one stereoisomer was included for each compound. Very low molecular weight (<58) compounds (single atoms, ions, very small functional groups). The resulting database included ~9K unique molecules. The compounds were protonated at pH 7 and energy-minimized using the MOE software to obtain low energy 3D structures. The compounds were saved on disk in sdf format and then converted to PDBQT format using AutoDock Tools^{139, 140}.

The ligands were docked to 10 clusters per receptor, each cluster corresponding to a different configuration of the binding pocket. The clusters corresponding to the first 100ns of the MD simulations have been uploaded on the publicly available structure repository <https://cmbcovid19.flywheelsites.com/data/> additional data from the complete 750ns T-REMD simulations is forthcoming. The residues used to determine the clusters fall into one of three categories: the protein active site, residues at the protein-protein interfaces (for complexes), and all the protein non-hydrogen atoms. Tables II and SI Table II list the receptors and binding sites we have screened so far.

Table II. List of proteins and binding sites used for ‘smaller database’ docking. PPI refers to a protein-protein interface. In some cases, FTMap was used to identify potential binding sites (see SI Table II)

Receptor / Binding Site	Receptor / Binding Site
MPro monomer / catalytic pocket	NSP15 monomer / catalytic pocket
MPro dimer / PPI	NSP15 dimer / PPI
NSP9 dimer / FTMap sites	NSP10 monomer / PPI to NSP16
Nucleocapsid phosphoprotein / RNA binding site	NSP16 monomer / PPI to NSP10
Nucleocapsid phosphoprotein / PPI	NSP10:NSP16 / PPI
Nucleocapsid tetramer / FTMap sites	NSP3 dimer / active site

Binding Sites for Docking

In general, we have two classes of potential binding sites: 1) catalytic pocket or substrate-binding site and 2) PPI. The first aims at identifying potential competitive inhibitors of the viral proteins, and the second aims at finding compounds potentially disrupting a viral protein-protein complex. Binding site definition requires manual intervention and cannot be easily automated. Examples of definitions are listed below for three viral proteins.

- a) In the main protease dimer (PDB: 6WQF), the docking box contains catalytic sites of chain A and PPI residues. The docking box was constructed to align with the peptide-binding groove on either side of the catalytic dyad of chain A, which extends outward to include the S3, S2, S1, S1', and S2' catalytic pockets.
- b) In the NSP10-NSP16 complex (PDB: 6W4H), the *S*-adenosyl methionine (SAM) binding site Asp6928 in NSP16 was considered⁹⁷. In addition, PPI residues such as Tyr4349, Val4295 to Leu4298 in NSP10, and Gln6885 in NSP16 were included. Tyr4349 and Gln6885 interact with each other in SARS-CoV virus⁹⁷, and Val4295 to Leu4298 are hot spot residues in the SARS-CoV-2 virus computationally predicted using the crystal structure along with the KFC2 method¹⁴¹, which is based on a machine learning predictive model (<https://kfc.mitchell-lab.org>). Hot spot residues are the fraction of PPI residues that account significantly for the overall protein-protein binding affinity, and they are typically determined experimentally using alanine scanning mutagenesis¹⁴².
- c) In the N-terminal domain of nucleocapsid protein tetramer (PDB: 6VYO), three critical RNA-binding residues on the beta-sheet core were included in docking: Arg88, Arg92, and Arg107^{80, 143, 144}.

f.2 Billion-compound supercomputer docking with Enamine Real database

A major aim of this exercise was to see whether it would be possible to dock a billion compounds with full ligand optimization on the OLCF Summit supercomputer in 24 hours of wall-clock time. To perform efficient ensemble docking, we modified AutoDock-GPU^{39, 41}, to enable it to run at peak efficiency on the Summit system. For compatibility, OpenCL kernels were re-written in CUDA, and file input and output were streamlined to enable it to keep up with the GPU's speed. These modifications, together with the size of the Summit supercomputer, indeed allow over 1 billion compounds to be docked within 24 hours. This capability will enable giga-compound docking for a number of proteins in the viral proteome and beyond.

We performed initial docking tests using this framework on NSP15 (NendoU) and the main protease (MPro). For NSP15 we used a 9,000 compound dataset composed of the SWEETLEADS¹³³ database with additional ligands, and also a trimmed version of this dataset containing only ligands containing less than 11 rotatable bonds, consisting of about 5,000 ligands. For tests with MPro, we used a 90,000 ligand subset of the Enamine REAL database¹⁴⁵. All ligands were prepared with AutoDockTools^{140, 146}, and the receptor grids were generated with the program *autogrid* with a grid spacing of 0.375 Å. We tested a set of search box sizes: 40, 25, 20 and 15 Å³, and different settings for the number of runs, *nruns*, which defines how many separate instances of the genetic algorithm are executed. For the trimmed dataset, we also performed docking with AutoDock Vina with exhaustiveness of 10 to compare results. These results provided us with the confidence to dock over 1 billion compounds from the Enamine real database to two different MPro

crystal structures, 5R84 and 6WQF¹⁰⁸, with a search space 25 Å large on each side, centered on the active site. The analysis of this dataset is ongoing. Due to the documented inaccuracies of force field-based scoring functions in the task of screening and affinity prediction of compounds,¹⁴⁷ rescoring of at least 1 percent of the billion compounds is being performed using the accurate, yet highly computationally efficient machine learning-based rescoring method known as RF-Score-3¹⁴⁸. Also, at least 50% of those compounds re-scored with RF-Score-3 will be further filtered using recently developed rescoring described below in [Future Directions and Preliminary Results from New Methodologies](#), subsection *Protein-ligand rescoring using machine learning*.

Sequence analysis and mutational entropy calculations

We performed an analysis of available sequences of the SAR-CoV-2 virus to look for numbers of mutations and map these locations on the proteins we were using as drug discovery targets. All complete, high-coverage genomes labeled as human host SARS-CoV-2 were downloaded from GISIAID^{149, 150} on May 5, 2020, yielding a total of 16,252 genomes. Sequences were filtered to remove any genomes with greater than 3% ambiguous (N) nucleotides or were less than 29,000 nucleotides in length, resulting in 14,284 genomes. Multiple sequence alignment of the 14,284 genomes was performed using MAFFT¹⁵¹ v.7.464 with the --addfragments method using NC_045512.2 (EPI_ISL_402125) as the reference genome and removing insertions relative to the reference. Mature protein-coding sequences for each protein were extracted from the alignment using coordinates from the reference genome and translated using FAST¹⁵² v1.6, with protein sequences containing internal stop codons discarded from further analyses. Shannon entropy¹⁵³ was calculated for every column of each protein alignment using a custom script, disregarding ambiguous and gap characters using a custom script. Additionally, the frequency and types of substitutions with respect to the reference were recorded. For visualization of the mutation entropy per residue of the proteins studied in this paper, entropy values were color-coded in protein PDB structures. Known SARS-CoV and SARS-CoV-2 structures were downloaded from the Protein Data Bank, their sequences were aligned with the SARS-CoV-2 reference genome (NC_045512.2) using BLASTP, and the calculated entropy of the sequences was embedded in the PDB file in the place of the B factor column using a custom Python script.

Preliminary QM Refinement Protocol

Along with ML-based approaches, quantum mechanics (QM)-based refinement of classical docking results is being developed here as a tool to narrow down the list of promising inhibitor candidates¹⁵⁴. Until recently, the inclusion of QM electronic structure in high-throughput drug screening was deemed computationally intractable due to the enormous computational resources required even for density functional theory (DFT) calculations. The poor scaling of most quantum chemical methods further exacerbates the situation. A viable emergent alternative is the recently developed linear-scaling version of an approximate, yet remarkably accurate DFT method called “fragment molecular orbital density-functional tight-binding (FMO-DFTB)”¹⁵⁵. This method is implemented in the widely-utilized GAMESS quantum chemistry code¹⁵⁶. We here report preliminary calculations of FMO-DFTB with the so-called polarizable continuum model (PCM) of the solvent¹⁵⁷ for quantum mechanics-based evaluation of potential COVID-19 spike protein inhibitor drugs identified by re-clustering and re-docking to an extended simulation of the S protein, similar to the initial work by Smith & Smith²⁵. In addition, to improve the accuracy in describing non-covalent interactions, the D3 dispersion correction was employed. To obtain the refined binding energy of a given candidate, its unbound geometry, the unbound protein, and its corresponding complex were optimized using FMO-DFTB/PCM. While the unbound ligands were completely optimized, only selective residues in the binding pocket of the unbound protein, and in the protein-ligand complexes were locally optimized. The QM-refined binding energy is defined here as the difference between FMO-DFTB/PCM total energy of the complex and the sum of the total energy of unbound protein plus total energy of unbound ligand. In

preliminary work, the QM-refinement was carried out for the Vina top-10 best binders of each spike protein cluster. In total, 15 spike protein clusters were investigated, and the binding energies of 150 protein-ligands complexes were refined.

Results

We present here preliminary results obtained for members of the SARS-CoV-2 proteome. Naturally, ongoing refinements of the results are continually being undertaken, and the results are incomplete. However, they give a snapshot report on the state of delivery of the pipeline. At the moment of submission, 23 T-REMD simulations have been performed on nine members of the proteome, in various oligomerization and protonation states, for a total of 0.612 ms of MD aggregated over all replicas and $\sim 17.25\mu\text{s}$ aggregated overall lowest temperature windows. At present $\sim 2.07\text{M}$ physical docking calculations have been performed with the smaller database and on Summit 2.4 billion docking calculations with the Enamine REAL database. The preliminary results presented are general trends observed in the MD and docking runs and do not describe details of the candidate compounds or dynamical properties of individual proteins, which will be reserved for future publications. All results of MD and docking are available at the website <https://coronavirus-hpc.ornl.gov>

a. T-REMD scaling performance

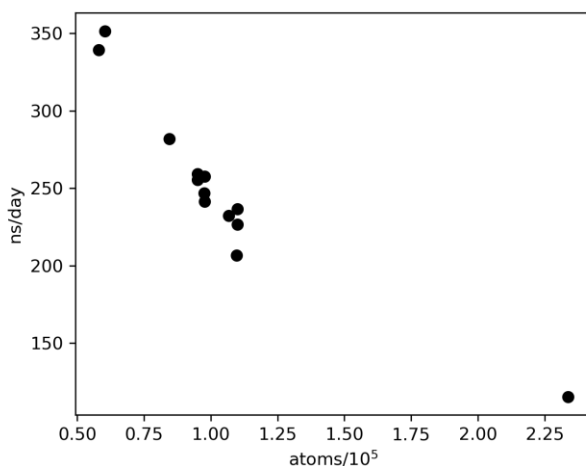


Figure 1. Simulation throughput per replica. Each point represents the performance achieved by replica-exchange MD simulations on a single protein/water system. Run parameters were one replica per node (each node has 6 GPUs), using between 24 and 40 replicas in a given system.

Figure 1 shows the performance per replica on Summit of T-REMD simulations for the majority of the simulations performed in this work using GROMACS version 2020.1. A few simulations were run with GROMACS 2018 and/or with different scheduling parameters and achieved only 20-50% of the performance shown above and are not included in the figure. We found that performance was maximized when running all bonded and nonbonded calculations on the GPUs (interatomic and both particle-mesh Ewald and pairwise Lennard-Jones). With the noted choices, performance saturates at around 100 ns/day for 250,000 atoms and above, even if more nodes are allocated per replica, for two reasons. First, the GPU-based fast Fourier transform is limited to a single GPU, and communication latencies between nodes slow down the calculation. However, throughput around 100 ns/day can still be achieved for simulations above 250,000 atoms if nodes are increased proportionately to system size.

b. T-REMD: Conformational Sampling of SARS-CoV-2 Proteins.

T-REMD simulations were performed with the number of replicas ranging from 20 to 60 for 750 ns each for an aggregate sampling of over 0.6 ms (Table I & SI Table I). Given the scaling data noted above, for the total 816 replicas simulated, the calculations (is performed simultaneously) used the equivalent of ~18% of the entire Summit supercomputer for ~3 days. The performance thus scales up to ~1 ms/day was the whole machine to be used for different proteins at the same time. For all systems, the replicate temperatures range from 310 K to ~350 K, and the average exchange probabilities were near 0.3.

From the simulations, structural diversity was quantified by calculating, when a binding site is known, the gyration tensor of the binding site residues, the solvent-accessible surface area (SASA) of the binding sites, and the construction of pairwise snapshot-snapshot root-mean-squared deviation (RMSD) matrices for the target temperature replica, i.e., the replica with the temperature set to 310 K (see *Methods* for calculation details). Additionally, using the gyration tensor, the shape anisotropy of the pockets was also obtained.

Linkage-based RMSD clustering, using the pairwise RMSD matrices was performed to gauge the overall structural diversity of the proteins. Figure 2 provides example conformations and calculated quantities for one example target, the neutral variant of the PL-Protease (PLPro). Similar plots for the other simulated systems are provided as supplemental material and on <https://coronavirus-hpc.ornl.gov>.

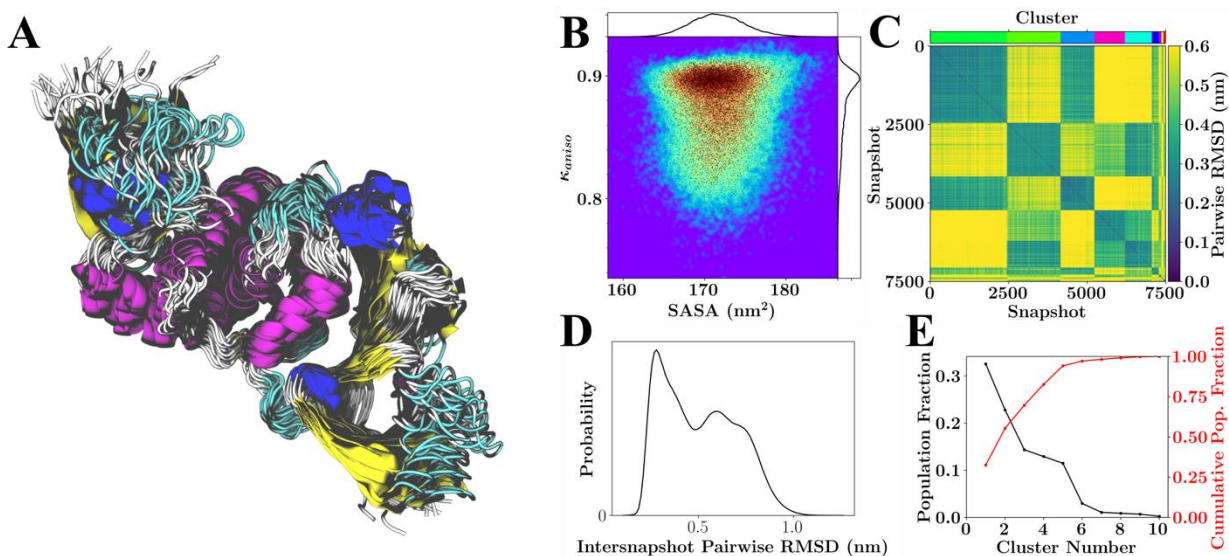


Figure 2. Configurational variability of PLPro with neutral HIS protonation states. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replica spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).

From Figure 2A (and subfigure A of the SI Figures 1 through 22), it is clear that the simulations generate a diverse ensemble of states with varied loop structures. For the case of the neutral variant PL-Protease, Figs.

3C and 3E indicate the existence of a number of dominant conformational states. Figs. 3D and 3C further suggest that, although six dominant states exist, these states could be grouped into two ‘super-states,’ which may indicate a switching like behavior or the potential existence of a ‘hinge.’ Finally, subfigure B shows a significant amount of sampling of rod-like geometries (anisotropies near 1); however, there are states that have a correlated reduction in SASA and shape anisotropy, which would correspond with a nearly continuous transition between rod-like structures and spheroid-like structures.

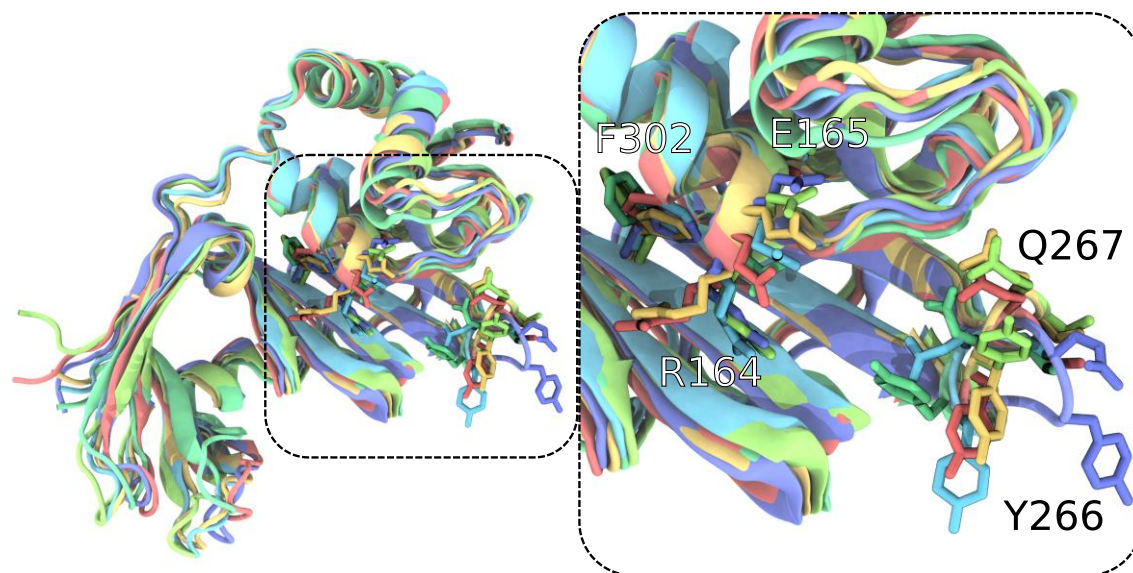


Figure 3. Configurational variability of the PLPro active site region generally bounded by the black dashed lines and the next step in analysis after Figure 3. Each of the differently colored aligned protein models represents the center of a populous cluster, as defined by active site conformation RMSD. Residues such as R164, E165, Y266, Q267, and F302 vary in conformation substantially and highlight the conformational variation within the ensemble created through T-REMD. For clearer visualization, only residues 91 and onward for PLPro are shown, as this selection was used for active site alignment. Within the VMD¹³⁰ rendering, sidechains are displayed without their hydrogens.

The general conformation variation highlighted by Figs. 2 and S1-22, to some degree, masks the conformational variation within binding sites; however, when for docking to the individual binding sites, clusters within the T-REMD trajectory are identified and demonstrate significant variability within the active site region (Fig. 4). While not specifically active site residues, residue variability at the tip of the loop centered on Y266 and the charged residue pair R164-E165 near the active site imply that accounting for the protein conformational ensemble is essential. Otherwise, the docking calculations would be strongly biased by the rotameric states present in the single static structure used in typical single-structure docking calculations.

c. Smaller Database Docking

A preliminary analysis was performed of general trends seen in docking the smaller database to the 23 SARA-CoV-2 protein systems. For each protein target, all the docking results from each of the 10 cluster configurations were combined, and the top 500 scoring compounds extracted. The selectivity of the compounds for any given target varies considerably (Table III) with the number of compounds present on any two different top 500 lists as low as 171 or as high as 283. In comparison, from two random selections of 500 items out of 9,014 items (see SI Figure 23, 5% percentile = 19 compounds, 95% percentile = 36 compounds), 27 identical compounds would be expected on average. Thus, the high number of identical top-scoring compounds observed between any two targets indicates a non-random selection of these duplicate compounds.

computational bias of some compounds based on other criteria than their good fit to the targets. On the other hand, such high numbers could correctly identify promiscuous binding sites that do not display marked structural specificity and hence could be indeed targeted by similar compounds. It is outside of the scope of the present work to assertively differentiate between these two possibilities. However, the number of duplicates varies greatly across several pairs of targets, which renders unlikely a systematic bias in the docking (because of, say, molecular weight or other ligand properties independent of the target's binding site).

Only ~55% of the top 500 compounds were the same in the docking results from the 100 ns and 750 ns clusters. Thus, extending the T-REMD simulation time by a factor of 7.5 nearly doubled the chemical diversity. Future analysis will be needed to indicate if the compounds that are identical in both sets of docking calculations are promiscuous compounds that would bind to many protein structures or if many of the clusters from the MD trajectories end up being selected by the same compounds.

d. Billion-Compound Supercomputing Screens.

We found that for the ligands with fewer numbers of rotatable bonds, such as found in the Enamine dataset, a docking calculation using 20 repeated runs could be performed in 0.5-2.5 seconds when using the Summit GPU (Fig. 5). The same set of ligands docked with Vina on Summit's CPUs showed a large spread of timings, with some ligands requiring nearly five minutes to complete (Fig. 5). In practice, this means that with GPU-enabled docking, it is feasible to flexibly dock a billion compounds in about a day on modern supercomputers, whereas with Vina, a similar calculation would require a multi-year effort on a university cluster. We confirmed that for ligands with less than 11 rotations, the Solis-Wets algorithm in AD-GPU provides equivalent results to the new ADADELTA algorithm⁴⁰. For the trimmed dataset, the top 5% of scores obtained with AD-GPU using the Solis-Wets algorithm formed an intersection with the top 5% of scores from Vina consisting of 18% of each top 5% set. Analysis of the full billion ligand sets is currently ongoing.

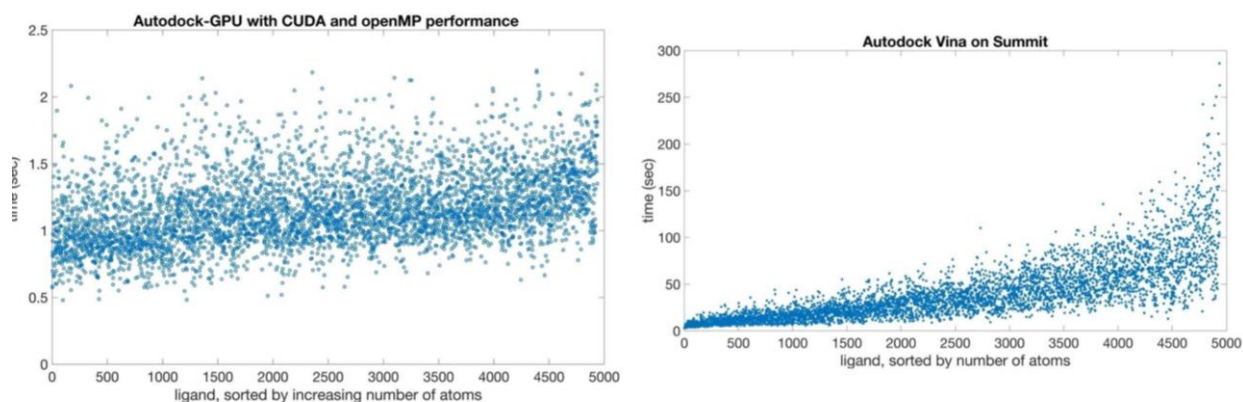


Figure 5. General benchmarking of Autodock-GPU and Autodock Vina performance against subset of Enamine database.

e. Mutation Analysis

The mutation frequency of the proteins simulated in this study is generally low. It should be noted, however, that it is as yet early in the history of SAR-CoV-2, and thus increased relative variability of residues along the proteome may indicate the propensity of those residues for future mutations. We did find higher variability, given by the entropy values, in other SARS-CoV-2 proteins not included in this study; in particular, the Spike mutation D614G noted in other reports continues to be seen with high frequency since

being described in a recent preprint that performed an analysis on GISAID through April 13¹⁵⁸. We counted 9107 D614G mutations (up from 3577 found April 13) and calculated entropy of 0.94 for this residue. The NSP12 RdRp protein also shows a large mutation entropy at residue 323, with a mutation entropy of 0.95. This residue, P, has mutated to L 9078 times (and F 3 times). Note that not every protein was represented in all sequences used for entropy calculations. Other regions of the genome with higher entropy values (greater than 0.5) are residues 203 and 204 of orf9 (entropy 0.70 and 0.69, respectively), residue 85 of NSP2 (0.74), 37 of NSP6 (0.57), 57 of orf3a (0.81), and 84 of orf8 (0.56).

The highest entropy found among the structures in this study was in the main protease, with an entropy of 0.13 for residue 15, a glycine. We found 261 G15S mutations and one G15D mutation in our dataset. The MPro also has a number of other residues with relatively high entropies, including residue 90, with entropy 0.07 and 117 K90R mutations, and 266 with entropy 0.04 and 64 A266V mutations. After this, the next highest entropy was 0.06 for the N-terminal region of the N protein and also for NSP15. These are displayed in Figure 6. An entropy of 0.04 was also found for domain X of NSP3 (glycine 76). A lower mutation entropy was found in PLPro, compared to MPro, with the highest value being 0.03. These mutations are important to consider when choosing targets for drug discovery, in that a protein that seems to be more rapidly mutating could potentially lead to an ineffective therapeutic if mutations alter the shape of the drug-binding site. In the case of MPro, the highest entropy mutations were not found in the active site; however, it is possible that they may still affect its conformation indirectly. The reduced mutation entropy for PLPro may indicate that an effort to target a protease could meet with fewer mutation-related problems if targeting PLPro rather than MPro.

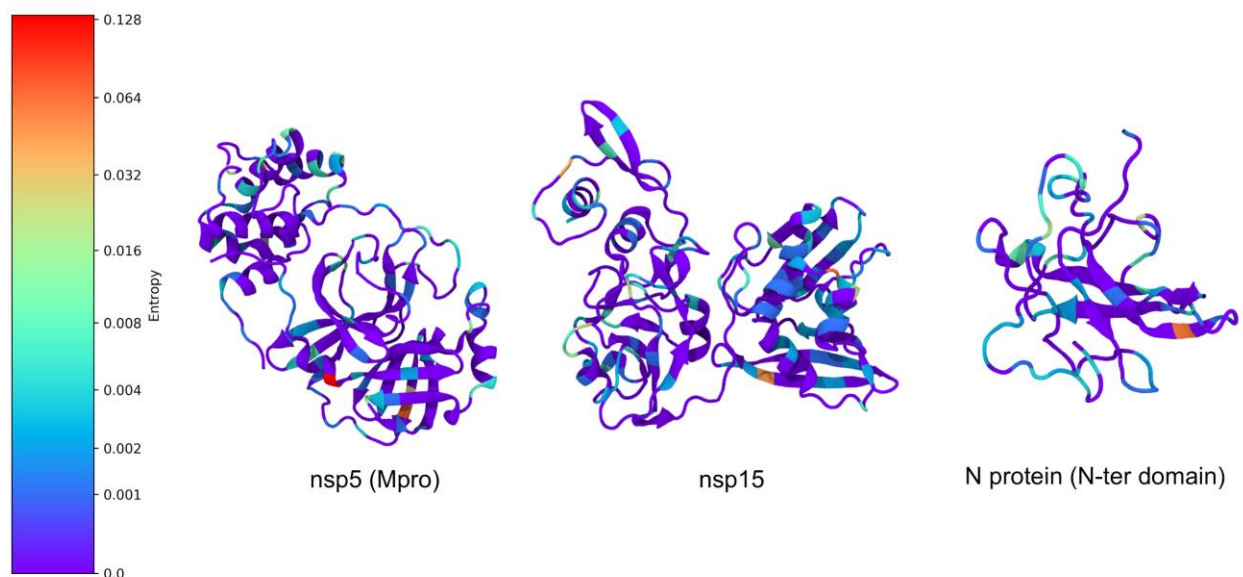


Figure 6. Example mutational entropy analysis. Residues are colored by entropy, with redder colors corresponding to greater entropy.

Future Directions and Preliminary Results from New Methodologies.

As emphasized above, this article is a progress report on an ongoing project. The development of the pipeline is continuing with advances being made in several directions. Notably, we are incorporating

artificial intelligence and machine learning into rescoring ligand ranking and clustering the MD trajectories. Further, we are developing methods to rescore docking using quantum chemical approaches. Although these developments have not been incorporated into the pipeline at the time of writing and were not applied to generate the results described above, we report on progress with them here.

a. Clustering MD trajectories using Deep Learning and AI.

The deluge of data generated from simulations such as the T-REMD runs reported here can make traditional approaches of machine learning and clustering approaches (based on measures of similarity in the RMSD-space, or other metrics) quite challenging. Often, practical aspects of computing dictate the use of subsample tracts of the MD data itself or use of prior knowledge about these datasets (e.g., knowing that the ligand binds only in a certain orientation) to filter such datasets. Deep learning techniques can be particularly valuable in ‘sifting’ through large datasets and can be powerful for clustering T-REMD simulations. We are investigating the use of a variational autoencoder with convolutional filters (CVAE), previously developed to cluster protein folding trajectories^{159, 160}, to cluster the T-REMD simulations of NSP15. As shown in Fig. 7, we find that the latent dimensions learned from the simulations indeed cluster the simulation data into a small number of conformational states. These states correspond to transitions observed in the simulations, as seen from various measured observables from the data such as the binding site RMSD, SASA, and the radius of gyration.

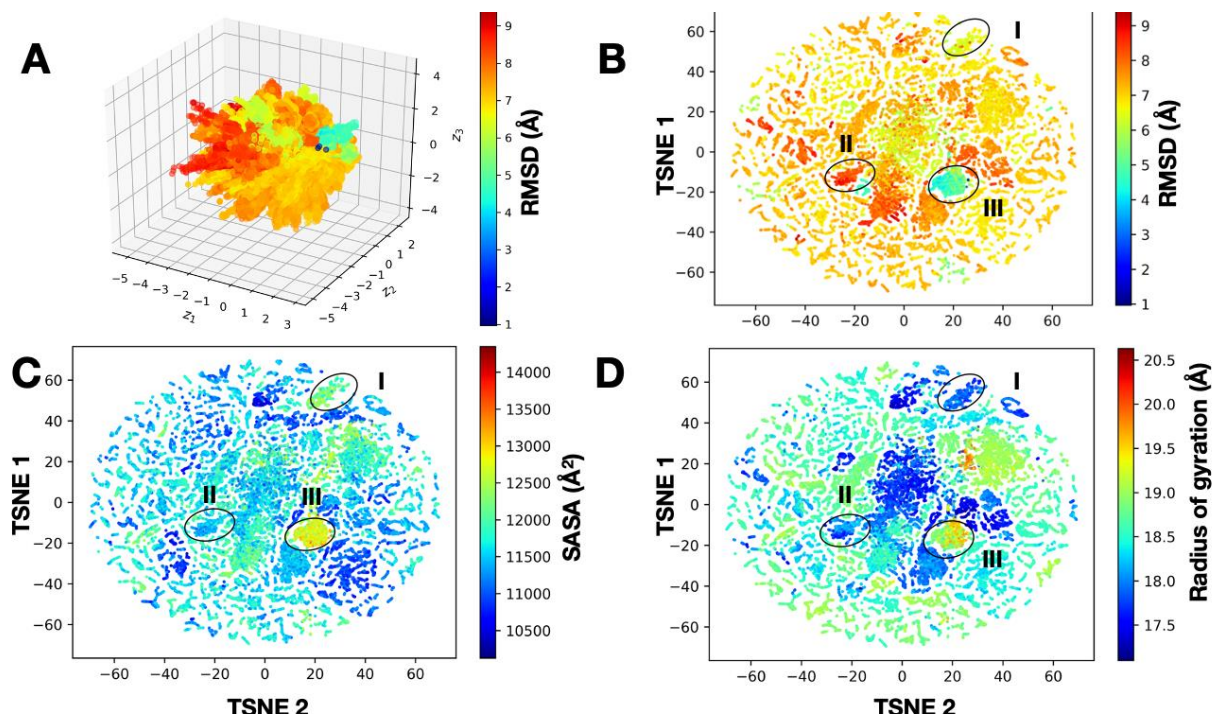


Figure 7. Deep learning clusters T-REMD simulations of the NSP15 hexameric complex into conformational states that are potentially relevant for docking studies. (A) A 3D-representation of the CVAE learned from the T-REMD simulations shows the presence of multiple conformational states. Each conformation from the simulation is painted using the RMSD to the starting structure and shows the presence of distinct directions in the conformational landscape where low- and high-RMSD structures are distributed. To understand this representation better, we use an at-stochastic neighbor embedding (t-SNE) algorithm to embed the data into a low-dimensional space, where we can clearly visualize how the conformational landscape is organized. In this two-dimensional space, we visualize various observables from the simulations, including (B) RMSD to the native structure, (C) SASA, and (D) radius of gyration. In each of these cases, we can observe the presence of at least three dominant sub-states with distinct structural characteristics, which can be further used for docking simulations.

The outcomes from the clustering provide insights into aspects of how the T-REMD simulations have sampled the conformational landscape - for example, in the case of this protein, as observed in Panel C, there is only one conformational state which has sampled a large SASA, indicating a potentially open state (which has only a minor change in the overall RMSD, Panel B). This information can be particularly helpful for selecting conformations and identifying metastable states for docking simulations¹⁶¹.

b. Protein-ligand rescoring using machine-learning.

The computational identification of drug compounds, and small molecules in general, that bind to a protein consists of three distinct tasks: 1) identifying a putative conformation of the protein-ligand complex (the docking problem); 2) given a docked conformation, determining whether or not the ligand is a true binder (the screening problem); and 3) is determined to be a true binder, ascertaining a relative, or better yet, absolute binding affinity (the affinity, or scoring, problem). In principle, one could perform the screening and affinity prediction problems using molecular dynamics techniques such as free energy perturbation, thermodynamic integration, or more approximate methods such as MM/PB(GB)SA (molecular mechanics/Poisson-Boltzmann[generalized Born]-surface area). However, this is computationally intractable for large numbers of compounds, even with supercomputers, and the accuracy can often be poor. Furthermore, these rigorous first principles-based methods assume a putative binding site, and cannot be applied to cases where the binding site is unknown. While the score or energy given by computational docking programs such as AutoDock Vina is reasonably well-suited for docking pose prediction, improvements are possible on the screening and affinity problems, and for this, we use here machine learning.

There is an ongoing need for the development of computationally tractable models that can be easily validated on benchmark docking data sets. To this end, accurate, physics-based, machine-learned models for the docking and affinity have been trained using the PDBbind database, a dataset consisting of experimentally determined protein-ligand complex structures with accurate experimental binding affinities¹⁶²⁻¹⁶⁵. On an independent data set, the CASF-2013 benchmark^{166, 167}, affinity prediction (random forest-based) models achieve, at best, a Pearson correlation (R^2) of 0.86, and docking pose prediction classifiers achieve an area under the curve of the receiver operator characteristic (AUC of ROC) of 0.91 using support vector machines (Demerdash *et al.*, *In Review*). While the random forest model trained on unnormalized features achieved the best R^2 at 0.86, a range of additional models (trained with random forest, gradient boosted trees, or support-vector machines using normalized or unnormalized features) achieved R^2 of 0.81-0.85. Regarding the docking pose prediction, the model used here achieves greater enrichment for native-like structures (78%) than AutoDock Vina (63%) (Demerdash *et al.*, *In Review*).

A model dedicated to virtual screening as a first step in triaging candidate molecules was developed. Once again, as with affinity and pose prediction, this model requires docked structures as input. This model is trained to discriminate between active and inactive compounds, and affinity ranking is performed as a second step only on the true active compounds. To this end, a support-vector machine-based model using the Dataset of Useful Decoys-Enhanced, a database of 102 proteins with experimentally verified active and inactive compounds, has been trained^{168, 169}. Preliminary performance on an independent validation set is encouraging, achieving AUC of ROC of 0.80 and recall of 0.76; that is, 76% of experimentally validated true positives were predicted positively by the model. This model is currently being subjected to further optimization, primarily through the calculation of additional physicochemical descriptors (features) and the optimization of hyperparameters.

Due to the urgent nature of the Covid-19 drug discovery campaign, computational expediency precluded calculating features on all docked structures for a given compound and, in turn, precluded running the docking pose classifier on the output of AutoDock Vina. Therefore, we relied on AutoDock Vina's ranking and not the machine-learned docking pose classifier, thereby reducing the number of feature calculations

that must be performed and increasing the throughput. (Parallelization efforts and code optimization are underway, so that feature calculation on all docked poses and subsequent application of the docking pose classifier becomes less onerous.) The virtual screening model was applied to these top-scoring structures from AutoDock Vina, generating a “binder” vs. “non-binder” classification. Subsequently, affinity prediction models were applied to just those complex structures classified as “binder.” The affinity prediction was performed using the range of high-performing models on docked structures corresponding to each MD cluster representative used in ensemble docking (See *Ensemble Docking with HPC Methods* for details.). This results in affinity predictions on typically 10 cluster representatives, each with affinity predictions from 5 machine-learned models (1 SVM, 2 boosted tree, and two random forest approaches), resulting in 50 “cases.” For each case, the top-500 ligands in terms of predicted affinity, were obtained. Molecules that appeared in the top 500 in at least 25 of the 50 cases were deemed hits and are presently undergoing experimental testing.

QM analysis of S-protein docking results

In a preliminary evaluation of the accuracy of FMO-DFTB/PCM in describing the interactions between ligands and the S-protein, we compare FMO-DFTB/PCM pair interaction energy (PIE) to that of the higher-level, but more expensive FMO-MP2/PCM method. The PIEs were calculated for ligands binding to the S-protein in the binding pocket. Figure 7 shows that FMO-DFTB/PCM interaction energies agree very well with high-level ab initio FMO-MP2/PCM data with the R correlation coefficient; in this case, it is 0.984. The high correlation between FMO-DFTB/PCM PIE and FMO-MP2/PCM PIE indicates that FMO-DFTB/PCM may be a fast and reliable QM-based method for interaction energy calculations.

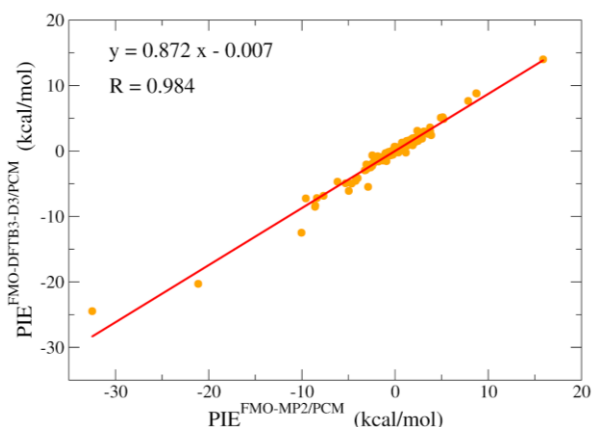


Figure 8. Pair interaction energy (PIE) decomposition analysis for FMO-DFTB/PCM plotted against FMO-MP2/3-21G/PCM data.

An updated preliminary homology model and T-REMD simulation similar to that reported by Smith & Smith of the S-protein RBD and re-docking to new clusters from followed by a re-evaluation of the top scoring 150 complexes with FMO-DFTB/PCM was performed as follows: 15 protein conformations and their ten strongest binding ligands predicted by AutoDock Vina were selected, and geometry optimizations were performed at the FMO-DFTB3-D3(BJ)/3ob/PCM level of theory. The binding energies of the top-3 best binders ranked by FMO-DFTB/PCM are listed in Table IV. According to our preliminary results, although FMO-DFTB/PCM agrees well with Vina in categorizing the strong binders, with all Vina, top-10 ligands have considerably stronger FMO-DFTB/PCM binding energies ($\Delta E^{\text{bind}} < -12$ kcal/mol) and the QM-based ranking is significantly different from the Vina ranking.

Table IV. The binding energy of the top-3 best ranked by FMO-DFTB/PCM and their binding free energy predicted by Autodock Vina.

Ligand SWEETLEAD ID	Protein Cluster ID	FMO-DFTB/PCM ΔE^{bind} (kcal/mol)	Vina ΔG^{bind} (kcal/mol)
4752	7	-67.75	-5.40
7055	11	-66.78	-7.60
4698	12	-66.41	-7.60

It is important to note that the current FMO-DFTB/PCM energy, which is based on solvent-corrected binding interaction energies is not the binding free energy. Various additional contributions to the binding free energy can be separately evaluated, and work is underway in this regard. For example, an entropic contribution can be estimated from vibrational frequencies once the requisite Hessian matrix is available.

Conclusions

The present manuscript reports on the establishment of a supercomputer-based virtual high-throughput screening ensemble-docking pipeline that takes into account the dynamic properties of protein targets, as well as preliminary results on simulations and docking screens to a number of protein targets from SARS-CoV-2.

The speed at which structural data have been derived experimentally for the SARS-CoV-2 proteome means that several of the simulations reported above were ‘out of date’ almost immediately. By this, we mean that the simulations were performed using models derived from experiments that had been superseded by higher-resolution or more complete data. Examples of these are the S-protein, MPro, N Protein, and NSP9. Clearly, as information on structures increases in quality, simulations will be further repeated. Furthermore, the complexity of the structural models derived is expected to increase. For example, models of the S protein interacting with the viral envelope or extending up to the complete virion can be envisaged and, in principle at least, incorporated into drug screening protocols.

The present results provide comprehensive simulation models for 8 of the viral proteins in 23 molecular systems. T-REMD is well suited for massively parallel supercomputing because many replicas are run simultaneously, and they need to communicate with each other. In the present tests, 350ns/day/replica was obtained for the smaller, and this, therefore, scales up to about 1.5ms/day of aggregate MD time, given the hypothetical situation that one had about 100 different proteins to run of roughly the same size. For bigger systems, with 10^5 atoms, the throughput is lower, about 1.0 ms/day. Nevertheless, it is clear that extensive simulation data can be obtained on many proteins with a short time-to-solution on this machine. As one possible future direction, one might envisage running T-REMD on the 44 drug targets that have been suggested as a minimal screen for the toxicity effects in human drug trials¹⁷⁰.

The ensemble docking performed so far mostly involves repurposing databases and therefore is limited to about 10k compounds. Many of these compounds are predicted to be quite promiscuous in binding to the targets. Two of the compounds identified in the top 1% of our preliminary S-protein screen have been reported to be in two registered clinical trials (quercetin and hypericin). Further, several compounds from the screens reported above show activity in reducing live viral infectivity: these results will be reported elsewhere.

The docking results using the smaller database were not run on Summit, because of the fact that for Summit code running on GPUs is preferred. However, as COVID-19 therapeutic research moves beyond repurposing to the discovery of novel compounds, there is a need to quickly screen many more compounds. Therefore, we have installed Autodock-GPU and demonstrated that it is capable of screening 1 billion compounds on Summit in 12 hours when scaled to the whole machine. Although several other groups have reported billion-compound screens, these have been using AI approaches or rigid docking without pose optimization^{171, 172}. The present billion-compound screen calculations, therefore, represent a potential supercomputer-driven paradigm shift in computational drug discovery and can be envisaged to be performed on dozens of proteins in a single day when the exascale era of supercomputing arrives, as planned for 2021.

Acknowledgments

This work was made possible in part by a grant of high-performance computing resources and technical support from the Alabama Supercomputer Authority to JB and KB.

CJC was supported by a National Science Foundation Graduate Research Fellowship under Grant No. 2017219379.

This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725 and National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility operated under Contract No. DE-AC02-05CH11231.

This research was supported by the Cancer Research Informatics Shared Resource Facility of the University of Kentucky Markey Cancer Center (P30CA177558) and the University of Kentucky's Center for Computational Sciences (CCS) high-performance computing resources.

References

1. DiMasi, J. A.; Grabowski, H. G.; Hansen, R. W., Innovation in the pharmaceutical industry: New estimates of R&D costs. *J Health Econ* **2016**, *47*, 20-33.
2. Silver, L. L., Challenges of antibacterial discovery. *Clinical microbiology reviews* **2011**, *24* (1), 71-109.
3. Kaldor, S. W.; Kalish, V. J.; Davies, J. F.; Shetty, B. V.; Fritz, J. E.; Appelt, K.; Burgess, J. A.; Campanale, K. M.; Chirgadze, N. Y.; Clawson, D. K., Viracept (nelfinavir mesylate, AG1343): a potent, orally bioavailable inhibitor of HIV-1 protease. *J. Med. Chem.* **1997**, *40* (24), 3979-3985.
4. von Itzstein, M.; Wu, W.-Y.; Kok, G. B.; Pegg, M. S.; Dyason, J. C.; Jin, B.; Van Phan, T.; Smythe, M. L.; White, H. F.; Oliver, S. W., Rational design of potent sialidase-based inhibitors of influenza virus replication. *Nature* **1993**, *363* (6428), 418-423.
5. Gorgulla, C.; Boeszoermyeni, A.; Wang, Z. F.; Fischer, P. D.; Coote, P. W.; Padmanabha Das, K. M.; Malets, Y. S.; Radchenko, D. S.; Moroz, Y. S.; Scott, D. A.; Fackeldey, K.; Hoffmann, M.; Iavniuk, I.; Wagner, G.; Arthanari, H., An open-source drug discovery platform enables ultra-large virtual screens. *Nature* **2020**, *580* (7805), 663-668.

6. Amaro, R. E.; Baudry, J.; Chodera, J.; Demir, Ö.; McCammon, J. A.; Miao, Y.; Smith, J. C., Ensemble Docking in Drug Discovery. *Biophysical Journal* **2018**, *114* (10), 2271-2278.
7. Teague, S. J., Implications of protein flexibility for drug discovery. *Nature reviews Drug discovery* **2003**, *2* (7), 527-541.
8. Carlson, H. A.; Masukawa, K. M.; McCammon, J. A., Method for Including the Dynamic Fluctuations of a Protein in Computer-Aided Drug Design. *The Journal of Physical Chemistry A* **1999**, *103* (49), 10213-10219.
9. Amaro, R. E.; Baudry, J.; Chodera, J.; Demir, Ö.; McCammon, J. A.; Miao, Y.; Smith, J. C., Ensemble Docking in Drug Discovery. *Biophys J* **2018**, *114* (10), 2271-2278.
10. Pi, M.; Kapoor, K.; Wu, Y.; Ye, R.; Senogles, S. E.; Nishimoto, S. K.; Hwang, D.-J.; Miller, D. D.; Narayanan, R.; Smith, J. C.; Baudry, J.; Quarles, L. D., Structural and Functional Evidence for Testosterone Activation of GPRC6A in Peripheral Tissues. *Molecular Endocrinology* **2015**, *29* (12), 1759-1773.
11. Pi, M.; Kapoor, K.; Ye, R.; Nishimoto, S. K.; Smith, J. C.; Baudry, J.; Quarles, L. D., Evidence for osteocalcin binding and activation of GPRC6A in β -cells. *Endocrinology* **2016**, *157* (5), 1866-1880.
12. Evangelista, W.; Weir, R. L.; Ellingson, S. R.; Harris, J. B.; Kapoor, K.; Smith, J. C.; Baudry, J., Ensemble-based docking: From hit discovery to metabolism and toxicity predictions. *Bioorg. Med. Chem* **2016**, *24* (20), 4928-4935.
13. Xiao, Z.; Riccardi, D.; Velazquez, H. A.; Chin, A. L.; Yates, C. R.; Carrick, J. D.; Smith, J. C.; Baudry, J.; Quarles, L. D., A computationally identified compound antagonizes excess FGF-23 signaling in renal tubules and a mouse model of hypophosphatemia. *Science Signaling* **2016**, *9* (455), ra113.
14. Abdali, N.; Parks, J. M.; Haynes, K. M.; Chaney, J. L.; Green, A. T.; Wolloscheck, D.; Walker, J. K.; Rybenkov, V. V.; Baudry, J.; Smith, J. C.; Zgurskaya, H. I., Reviving Antibiotics: Efflux Pump Inhibitors That Interact with AcrA, a Membrane Fusion Protein of the AcrAB-TolC Multidrug Efflux Pump. *ACS Infect. Dis.* **2017**, *3* (1), 89-98.
15. Dale, J. B.; Smeesters, P. R.; Courtney, H. S.; Penfound, T. A.; Hohn, C. M.; Smith, J. C.; Baudry, J. Y., Structure-based design of broadly protective group a streptococcal M protein-based vaccines. *Vaccine* **2017**, *35* (1), 19-26.
16. Haynes, K. M.; Abdali, N.; Jhawar, V.; Zgurskaya, H. I.; Parks, J. M.; Green, A. T.; Baudry, J.; Rybenkov, V. V.; Smith, J. C.; Walker, J. K., Identification and Structure–Activity Relationships of Novel Compounds that Potentiate the Activities of Antibiotics in *Escherichia coli*. *J Med Chem* **2017**, *60* (14), 6205-6219.
17. Velazquez, H. A.; Riccardi, D.; Xiao, Z.; Quarles, L. D.; Yates, C. R.; Baudry, J.; Smith, J. C., Ensemble docking to difficult targets in early-stage drug discovery: Methodology and application to fibroblast growth factor 23. *Chem Biol Drug Des* **2018**, *91* (2), 491-504.
18. Xiao, Z.; Baudry, J.; Cao, L.; Huang, J.; Chen, H.; Yates, C. R.; Li, W.; Dong, B.; Waters, C. M.; Smith, J. C., Polycystin-1 interacts with TAZ to stimulate osteoblastogenesis and inhibit adipogenesis. *J Clinical Invest* **2018**, *128* (1), 157-174.
19. Pi, M.; Kapoor, K.; Ye, R.; Hwang, D.-J.; Miller, D. D.; Smith, J. C.; Baudry, J.; Quarles, L. D., Computationally identified novel agonists for GPRC6A. *PloS one* **2018**, *13* (4).

20. Darzynkiewicz, Z. M.; Green, A. T.; Abdali, N.; Hazel, A.; Fulton, R. L.; Kimball, J.; Gryczynski, Z.; Gumbart, J. C.; Parks, J. M.; Smith, J. C., Identification of binding sites for efflux pump inhibitors of the AcrAB-TolC component AcrA. *Biophys J* **2019**, *116* (4), 648-658.
21. Evangelista Falcon, W.; Ellingson, S. R.; Smith, J. C.; Baudry, J., Ensemble Docking in Drug Discovery: How Many Protein Configurations from Molecular Dynamics Simulations are Needed To Reproduce Known Ligand Binding? *J Phys Chem* **2019**, *123* (25), 5189-5195.
22. Aranha, M. P.; Spooner, C.; Demerdash, O.; Czejdo, B.; Smith, J. C.; Mitchell, J. C., Prediction of peptide binding to MHC using machine learning with sequence and structure-based feature sets. *Biochim. Biophys. Acta-General Subjects* **2020**, 129535.
23. Micholas, S.; Jeremy C., S., Repurposing Therapeutics for COVID-19: Supercomputer-Based Docking to the SARS-CoV-2 Viral Spike Protein and Viral Spike Protein-Human ACE2 Interface. *ChemRxiv* **2020**, 10.26434/chemrxiv.11871402.v4.
24. Parks, J. M.; Smith, J. C., How to Discover Antiviral Drugs Quickly. *New England Journal of Medicine* **2020**.
25. Smith, M.; Smith, J., *Repurposing Therapeutics for COVID-19: Supercomputer-Based Docking to the SARS-CoV-2 Viral Spike Protein and Viral Spike Protein-Human ACE2 Interface*. 2020.
26. Evangelista Falcon, W.; Ellingson, S. R.; Smith, J. C.; Baudry, J., Ensemble Docking in Drug Discovery: How Many Protein Configurations from Molecular Dynamics Simulations are Needed To Reproduce Known Ligand Binding? *The Journal of Physical Chemistry B* **2019**, *123* (25), 5189-5195.
27. Abraham, M. J.; Murtola, T.; Schulz, R.; Páll, S.; Smith, J. C.; Hess, B.; Lindahl, E., GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX* **2015**, *1-2*, 19-25.
28. Kutzner, C.; Páll, S.; Fechner, M.; Esztermann, A.; de Groot, B. L.; Grubmüller, H., Best bang for your buck: GPU nodes for GROMACS biomolecular simulations. *Journal of Computational Chemistry* **2015**, *36* (26), 1990-2008.
29. Van Der Spoel, D.; Lindahl, E.; Hess, B.; Groenhof, G.; Mark, A. E.; Berendsen, H. J. C., GROMACS: Fast, flexible, and free. *Journal of Computational Chemistry* **2005**, *26* (16), 1701-1718.
30. Earl, D. J.; Deem, M. W., Parallel tempering: Theory, applications, and new perspectives. *Physical Chemistry Chemical Physics* **2005**, *7* (23), 3910-3916.
31. Hansmann, U. H. E., Parallel tempering algorithm for conformational studies of biological molecules. *Chem. Phys. Lett.* **1997**, *281* (1), 140-150.
32. Sugita, Y.; Okamoto, Y., Replica-exchange molecular dynamics method for protein folding. *Chemical Physics Letters* **1999**, *314* (1), 141-151.
33. Sugita, Y.; Kitao, A.; Okamoto, Y., Multidimensional replica-exchange method for free-energy calculations. *The Journal of Chemical Physics* **2000**, *113* (15), 6042-6051.
34. Bernardi, R. C.; Melo, M. C. R.; Schulten, K., Enhanced sampling techniques in molecular dynamics simulations of biological systems. *Biochimica et Biophysica Acta (BBA) - General Subjects* **2015**, *1850* (5), 872-877.
35. Tsai, T.-Y.; Chang, K.-W.; Chen, C. Y.-C., iScreen: world's first cloud-computing web server for virtual screening and de novo drug design based on TCM database@Taiwan. *Journal of Computer-Aided Molecular Design* **2011**, *25* (6), 525-531.

36. Dolezal, R.; Sobeslav, V.; Hornig, O.; Balik, L.; Korabecny, J.; Kuca, K. In *HPC Cloud Technologies for Virtual Screening in Drug Discovery*, Cham, Springer International Publishing: Cham, 2015; pp 440-449.
37. Huey, R.; Morris, G. M.; Olson, A. J.; Goodsell, D. S., A semiempirical free energy force field with charge-based desolvation. *Journal of computational chemistry* **2007**, *28* (6), 1145-1152.
38. Huey, R.; Goodsell, D. S.; Morris, G. M.; Olson, A. J., Grid-based hydrogen bond potentials with improved directionality. *Letters in Drug Design & Discovery* **2004**, *1* (2), 178-183.
39. El Khoury, L.; Santos-Martins, D.; Sasmal, S.; Eberhardt, J.; Bianco, G.; Ambrosio, F. A.; Solis-Vasquez, L.; Koch, A.; Forli, S.; Mobley, D. L., Comparison of affinity ranking using AutoDock-GPU and MM-GBSA scores for BACE-1 inhibitors in the D3R Grand Challenge 4. *Journal of Computer-Aided Molecular Design* **2019**, *33* (12), 1011-1020.
40. Solis-Vasquez, L.; Santos-Martins, D.; Koch, A.; Forli, S. In *Evaluating the Energy Efficiency of OpenCL-accelerated AutoDock Molecular Docking*, 2020 28th Euromicro International Conference on Parallel, Distributed and Network-Based Processing (PDP), IEEE: 2020; pp 162-166.
41. Santos-Martins, D.; Solis-Vasquez, L.; Koch, A.; Forli, S., Accelerating autodock4 with gpus and gradient-based local search. **2019**.
42. Fedorov, D. G., The fragment molecular orbital method: theoretical development, implementation in GAMESS , and applications. *WIREs Computational Molecular Science* **2017**, *7* (6).
43. Pronker, E. S.; Weenen, T. C.; Commandeur, H.; Claassen, E. H. J. H. M.; Osterhaus, A. D. M. E., Risk in Vaccine Research and Development Quantified. *PLOS ONE* **2013**, *8* (3), e57755.
44. Peeples, L., News Feature: Avoiding pitfalls in the pursuit of a COVID-19 vaccine. *PNAS* **2020**, *117* (15), 8218-8221.
45. Agostini, M. L.; Andres, E. L.; Sims, A. C.; Graham, R. L.; Sheahan, T. P.; Lu, X.; Smith, E. C.; Case, J. B.; Feng, J. Y.; Jordan, R., Coronavirus susceptibility to the antiviral remdesivir (GS-5734) is mediated by the viral polymerase and the proofreading exoribonuclease. *MBio* **2018**, *9* (2), e00221-18.
46. Brown, A. J.; Won, J. J.; Graham, R. L.; Dinnon III, K. H.; Sims, A. C.; Feng, J. Y.; Cihlar, T.; Denison, M. R.; Baric, R. S.; Sheahan, T. P., Broad spectrum antiviral remdesivir inhibits human endemic and zoonotic deltacoronaviruses with a highly divergent RNA dependent RNA polymerase. *Antiviral research* **2019**, *169*, 104541.
47. Sheahan, T. P.; Sims, A. C.; Leist, S. R.; Schäfer, A.; Won, J.; Brown, A. J.; Montgomery, S. A.; Hogg, A.; Babusis, D.; Clarke, M. O., Comparative therapeutic efficacy of remdesivir and combination lopinavir, ritonavir, and interferon beta against MERS-CoV. *Nat Commun* **2020**, *11* (1), 1-14.
48. Amanat, F.; Krammer, F., SARS-CoV-2 vaccines: status report. *Immunity* **2020**.
49. de Wit, E.; Feldmann, F.; Cronin, J.; Jordan, R.; Okumura, A.; Thomas, T.; Scott, D.; Cihlar, T.; Feldmann, H., Prophylactic and therapeutic remdesivir (GS-5734) treatment in the rhesus macaque model of MERS-CoV infection. *Proc Natl Acad Sci U S A* **2020**, *117* (12), 6771-6776.

50. Beigel, J. H.; Tomashek, K. M.; Dodd, L. E.; Mehta, A. K.; Zingman, B. S.; Kalil, A. C.; Hohmann, E.; Chu, H. Y.; Luetkemeyer, A.; Kline, S.; Lopez de Castilla, D.; Finberg, R. W.; Dierberg, K.; Tapson, V.; Hsieh, L.; Patterson, T. F.; Paredes, R.; Sweeney, D. A.; Short, W. R.; Touloumi, G.; Lye, D. C.; Ohmagari, N.; Oh, M. D.; Ruiz-Palacios, G. M.; Benfield, T.; Fatkenheuer, G.; Kortepeter, M. G.; Atmar, R. L.; Creech, C. B.; Lundgren, J.; Babiker, A. G.; Pett, S.; Neaton, J. D.; Burgess, T. H.; Bonnett, T.; Green, M.; Makowski, M.; Osinusi, A.; Nayak, S.; Lane, H. C.; Members, A.-S. G., Remdesivir for the Treatment of Covid-19 - Preliminary Report. *N Engl J Med* **2020**.
51. Johnson, R. M.; Vinetz, J. M., Dexamethasone in the management of covid -19. *BMJ* **2020**, *370*, m2648.
52. Omoto, S.; Speranzini, V.; Hashimoto, T.; Noshi, T.; Yamaguchi, H.; Kawai, M.; Kawaguchi, K.; Uehara, T.; Shishido, T.; Naito, A.; Cusack, S., Characterization of influenza virus variants induced by treatment with the endonuclease inhibitor baloxavir marboxil. *Sci Rep* **2018**, *8* (1), 9633-9633.
53. Kim, Y.; Jedrzejczak, R.; Maltseva, N. I.; Endres, M.; Godzik, A.; Michalska, K.; Joachimiak, A., Crystal structure of Nsp15 endoribonuclease NendoU from SARS-CoV-2. *bioRxiv* **2020**, 2020.03.02.968388.
54. Liu, C.; Zhou, Q.; Li, Y.; Garner, L. V.; Watkins, S. P.; Carter, L. J.; Smoot, J.; Gregg, A. C.; Daniels, A. D.; Jervey, S.; Albaiu, D., Research and Development on Therapeutic Agents and Vaccines for COVID-19 and Related Human Coronavirus Diseases. *ACS Cent. Sci* **2020**, *6* (3), 315-331.
55. Jiang, S.; Hillyer, C.; Du, L., Neutralizing Antibodies against SARS-CoV-2 and Other Human Coronaviruses. *Trends in Immunology* **2020**, *41* (5), 355-359.
56. Netherton, C. L.; Wileman, T., Virus factories, double membrane vesicles and viroplasm generated in animal cells. *Curr Opin Virol* **2011**, *1* (5), 381-387.
57. den Boon, J. A.; Diaz, A.; Ahlquist, P., Cytoplasmic viral replication complexes. *Cell host & microbe* **2010**, *8* (1), 77-85.
58. Hoffmann, M.; Kleine-Weber, H.; Pöhlmann, S., A Multibasic Cleavage Site in the Spike Protein of SARS-CoV-2 Is Essential for Infection of Human Lung Cells. *Molecular Cell* **2020**, *78* (4), 779-784.e5.
59. Hagemeyer, M. C.; Rottier, P. J.; De Haan, C. A., Biogenesis and dynamics of the coronavirus replicative structures. *Viruses* **2012**, *4* (11), 3245-3269.
60. Wu, A.; Peng, Y.; Huang, B.; Ding, X.; Wang, X.; Niu, P.; Meng, J.; Zhu, Z.; Zhang, Z.; Wang, J., Genome composition and divergence of the novel coronavirus (2019-nCoV) originating in China. *Cell host & microbe* **2020**.
61. Li, X.; Zhang, L.; Duan, Y.; Yu, J.; Wang, L.; Yang, K.; Liu, F.; You, T.; Liu, X.; Yang, X., Structure of Mpro from SARS-CoV-2 and discovery of its inhibitors. *Nature* **2020**.
62. Yin, W.; Mao, C.; Luan, X.; Shen, D.-D.; Shen, Q.; Su, H.; Wang, X.; Zhou, F.; Zhao, W.; Gao, M., Structural basis for inhibition of the RNA-dependent RNA polymerase from SARS-CoV-2 by remdesivir. *Science* **2020**.
63. Wang, X.; Cao, R.; Zhang, H.; Liu, J.; Xu, M.; Hu, H.; Li, Y.; Zhao, L.; Li, W.; Sun, X., The anti-influenza virus drug, arbidol is an efficient inhibitor of SARS-CoV-2 in vitro. *Cell Discovery* **2020**, *6* (1), 1-5.

64. Lei, J.; Hilgenfeld, R., RNA-virus proteases counteracting host innate immunity. *FEBS letters* **2017**, *591* (20), 3190-3210.
65. Lei, J.; Kusov, Y.; Hilgenfeld, R., Nsp3 of coronaviruses: Structures and functions of a large multi-domain protein. *Antiviral research* **2018**, *149*, 58-74.
66. Snijder, E.; Decroly, E.; Ziebuhr, J., The nonstructural proteins directing coronavirus RNA synthesis and processing. In *Advances in virus research*, Elsevier: 2016; Vol. 96, pp 59-126.
67. Owens, J., Determining druggability. *Nature Reviews Drug Discovery* **2007**, *6* (3), 187-187.
68. Wood, D. J.; Lopez-Fernandez, J. D.; Knight, L. E.; Al-Khawaldeh, I.; Gai, C.; Lin, S.; Martin, M. P.; Miller, D. C.; Cano, C.; Endicott, J. A.; Hardcastle, I. R.; Noble, M. E. M.; Waring, M. J., FragLites—Minimal, Halogenated Fragments Displaying Pharmacophore Doublets. An Efficient Approach to Druggability Assessment and Hit Generation. *Journal of Medicinal Chemistry* **2019**, *62* (7), 3741-3752.
69. Wu, C.; Liu, Y.; Yang, Y.; Zhang, P.; Zhong, W.; Wang, Y.; Wang, Q.; Xu, Y.; Li, M.; Li, X.; Zheng, M.; Chen, L.; Li, H., Analysis of therapeutic targets for SARS-CoV-2 and discovery of potential drugs by computational methods. *Acta Pharmaceutica Sinica B* **2020**, *10* (5), 766-788.
70. Richman, D. D., Antiviral drug resistance. *Antiviral Research* **2006**, *71* (2), 117-121.
71. Sanjuán, R.; Nebot, M. R.; Chirico, N.; Mansky, L. M.; Belshaw, R., Viral mutation rates. *J Virol* **2010**, *84* (19), 9733-9748.
72. Akand, E. H.; Downard, K. M., Ancestral and compensatory mutations that promote antiviral resistance in influenza N1 neuraminidase revealed by a phylonumerics approach. *Journal of Molecular Evolution* **2018**, *86* (8), 546-553.
73. Pybus, O. G.; Rambaut, A.; Belshaw, R.; Freckleton, R. P.; Drummond, A. J.; Holmes, E. C., Phylogenetic Evidence for Deleterious Mutation Load in RNA Viruses and Its Contribution to Viral Evolution. *Molecular Biology and Evolution* **2007**, *24* (3), 845-852.
74. van Dorp, L.; Acman, M.; Richard, D.; Shaw, L. P.; Ford, C. E.; Ormond, L.; Owen, C. J.; Pang, J.; Tan, C. C. S.; Boshier, F. A. T.; Ortiz, A. T.; Balloux, F., Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. *Infection, Genetics and Evolution* **2020**, *83*, 104351.
75. Gordon, D. E.; Jang, G. M.; Bouhaddou, M.; Xu, J.; Obernier, K.; White, K. M.; O'Meara, M. J.; Rezelj, V. V.; Guo, J. Z.; Swaney, D. L., A SARS-CoV-2 protein interaction map reveals targets for drug repurposing. *Nature* **2020**, 1-13.
76. Zhou, H.; Fang, Y.; Xu, T.; Ni, W. J.; Shen, A. Z.; Meng, X. M., Potential therapeutic targets and promising drugs for combating SARS-CoV-2. *British Journal of Pharmacology* **2020**.
77. Kim, Y.; Jedrzejczak, R.; Maltseva, N. I.; Wilamowski, M.; Endres, M.; Godzik, A.; Michalska, K.; Joachimiak, A., Crystal structure of Nsp15 endoribonuclease NendoU from SARS-CoV-2. *Protein Science* **2020**.
78. Jin, Z.; Du, X.; Xu, Y.; Deng, Y.; Liu, M.; Zhao, Y.; Zhang, B.; Li, X.; Zhang, L.; Peng, C., Structure of M pro from SARS-CoV-2 and discovery of its inhibitors. *Nature* **2020**, 1-5.
79. Lan, J.; Ge, J.; Yu, J.; Shan, S.; Zhou, H.; Fan, S.; Zhang, Q.; Shi, X.; Wang, Q.; Zhang, L., Structure of the SARS-CoV-2 spike receptor-binding domain bound to the ACE2 receptor. *Nature* **2020**, *581* (7807), 215-220.
80. Kang, S.; Yang, M.; Hong, Z.; Zhang, L.; Huang, Z.; Chen, X.; He, S.; Zhou, Z.; Zhou, Z.; Chen, Q., Crystal structure of SARS-CoV-2 nucleocapsid protein RNA binding domain reveals potential unique drug targeting sites. *Acta Pharmaceutica Sinica B* **2020**.

81. Maeda, Y.; Kinoshita, T., The acidic environment of the Golgi is critical for glycosylation and transport. In *Methods in enzymology*, Elsevier: 2010; Vol. 480, pp 495-510.
82. Wu, M., Grabe M, Adams S, Tsien RY, Moore HP, and Machen TE. *Mechanisms of pH regulation in the regulated secretory pathway. J Biol Chem* **2001**, *276*, 33027-33035.
83. Zumla, A.; Chan, J. F. W.; Azhar, E. I.; Hui, D. S. C.; Yuen, K.-Y., Coronaviruses — drug discovery and therapeutic options. *Nature Reviews Drug Discovery* **2016**, *15* (5), 327-347.
84. Chen, A. A.; Pappu, R. V., Parameters of Monovalent Ions in the AMBER-99 Forcefield: Assessment of Inaccuracies and Proposed Improvements. *The Journal of Physical Chemistry B* **2007**, *111* (41), 11884-11887.
85. Yoo, J.; Aksimentiev, A., Improved Parametrization of Li⁺, Na⁺, K⁺, and Mg²⁺ Ions for All-Atom Molecular Dynamics Simulations of Nucleic Acid Systems. *The Journal of Physical Chemistry Letters* **2012**, *3* (1), 45-50.
86. Ahlstrand, E.; Schpector, J. Z.; Friedman, R., Computer simulations of alkali-acetate solutions: Accuracy of the forcefields in difference concentrations. *The Journal of Chemical Physics* **2017**, *147* (19), 194102.
87. Jing, Z.; Liu, C.; Cheng, S. Y.; Qi, R.; Walker, B. D.; Piquemal, J.-P.; Ren, P., Polarizable Force Fields for Biomolecular Simulations: Recent Advances and Applications. *Annual Review of Biophysics* **2019**, *48* (1), 371-394.
88. Chang, C.-k.; Hou, M.-H.; Chang, C.-F.; Hsiao, C.-D.; Huang, T.-h., The SARS coronavirus nucleocapsid protein—forms and functions. *Antiviral research* **2014**, *103*, 39-50.
89. Astuti, I.; Ysrafil, Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2): An overview of viral structure and host response. *Diabetes & Metabolic Syndrome: Clinical Research & Reviews* **2020**, *14* (4), 407-412.
90. Báez-Santos, Y. M.; St John, S. E.; Mesecar, A. D., The SARS-coronavirus papain-like protease: structure, function and inhibition by designed antiviral compounds. *Antiviral research* **2015**, *115*, 21-38.
91. Michalska, K.; Kim, Y.; Jedrzejczak, R.; Maltseva, N. I.; Stols, L.; Endres, M.; Joachimiak, A., Crystal structures of SARS-CoV-2 ADP-ribose phosphatase (ADRP): from the apo form to ligand complexes. *bioRxiv* **2020**, 2020.05.14.096081.
92. Lei, J.; Hilgenfeld, R., RNA-virus proteases counteracting host innate immunity. *FEBS Letters* **2017**, *591* (20), 3190-3210.
93. 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, L. W.; Xu, W.; Xiao, G.; Qin, C.; Shi, Z.; Jiang, H.; Rao, Z.; Yang, H., Structure of Mpro from SARS-CoV-2 and discovery of its inhibitors. *Nature* **2020**, *582* (7811), 289-293.
94. Littler, D. R.; Gully, B. S.; Colson, R. N.; Rossjohn, J., Crystal Structure of the SARS-CoV-2 Non-structural Protein 9, Nsp9. *iScience* **2020**, *23* (7), 101258.
95. Kim, Y.; Jedrzejczak, R.; Maltseva, N. I.; Wilamowski, M.; Endres, M.; Godzik, A.; Michalska, K.; Joachimiak, A., Crystal structure of Nsp15 endoribonuclease NendoU from SARS-CoV-2. *Protein Science* **2020**, *29* (7), 1596-1605.
96. Rosas-Lemus, M.; Minasov, G.; Shuvalova, L.; Inniss, N. L.; Kiryukhina, O.; Wiersum, G.; Kim, Y.; Jedrzejczak, R.; Enders, M.; Jaroszewski, L.; Godzik, A.; Joachimiak, A.; Satchell, K.

- J. F., The crystal structure of nsp10-nsp16 heterodimer from SARS-CoV-2 in complex with S-adenosylmethionine. *bioRxiv* **2020**, 2020.04.17.047498.
97. Chen, Y.; Su, C.; Ke, M.; Jin, X.; Xu, L.; Zhang, Z.; Wu, A.; Sun, Y.; Yang, Z.; Tien, P.; Ahola, T.; Liang, Y.; Liu, X.; Guo, D., Biochemical and Structural Insights into the Mechanisms of SARS Coronavirus RNA Ribose 2'-O-Methylation by nsp16/nsp10 Protein Complex. *PLOS Pathogens* **2011**, *7* (10), e1002294.
98. Menachery, V. D.; Debbink, K.; Baric, R. S., Coronavirus non-structural protein 16: Evasion, attenuation, and possible treatments. *Virus Research* **2014**, *194*, 191-199.
99. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J., Efficient Estimation of Word Representations in Vector Space. *arXiv* **2013**, 1301.3781.
100. Robertson, S., The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends® in Information Retrieval* **2010**, *3* (4), 333-389.
101. Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. In *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*, Association for Computational Linguistics, Minneapolis, Minnesota, Minneapolis, Minnesota, 2019.
102. Jo, S.; Kim, T.; Iyer, V. G.; Im, W., CHARMM-GUI: a web-based graphical user interface for CHARMM. *J Comput Chem* **2008**, *29* (11), 1859-65.
103. Zhang, J.; Yang, W.; Piquemal, J.-P.; Ren, P., Modeling Structural Coordination and Ligand Binding in Zinc Proteins with a Polarizable Potential. *Journal of Chemical Theory and Computation* **2012**, *8* (4), 1314-1324.
104. Calimet, N.; Simonson, T., CysHis₂-Zn²⁺ interactions: Possibilities and limitations of a simple pairwise force field. *Journal of Molecular Graphics and Modelling* **2006**, *24* (5), 404-411.
105. Wambo, T. O.; Chen, L. Y.; McHardy, S. F.; Tsin, A. T., Molecular dynamics study of human carbonic anhydrase II in complex with Zn²⁺ and acetazolamide on the basis of all-atom force field simulations. *Biophysical Chemistry* **2016**, *214-215*, 54-60.
106. Fan, K.; Wei, P.; Feng, Q.; Chen, S.; Huang, C.; Ma, L.; Lai, B.; Pei, J.; Liu, Y.; Chen, J., Biosynthesis, purification, and substrate specificity of severe acute respiratory syndrome coronavirus 3C-like proteinase. *J. Biol. Chem.* **2004**, *279* (3), 1637-1642.
107. Zhang, L.; Lin, D.; Sun, X.; Curth, U.; Drosten, C.; Sauerhering, L.; Becker, S.; Rox, K.; Hilgenfeld, R., Crystal structure of SARS-CoV-2 main protease provides a basis for design of improved α -ketoamide inhibitors. *Science* **2020**, *368* (6489), 409-412.
108. Kneller, D. W.; Phillips, G.; O'Neill, H. M.; Jedrzejczak, R.; Stols, L.; Langan, P.; Joachimiak, A.; Coates, L.; Kovalevsky, A., Structural plasticity of SARS-CoV-2 3CL M(pro) active site cavity revealed by room temperature X-ray crystallography. *Nat Commun* **2020**, *11* (1), 3202.
109. Yang, H.; Yang, M.; Ding, Y.; Liu, Y.; Lou, Z.; Zhou, Z.; Sun, L.; Mo, L.; Ye, S.; Pang, H., The crystal structures of severe acute respiratory syndrome virus main protease and its complex with an inhibitor. *PNAS* **2003**, *100* (23), 13190-13195.
110. Wang, F.; Chen, C.; Tan, W.; Yang, K.; Yang, H., Structure of Main Protease from Human Coronavirus NL63: Insights for Wide Spectrum Anti-Coronavirus Drug Design. *Sci Rep* **2016**, *6* (1), 22677.
111. Anand, K.; Palm, G. J.; Mesters, J. R.; Siddell, S. G.; Ziebuhr, J.; Hilgenfeld, R., Structure of coronavirus main proteinase reveals combination of a chymotrypsin fold with an extra α -helical domain. *EMBO J* **2002**, *21* (13), 3213-3224.

112. Paasche, A.; Zipper, A.; Schäfer, S.; Ziebuhr, J.; Schirmeister, T.; Engels, B., Evidence for substrate binding-induced zwitterion formation in the catalytic Cys-His dyad of the SARS-CoV main protease. *Biochemistry* **2014**, *53* (37), 5930-5946.
113. Katarzyna, S.; Vicent, M., Revealing the Molecular Mechanisms of Proteolysis of SARS-CoV-2 Mpro from QM/MM Computational Methods. *ChemRxiv* **2020**, 10.26434/chemrxiv.12283967.v1.
114. Świderek, K.; Moliner, V., Revealing the molecular mechanisms of proteolysis of SARS-CoV-2 Mpro by QM/MM computational methods. *Chemical Science* **2020**.
115. Dolinsky, T. J.; Czodrowski, P.; Li, H.; Nielsen, J. E.; Jensen, J. H.; Klebe, G.; Baker, N. A., PDB2PQR: expanding and upgrading automated preparation of biomolecular structures for molecular simulations. *Nucleic Acids Res* **2007**, *35* (suppl_2), W522-W525.
116. Vermaas, J. V.; Hardy, D. J.; Stone, J. E.; Tajkhorshid, E.; Kohlmeyer, A., TopoGromacs: Automated Topology Conversion from CHARMM to GROMACS within VMD. *J. Chem. Inf. Model.* **2016**, *56* (6), 1112-1116.
117. Hess, B.; Kutzner, C.; van der Spoel, D.; Lindahl, E., GROMACS 4: Algorithms for Highly Efficient, Load-Balanced, and Scalable Molecular Simulation. *J. Chem. Theory Comput.* **2008**, *4* (3), 435-447.
118. Huang, J.; Rauscher, S.; Nawrocki, G.; Ran, T.; Feig, M.; de Groot, B. L.; Grubmüller, H.; MacKerell, A. D., CHARMM36m: an improved force field for folded and intrinsically disordered proteins. *Nature methods* **2017**, *14* (1), 71-73.
119. Abraham, M. J.; Gready, J. E., Optimization of parameters for molecular dynamics simulation using smooth particle-mesh Ewald in GROMACS 4.5. *J Comput Chem* **2011**, *32* (9), 2031-2040.
120. Hess, B., P-LINCS: A parallel linear constraint solver for molecular simulation. *J. Chem. Theory Comput.* **2008**, *4* (1), 116-122.
121. Berendsen, H. J. C.; Postma, J. P. M.; Gunsteren, W. F. v.; DiNola, A.; Haak, J. R., Molecular dynamics with coupling to an external bath. *The Journal of Chemical Physics* **1984**, *81* (8), 3684-3690.
122. Bernardi, R. C.; Melo, M. C. R.; Schulten, K., Enhanced sampling techniques in molecular dynamics simulations of biological systems. *Biochim Biophys Acta* **2015**, *1850* (5), 872-877.
123. Hansmann, U. H. E., Parallel tempering algorithm for conformational studies of biological molecules. *Chemical Physics Letters* **1997**, *281* (1), 140-150.
124. Nymeyer, H., How Efficient Is Replica Exchange Molecular Dynamics? An Analytic Approach. *J. Chem. Theory Comput.* **2008**, *4* (4), 626-636.
125. Yang, Y. I.; Shao, Q.; Zhang, J.; Yang, L.; Gao, Y. Q., Enhanced sampling in molecular dynamics. *J. Chem. Phys.* **2019**, *151* (7), 070902.
126. Patriksson, A.; van der Spoel, D., A temperature predictor for parallel tempering simulations. *Phys Chem Chem Phys* **2008**, *10* (15), 2073-2077.
127. Bussi, G.; Donadio, D.; Parrinello, M., Canonical sampling through velocity rescaling. *J. Chem. Phys.* **2007**, *126* (1).
128. Parrinello, M.; Rahman, A., polymorphic transitions in single-crystals - a new molecular dynamics method *J. Appl. Phys.* **1981**, *52* (12), 7182-7190.
129. Nosé, S.; Klein, M., Constant pressure molecular dynamics for molecular systems. *Molecular Physics* **1983**, *50* (5), 1055-1076.

130. Humphrey, W.; Dalke, A.; Schulten, K., VMD: visual molecular dynamics. *Journal of molecular graphics* **1996**, *14* (1), 33-38.
131. Bressert, E., *SciPy and NumPy: an overview for developers*. " O'Reilly Media, Inc.": 2012.
132. Theobald, D. L., Rapid calculation of RMSDs using a quaternion-based characteristic polynomial. *Acta Crystallographica Section A: Foundations of Crystallography* **2005**, *61* (4), 478-480.
133. Novick, P. A.; Ortiz, O. F.; Poelman, J.; Abdulhay, A. Y.; Pande, V. S., SWEETLEAD: an in silico database of approved drugs, regulated chemicals, and herbal isolates for computer-aided drug discovery. *PLoS One* **2013**, *8* (11), e79568.
134. Novick, P. A.; Ortiz, O. F.; Poelman, J.; Abdulhay, A. Y.; Pande, V. S., SWEETLEAD: an in silico database of approved drugs, regulated chemicals, and herbal isolates for computer-aided drug discovery. *PLoS One* **2013**, *8* (11).
135. Goede, A.; Dunkel, M.; Mester, N.; Frommel, C.; Preissner, R., SuperDrug: a conformational drug database. *Bioinformatics* **2005**, *21* (9), 1751-1753.
136. Siramshetty, V. B.; Eckert, O. A.; Gohlke, B.-O.; Goede, A.; Chen, Q.; Devarakonda, P.; Preissner, S.; Preissner, R., SuperDRUG2: a one stop resource for approved/ marketed drugs. *Nucleic Acids Res* **2018**, *46* (D1), D1137-D1143.
137. Repositories, N., Diversity Set Information. 2013.
138. Ellingson, S. R.; Smith, J. C.; Baudry, J., VinaMPI: Facilitating multiple receptor high-throughput virtual docking on high-performance computers. *J Comput Chem* **2013**, *34* (25), 2212-2221.
139. El-Hachem, N.; Haibe-Kains, B.; Khalil, A.; Kobeissy, F. H.; Nemer, G., AutoDock and AutoDockTools for protein-ligand docking: Beta-site amyloid precursor protein cleaving enzyme 1 (BACE1) as a case study. In *Neuroproteomics*, Springer: 2017; pp 391-403.
140. Morris, G. M.; Huey, R.; Olson, A. J., Using autodock for ligand-receptor docking. *Current protocols in bioinformatics* **2008**, *24* (1), 8.14. 1-8.14. 40.
141. Zhu, X.; Mitchell, J. C., KFC2: a knowledge-based hot spot prediction method based on interface solvation, atomic density, and plasticity features. *Proteins: Structure, Function, and Bioinformatics* **2011**, *79* (9), 2671-2683.
142. Macalino, S. J. Y.; Basith, S.; Clavio, N. A. B.; Chang, H.; Kang, S.; Choi, S., Evolution of In Silico Strategies for Protein-Protein Interaction Drug Discovery. *Molecules* **2018**, *23* (8).
143. Lagzian, M.; Valadan, R.; Saeedi, M.; Roozbeh, F.; Hedayatizadeh-Omran, A.; Amanlou, M.; Alizadeh-Navaei, R., Repurposing naproxen as a potential antiviral agent against SARS-CoV-2. **2020**, DOI: 10.21203/rs.3.rs-21833/v1.
144. Dinesh, D. C.; Chalupska, D.; Silhan, J.; Veverka, V.; Boura, E., Structural basis of RNA recognition by the SARS-CoV-2 nucleocapsid phosphoprotein. *bioRxiv* **2020**, 2020.04.02.022194.
145. Shivanyuk, A.; Ryabukhin, S.; Tolmachev, A.; Bogolyubsky, A.; Mykytenko, D.; Chupryna, A.; Heilman, W.; Kostyuk, A., Enamine real database: Making chemical diversity real. *Chemistry today* **2007**, *25* (6), 58-59.
146. Forli, W.; Halliday, S.; Belew, R.; Olson, A. J., AutoDock Version 4.2. Citeseer: 2012.
147. Gaillard, T., Evaluation of AutoDock and AutoDock Vina on the CASF-2013 Benchmark. *J Chem Inf Model* **2018**, *58* (8), 1697-1706.

148. Li, H.; Leung, K. S.; Wong, M. H.; Ballester, P. J., Improving AutoDock Vina Using Random Forest: The Growing Accuracy of Binding Affinity Prediction by the Effective Exploitation of Larger Data Sets. *Mol Inform* **2015**, *34* (2-3), 115-26.
149. GISAID. <https://www.gisaid.org/>.
150. GISAID, GISAID-Homepage.
151. Katoh, K.; Toh, H., Parallelization of the MAFFT multiple sequence alignment program. *Bioinformatics* **2010**, *26* (15), 1899-1900.
152. Lawrence, T. J.; Kauffman, K. T.; Amrine, K. C. H.; Carper, D. L.; Lee, R. S.; Becich, P. J.; Canales, C. J.; Ardell, D. H., FAST: FAST Analysis of Sequences Toolbox. *Front Genet* **2015**, *6*, 172.
153. Shannon, C. E., A Mathematical Theory of Communication. *Bell System Technical Journal* **1948**, *27* (3), 379-423.
154. Morao, I.; Heifetz, A.; Fedorov, D. G., Accurate Scoring in Seconds with the Fragment Molecular Orbital and Density-Functional Tight-Binding Methods. In *Quantum Mechanics in Drug Discovery*, Springer: 2020; pp 143-148.
155. Nishimoto, Y.; Fedorov, D. G.; Irle, S., Density-functional tight-binding combined with the fragment molecular orbital method. *J. Chem. Theory Comput.* **2014**, *10* (11), 4801-4812.
156. Barca, G. M. J.; Bertoni, C.; Carrington, L.; Datta, D.; De Silva, N.; Deustua, J. E.; Fedorov, D. G.; Gour, J. R.; Gunina, A. O.; Guidez, E.; Harville, T.; Irle, S.; Ivanic, J.; Kowalski, K.; Leang, S. S.; Li, H.; Li, W.; Lutz, J. J.; Magoulas, I.; Mato, J.; Mironov, V.; Nakata, H.; Pham, B. Q.; Piecuch, P.; Poole, D.; Pruitt, S. R.; Rendell, A. P.; Roskop, L. B.; Ruedenberg, K.; Sattasathuchana, T.; Schmidt, M. W.; Shen, J.; Slipchenko, L.; Sosonkina, M.; Sundriyal, V.; Tiwari, A.; Galvez Vallejo, J. L.; Westheimer, B.; Wloch, M.; Xu, P.; Zahariev, F.; Gordon, M. S., Recent developments in the general atomic and molecular electronic structure system. *J Chem Phys* **2020**, *152* (15), 154102.
157. Fedorov, D. G.; Kitaura, K.; Li, H.; Jensen, J. H.; Gordon, M. S., The polarizable continuum model (PCM) interfaced with the fragment molecular orbital method (FMO). *J. Comp. Chem.* **2006**, *27* (8), 976-985.
158. Korber, B.; Fischer, W.; Gnanakaran, S. G.; Yoon, H.; Theiler, J.; Abfalterer, W.; Foley, B.; Giorgi, E. E.; Bhattacharya, T.; Parker, M. D., Spike mutation pipeline reveals the emergence of a more transmissible form of SARS-CoV-2. *bioRxiv* **2020**.
159. Bhowmik, D.; Gao, S.; Young, M. T.; Ramanathan, A., Deep clustering of protein folding simulations. *BMC Bioinformatics* **2018**, *19* (18), 484.
160. Romero, R.; Ramanathan, A.; Yuen, T.; Bhowmik, D.; Mathew, M.; Munshi, L. B.; Javaid, S.; Bloch, M.; Lizneva, D.; Rahimova, A.; Khan, A.; Taneja, C.; Kim, S.-M.; Sun, L.; New, M. I.; Haider, S.; Zaidi, M., Mechanism of glucocerebrosidase activation and dysfunction in Gaucher disease unraveled by molecular dynamics and deep learning. *Proceedings of the National Academy of Sciences* **2019**, *116* (11), 5086-5095.
161. Lee, H.; Turilli, M.; Jha, S.; Bhowmik, D.; Ma, H.; Ramanathan, A. In *DeepDriveMD: Deep-Learning Driven Adaptive Molecular Simulations for Protein Folding*, 2019 IEEE/ACM Third Workshop on Deep Learning on Supercomputers (DLS), 17-17 Nov. 2019; 2019; pp 12-19.
162. Wang, R.; Fang, X.; Lu, Y.; Wang, S., The PDBbind database: Collection of binding affinities for protein–ligand complexes with known three-dimensional structures. *J. Med. Chem.* **2004**, *47* (12), 2977-2980.

163. Wang, R.; Fang, X.; Lu, Y.; Yang, C.-Y.; Wang, S., The PDBbind database: methodologies and updates. *Journal of medicinal chemistry* **2005**, *48* (12), 4111-4119.
164. Cheng, T.; Li, X.; Li, Y.; Liu, Z.; Wang, R., Comparative assessment of scoring functions on a diverse test set. *J. Chem. Inf. Model.* **2009**, *49* (4), 1079-1093.
165. Li, Y.; Han, L.; Liu, Z.; Wang, R., Comparative assessment of scoring functions on an updated benchmark: 2. Evaluation methods and general results. *J. Chem. Inf. Model.* **2014**, *54* (6), 1717-1736.
166. Li, Y.; Han, L.; Liu, Z.; Wang, R., Comparative assessment of scoring functions on an updated benchmark: 2. Evaluation methods and general results. *J Chem Inf Model* **2014**, *54* (6), 1717-36.
167. Li, Y.; Liu, Z.; Li, J.; Han, L.; Liu, J.; Zhao, Z.; Wang, R., Comparative assessment of scoring functions on an updated benchmark: 1. Compilation of the test set. *J Chem Inf Model* **2014**, *54* (6), 1700-16.
168. Huang, N.; Shoichet, B. K.; Irwin, J. J., Benchmarking Sets for Molecular Docking. *J. Med. Chem.* **2006**, *49* (23), 6789-6801.
169. Mysinger, M. M.; Carchia, M.; Irwin, J. J.; Shoichet, B. K., Directory of Useful Decoys, Enhanced (DUD-E): Better Ligands and Decoys for Better Benchmarking. *J. Med. Chem* **2012**, *55* (14), 6582-6594.
170. Bowes, J.; Brown, A. J.; Hamon, J.; Jarolimek, W.; Sridhar, A.; Waldron, G.; Whitebread, S., Reducing safety-related drug attrition: the use of in vitro pharmacological profiling. *Nature Reviews Drug Discovery* **2012**, *11* (12), 909-922.
171. Ton, A. T.; Gentile, F.; Hsing, M.; Ban, F.; Cherkasov, A., Rapid identification of potential inhibitors of SARS-CoV-2 main protease by deep docking of 1.3 billion compounds. *Molecular informatics* **2020**.
172. Gorgulla, C.; Boeszoermyeni, A.; Wang, Z.-F.; Fischer, P. D.; Coote, P. W.; Das, K. M. P.; Malets, Y. S.; Radchenko, D. S.; Moroz, Y. S.; Scott, D. A., An open-source drug discovery platform enables ultra-large virtual screens. *Nature* **2020**, *580* (7805), 663-668.

Supplementary Information

SI Table I. Summary of simulated systems and the number of replicas per system.

PDB	Name	Oligomerization State / Modeling Notes	# Replicas
6W41	<i>spike protein</i>	<i>Receptor binding domain only ("Apo")</i>	30
6W41	<i>spike protein</i>	<i>Receptor binding domain in complex with ACE2</i>	60

6Y2E	<i>MPro (main protease)</i>	<i>Monomer with default CHARMM protonation (all HIS to HSD)</i>	38
6Y2E	<i>MPro (main protease)</i>	<i>Dimer with default CHARMM-GUI protonation (all HIS to HSD)</i>	39
6WQF	<i>MPro (main Protease)</i>	<i>Charged dimer. Protonation states: HSD41, HSP64, HSP80, HSP163, HSP164, HSE172, HSE246</i>	40
6WQF	<i>MPro (main protease)</i>	<i>Dimer. Protonation states: HSD41, HSD64, HSD80, HSE163, HSE164, HSE172, HSE246</i>	40
6WQF	<i>MPro (main protease)</i>	<i>Monomer. Protonation states: HSD41, HSD64, HSD80, HSE163, HSE164, HSE172, HSE246</i>	38
6WQF	<i>MPro (main protease)</i>	<i>Dimer. Protonation states: HSE41, HSD64, HSD80, HSE163, HSE164, HSE172, HSE246</i>	40
6WQF	<i>MPro (main protease)</i>	<i>Monomer. Protonation states: HSE41, HSD64, HSD80, HSE163, HSE164, HSE172, HSE246</i>	38
6W4H	<i>NSP10</i>	<i>Monomer, explicit Zn(Cys)₄ coordination</i>	30
6W4H	<i>NSP16</i>	<i>Monomer</i>	27
6W4H	<i>NSP10/NSP16, Methyltransferase complex</i>	<i>Heterodimer with explicit Zn(Cys)₄ coordination</i>	40
6VWW	<i>NSP15, Endoribonuclease</i>	<i>Monomer, His tags removed</i>	36

6VWW	<i>NSP15, Endoribonuclease</i>	<i>Hexamer, His-tags removed</i>	57
6M3M	<i>N protein, N-terminal RNA binding domain</i>	<i>Monomer, no RNA, no Zn</i>	25
6M3M	<i>N protein, N-terminal RNA binding domain</i>	<i>Tetramer, no RNA, no Zn</i>	43
6VYO	<i>N protein, N-terminal binding domain (Phosphoprotein)</i>	<i>Monomer, Zn²⁺ present in crystal, not present in simulation,</i>	25
6VYO	<i>N protein, N-terminal RNA binding domain (Phosphoprotein)</i>	<i>Tetramer, Zn²⁺ Complexed w/explicit bonds (patched)</i>	39
6W02	<i>NSP3, phosphatase domain</i>	<i>Asymmetric unit (dimer), tags not removed, Apo</i>	35
6W4B	<i>NSP9</i>	<i>Monomer</i>	23
6W4B	<i>NSP9</i>	<i>Dimer</i>	33
6W9C	<i>PL protease</i>	<i>Charged Monomer, computationally refined structure. Zn²⁺ complexed w/explicit bonds (patched). PropKa based assignment Protonation States: HSP15, HSP45, HSP48, HSP71, HSP87, HSD173, HSP253, HSP270, HSD273,</i>	38
6WRH	<i>PL protease</i>	<i>Neutral Monomer, Zn²⁺ complexed w/explicit bonds (patched). Neutral Protonation states (manual assignment): HSE17, HSE47, HSE50, HSE73, HSD89, HSD175, HSE255, HSE272, HSD275</i>	38

SI Table II: Summary of the docking performed. The docking was performed on the top 10 conformations of the clusters calculated from 750ns TREMD simulations. The table provides information on the oligomeric state, the docking region, the residues used for clustering, the center of the docking box, and the docking box size.

Simulation	Docking region	Clustering residues	Atom at the center of docking box	Docking box size
NSP10 monomer (PDB Id: 6W4H)	Binding interface to NSP16	4253 to 4267, 4295 to 4304, 4309 to 4314, 4321 to 4324, 4337 to 4340, 4347 to 4349	OH TYR 4349	
			O LYS P4296	
NSP16 monomer (PDB Id: 6W4H)	Binding interface to NSP10	6834 to 6839, 6874 to 6878, 6880 to 6882, 6884 to 6886, 6888 to 6891, 6900 to 6907, 7042 to 7047	HE21 GLN 6885	
	SAM-binding site	6841, 6845, 6867, 6869 to 6872, 6877 to 6879, 6897 to 6899, 6911 to 6913, 6928 to 6931, 6947, 6968	CG ASP 6928	
Complex NSP10-NSP16 (PDB Id: 6W4H)	Complex interface	6832 to 6842, 6871 to 6886, 6888 to 6891, 6900 to 6907, 7041 to 7049, 7085 to 7095, 4253 to 4267, 4295 to 4304, 4309 to 4314, 4321 to 4324, 4337 to 4340, 4347 to 4349	OH TYR 4349	
			OH TYR 4349	

Tetramer N-terminal domain of nucleocapsid (N) protein (PDB IDs: 6VYO)	RNA binding sites	88, 92, 107	geometric center of ARG 88, ARG 92, and ARG 107	30 x 30 x 30 Å
Dimer Main Protease (All HSD) (PDB Id: 6Y2E)	dimer interface		-23.53, -3.61, -10.42	30 x 29 x 28 Å
Tetramer Nucleocapsid (N) Protein Complex (PDB Id: 6VYO)	Three binding sites predicted by FTMap		-23.53, -3.61, -10.42	30 x 29 x 28 Å
			-23.53, -3.61, -10.42	30 x 29 x 28 Å
			-23.53, -3.61, -10.42	30 x 29 x 28 Å
Nucleocapsid Monomer (PDB Id: 6KL5)			34.42, 28.40, 29.28	28 x 28 x 28 Å
MPro Monomer (three different protonation variants) (PDB Ids: 6WQF and 6Y2E)	Catalytic pocket		CA Tyr 37	30 x 30 x 30 Å

MPro Dimer (PDB Ids: 6WQF & 6Y2E)	protein-protein interface These residues were close to the two cavities/pocket identified by the MOE site finder in the PPI of the dimer	Chain1: SER1 GLY2 PHE3 ARG4 LYS5 MET6 ALA7 VAL125 TYR126 GLN127 CYS128 ALA129 ARG131 THR135 ILE136 LYS137 GLY138 SER139 GLY170 VAL171 HSD172 ASP197 THR198 THR199 VAL204 TRP207 ILE213 ASN214 ASP216 TRP218 TYR239 ILE281 LEU282 GLY283 SER284 ALA285 LEU286 LEU287 GLU288 ASP289 GLU290 PHE291 Chain 2:PHE3 ARG4 LYS5 MET6 ALA7 VAL125 TYR126 GLN127 CYS128 ALA129 ARG131 THR135 ILE136 LYS137 GLY138 SER139 VAL171 ALA193 ALA194 GLY195 THR196 ASP197 THR198 THR199 VAL204 TRP207 ASN214 TYR237 ASN238 TYR239 LEU272 LEU282 GLY283 SER284 ALA285 LEU286 LEU287 GLU288 ASP289 GLU290 PHE291	CA TYR 126 (chain A)	40 x 40 x 40 Å
		Chain1: GLY11 LYS12 VAL13 GLU14 LYS97 PRO99 LYS100 ASP155 CYS156 VAL157 PHE305 GLN306 Chain2: GLY11 LYS12 GLU14 GLY15 MET17 VAL18 GLN19 GLN69 ALA70 GLY71 LYS97 ASN119 GLY120 SER121 PRO122	CA SER 10 (chain A)	40 x 40 x 40 Å
NSP9 Monomer (PDB Id: 6W4B)	Protein-Protein Interface	6 to 9, 97 98 100 102 104 105 106 108 109 110 112	2.99, 15.23, 2.21	112 x 36 x 76 Å

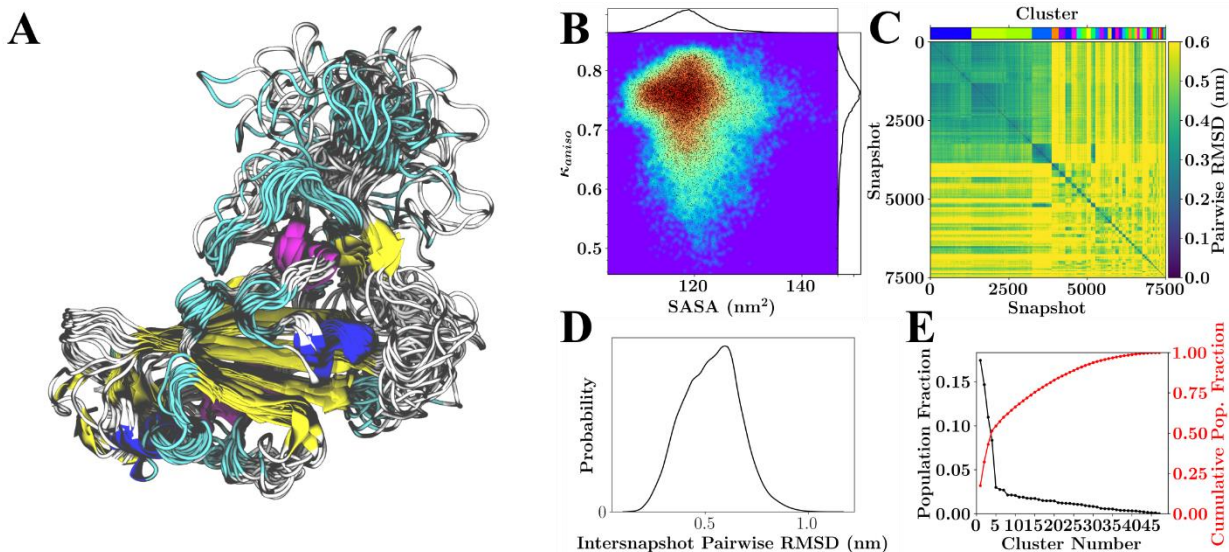
NSP9 Dimer (PDB Id: 6W4B)	Whole protein		CA MET 102	40 x 40 x 40 Å
NSP15 Monomer (PDB Id: 6VWW)	Large region including the catalytic pocket		CA VAL 166	70 x 40 x 40 Å
NSP15 Dimer (PDB Id: 6VWW)	Centered on hydrophobic core of the protein and the protein-protein interface		CA ASN 1413 (chain E)	70 x 50 x 50 Å
PLPro (PDB Id: 6W02)	Dimer	LEU160 GLY161 ASP162 VAL163 ARG164 MET206 SER243 ALA244 PRO245 PRO246 TYR262 TYR266 GLN267 CYS268 TYR271 THR299 ASP300 SER101 CYS109 ALA112 THR113 LEU116 HSP270 TYR271 LYS272 HSD273 ILE283 ASP284 GLY285 TYR 265 TYR 269 ASP 165	CA TYR 262	40 x 40 x 40 Å
NSP3 Monomer (PDB Id: 6W9C)	Catalytic pocket	21 to 24, 44 to 52, 38 to 40, 125 to 133, 154 to 160	-18.67, 1.68, 10.08	20 x 24 x 12 Å

S Protein RBD Apo PDB Id:	Spike-ACE2 Interface	613 614 615 636 637 638 639 664 665 666 667 668 669 670 671 672 673 674 678 679 680 681 682 683 684 685 686 687 688 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773		
S Protein RBD with ACE2	Spike ACE2 Interface	0 1 10 11 12 13 14 15 16 17 18 19 2 20 21 22 23 24 25 26 27 28 29 290 293 3 30 302 303 304 305 306 307 308 309 31 310 311 312 313 314 315 317 32 321 328 329 33 330 331 332 333 334 335 336 337 338 339 34 340 357 359 36 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 38 39 4 40 41 42 43 44 45 46 47 48 49 5 50 51 52 53 536 539 54 540 541 542 543 544 55 56 57 58 59 6 60 61 613 614 615 62 63 636 637 638 639 64 65 66 664 665 666 667 668 669 67 670 671 672 673 674 678 679 68 680 681 682 683 684 685 686 687 688 69 7 70 700 701 702 703 704 705 706 707 708 709 71 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 73 731 732 733 734 735 736 737 738 739 74 740 741 742 743 744 745 746 747 748 749 75 750 751 752 753 754 755 756 757 758 759 76 760 761 762 763 764 765 766 767 768 769 77 770 771 772 773 78 79 8 80 81 82 83 84 87 9		

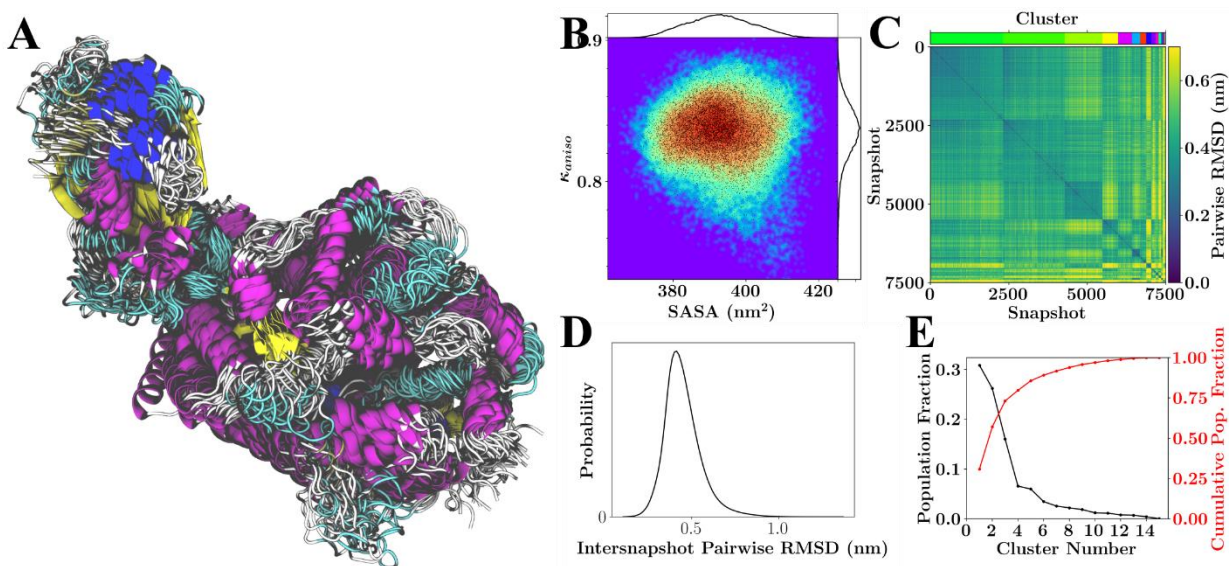
SI Tables III & IV. Identity between the compounds predicted to be in the top 500 lists from the 100ns and 750 ns MD simulation on selected targets.

SI Table III. Identity between compounds predicted to be in the top 500 lists from clusters derived from the first 100ns of T-REMD clusters.

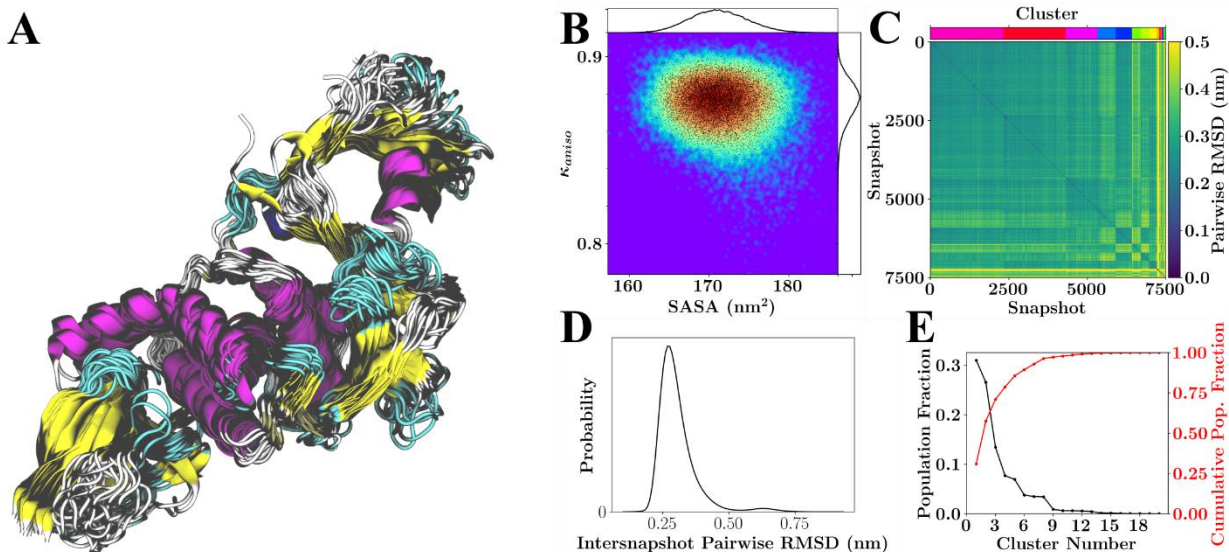
target	overlap in top 500 (%)
Nsp15, endoribonuclease hexamer	62.4
Nsp15, endoribonuclease monomer	59.4
MPro dimer interface	55.4
Nsp10/Nsp16 complex patched	45.2
Nsp9 dimer	56.8
N-protein tetramer	67.2
PLPro (neutral)	51.2
S-protein apo	60
S-protein Ace-2 complex	46.8
N-protein	61.6



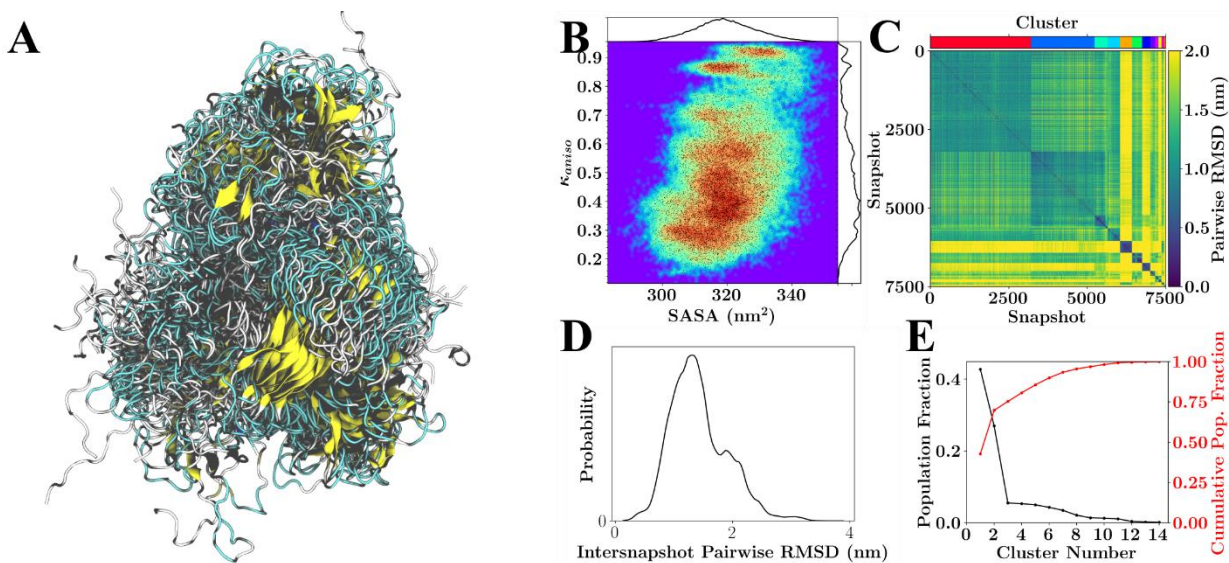
SI Figure 1) General results for simulations of S Protein RBD without ACE2. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



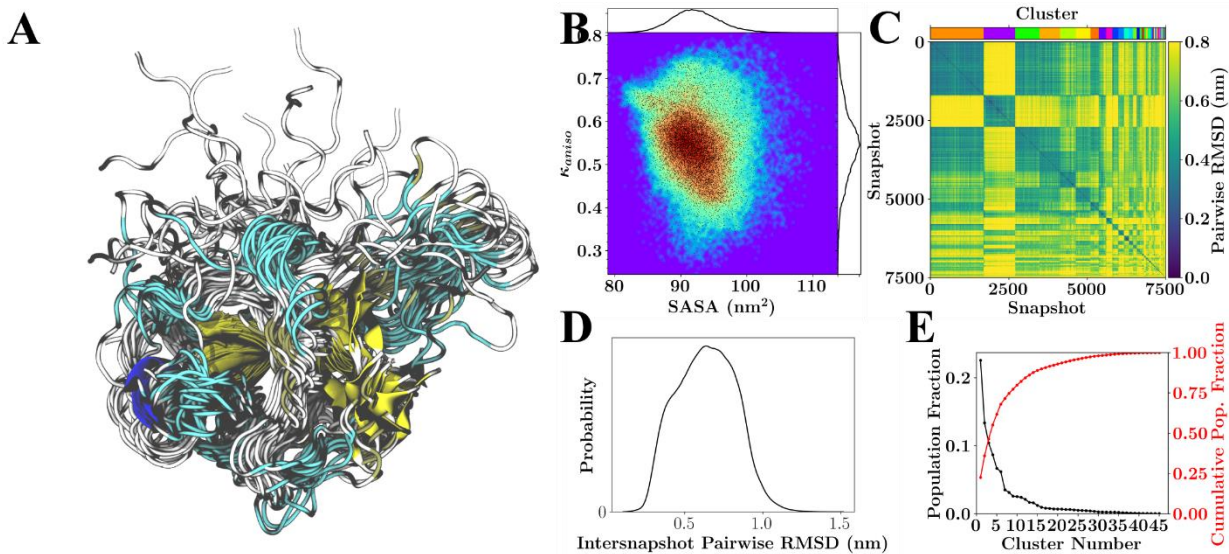
SI Figure 2. General results for simulations of S Protein in complex with ACE2. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



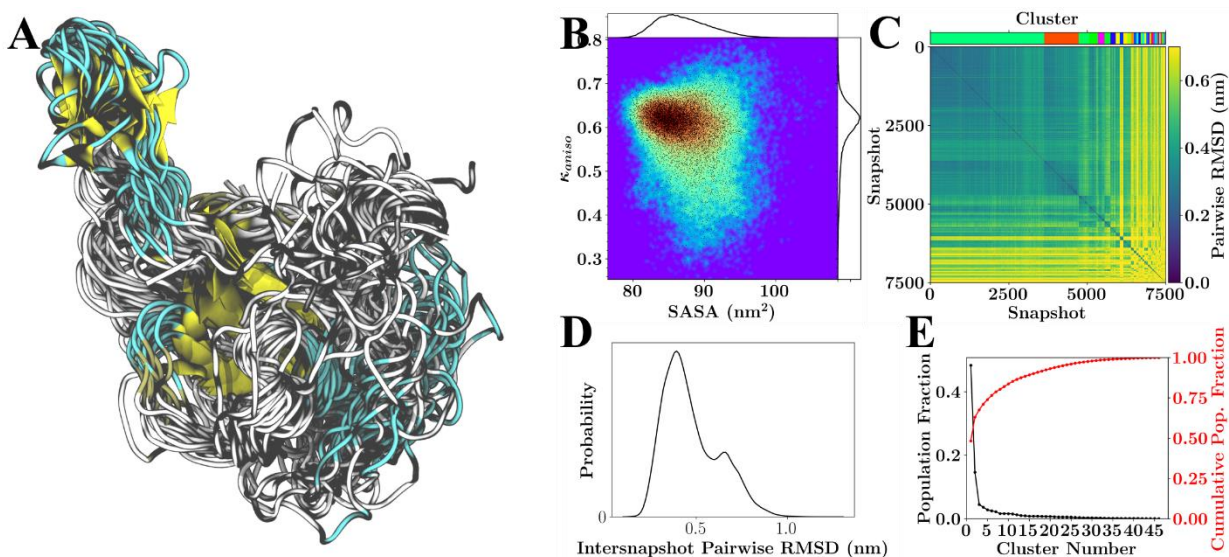
SI Figure 3. General results for simulations of PLPro with charged HIS protonation states. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



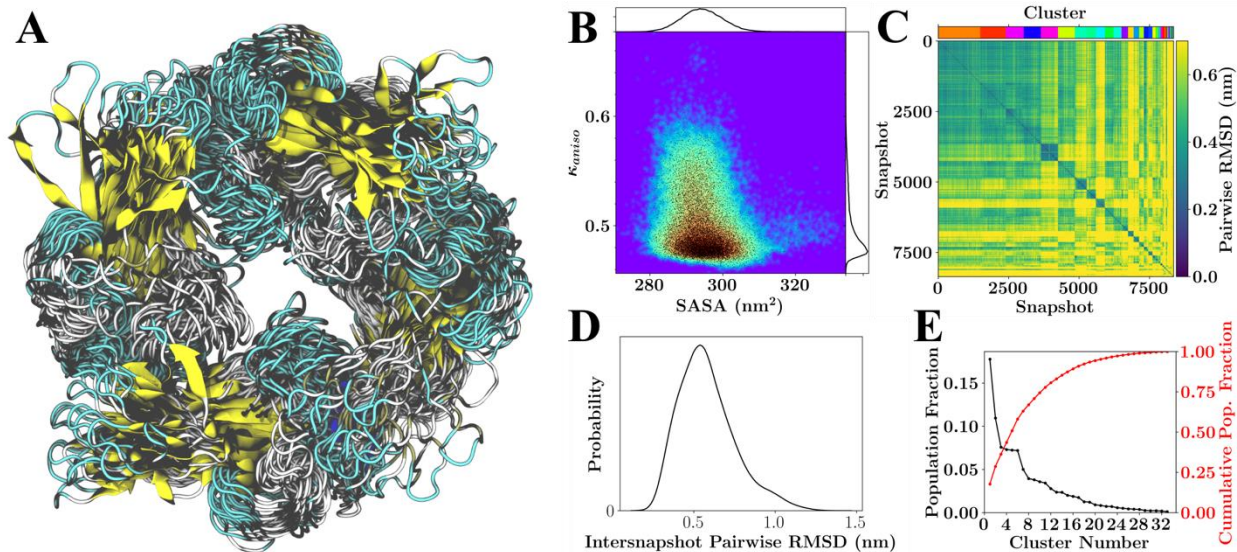
SI Figure 4. General results for simulations of the tetrameric complex of the N-terminal of the N Protein in the absence of Zn. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



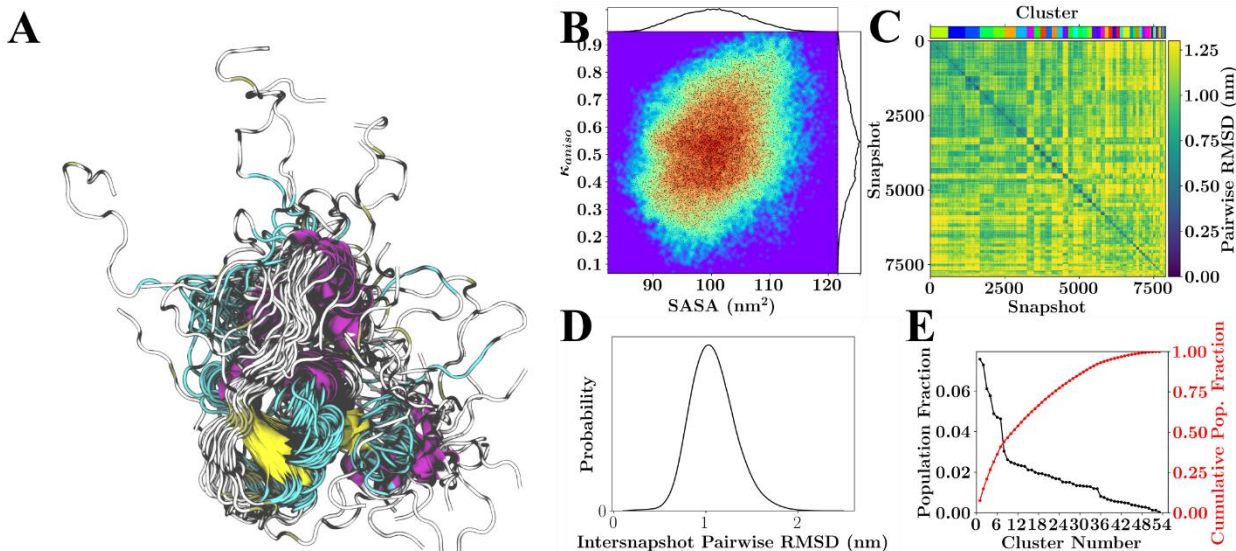
SI Figure 5. General results for simulations of N Protein N-terminus monomer using N Protein N-terminus without bound Zn as the initial structure. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



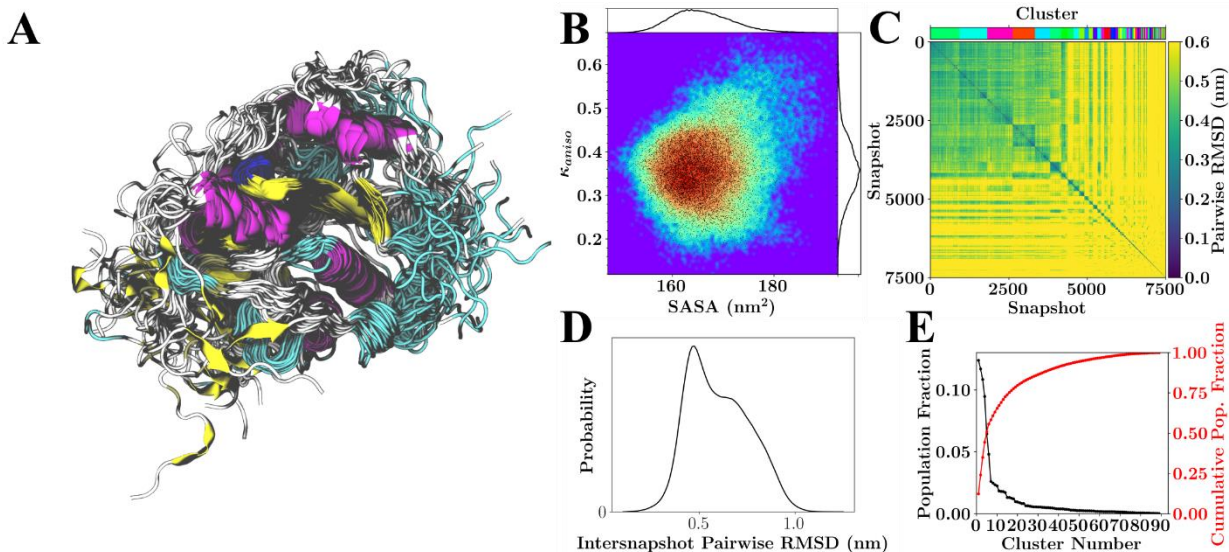
SI Figure 6. General results for simulations of N Protein Phosphoprotein (N-terminus) monomer states with initial conditions derived from the N Protein complex crystallized with Zn (though no Zn is bound in the monomer, derived from PDB: 6VYO). (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



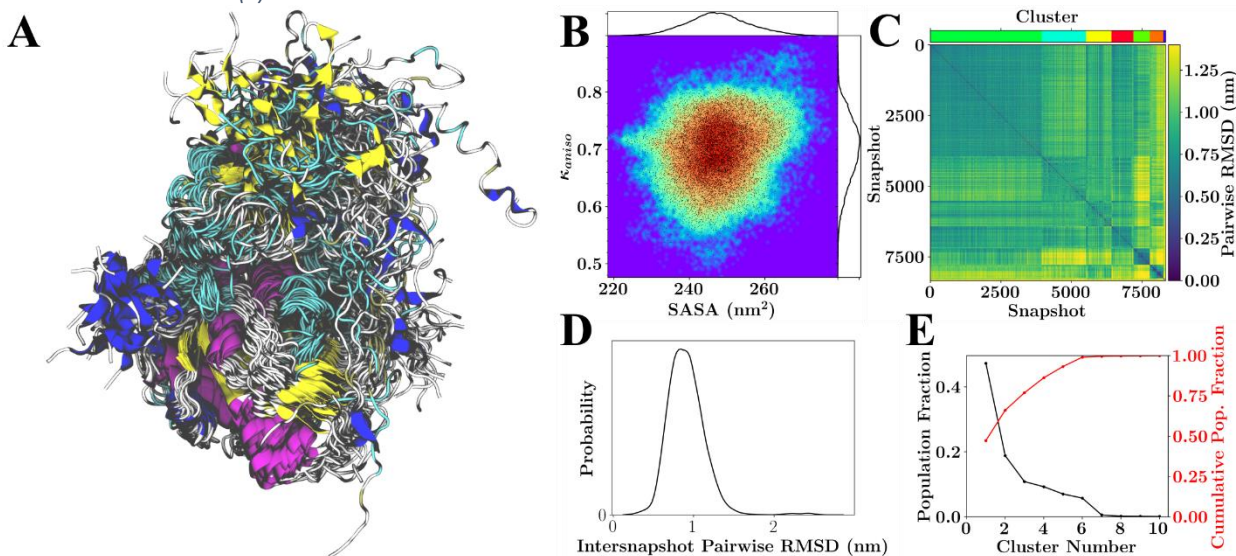
SI Figure 7. General results for simulations of N Protein Phosphoprotein (N-terminus) complex states with initial conditions derived from the N Protein complex crystallized in the presence of Zn (PDB 6VYO). (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



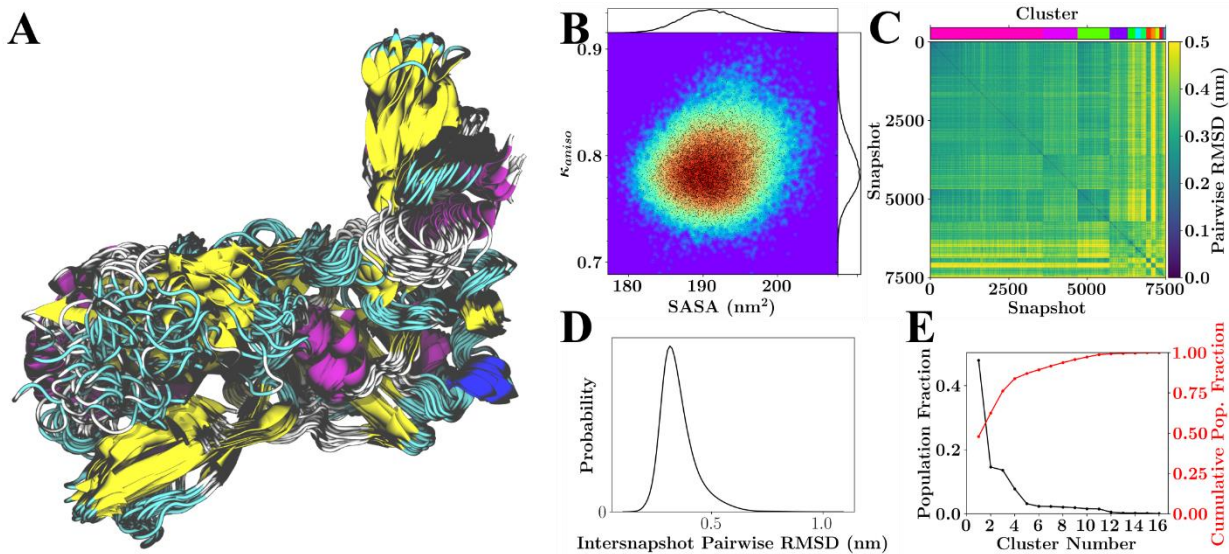
SI Figure 8. General results for simulations of NSP10 monomer. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



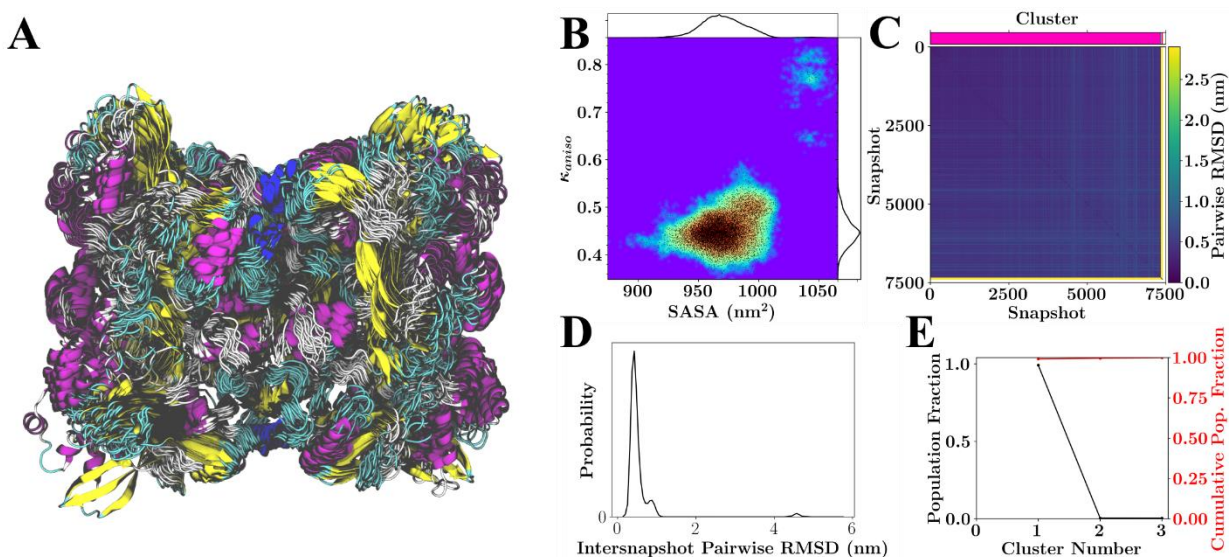
SI Figure 9. General results for simulations of NSP16 monomer. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



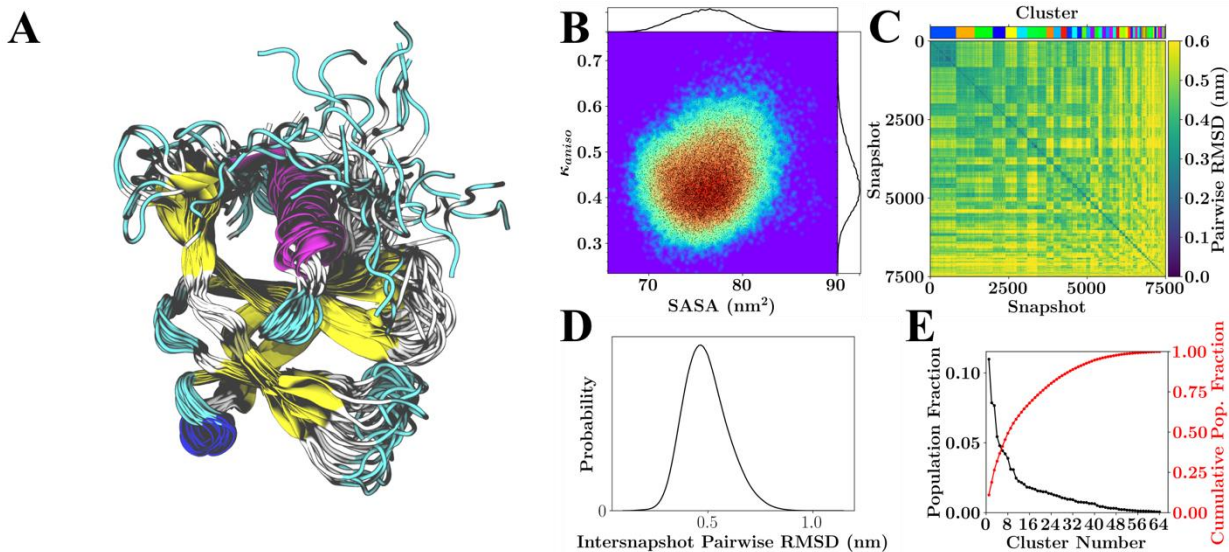
SI Figure 10. General results for simulations of NSP10+NSP16 Heterodimer. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



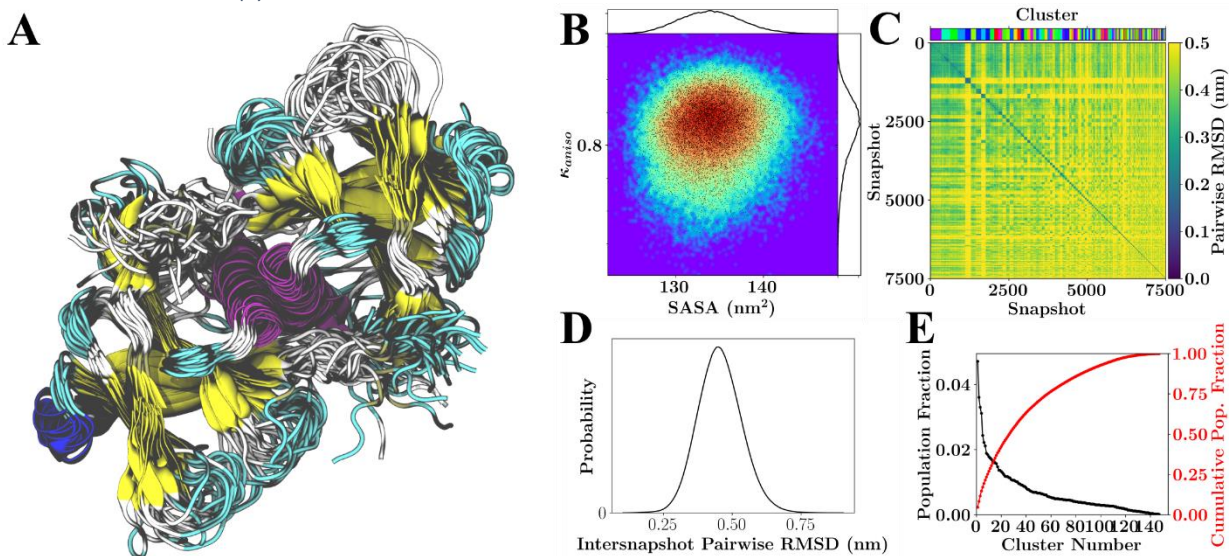
SI Figure 11. General results for simulations of NSP15 monomer. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



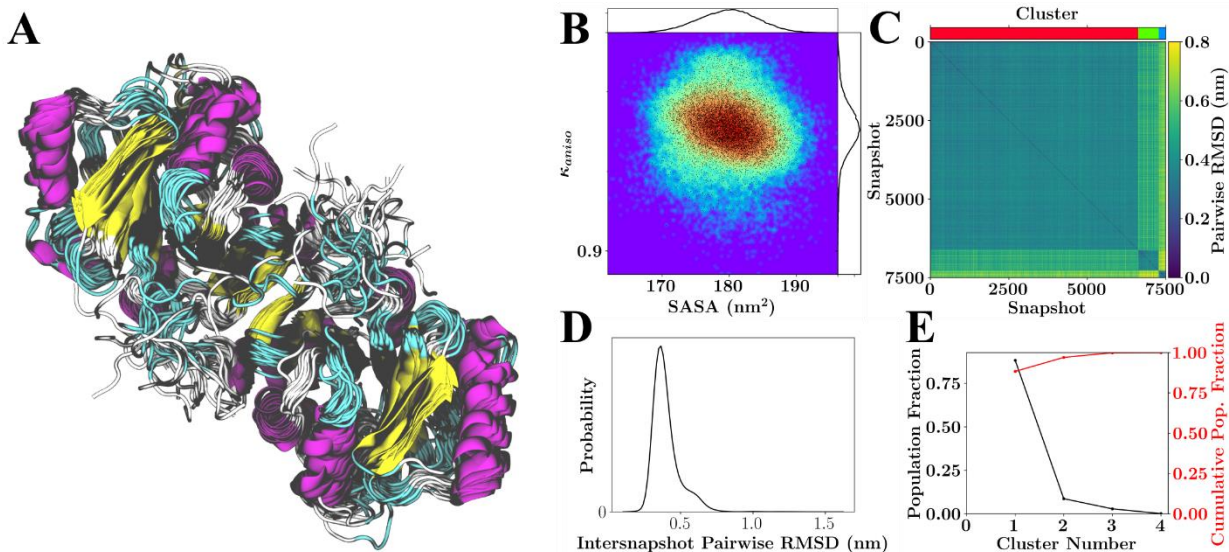
SI Figure 12. General results for simulations of NSP15 hexamer. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



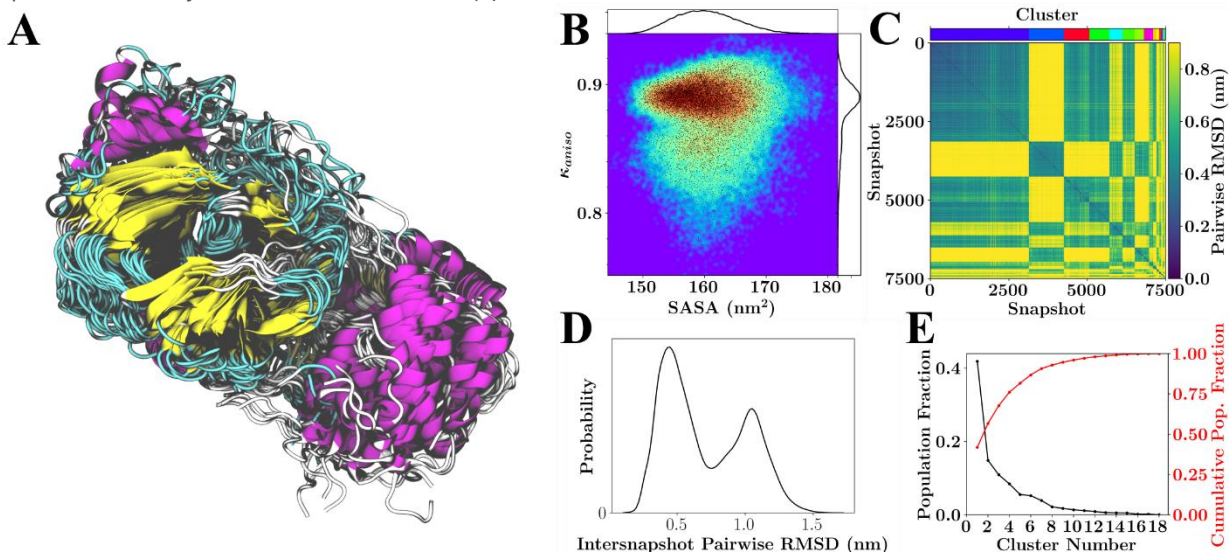
SI Figure 13. General results for simulations of NSP9 monomer. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



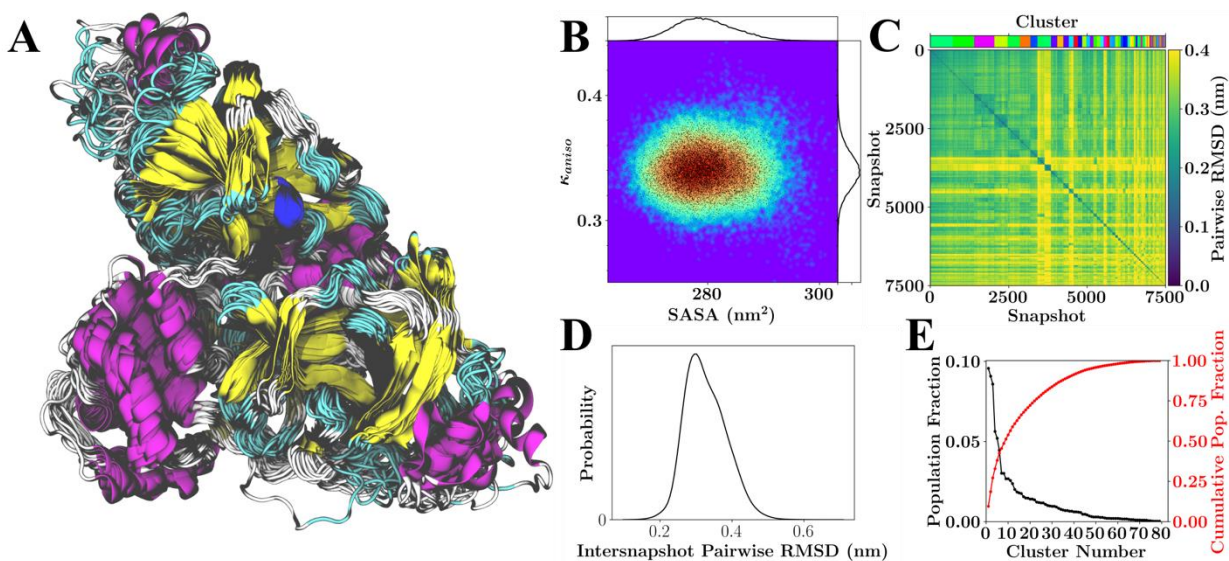
SI Figure 14. General results for simulations of NSP9 Dimer. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



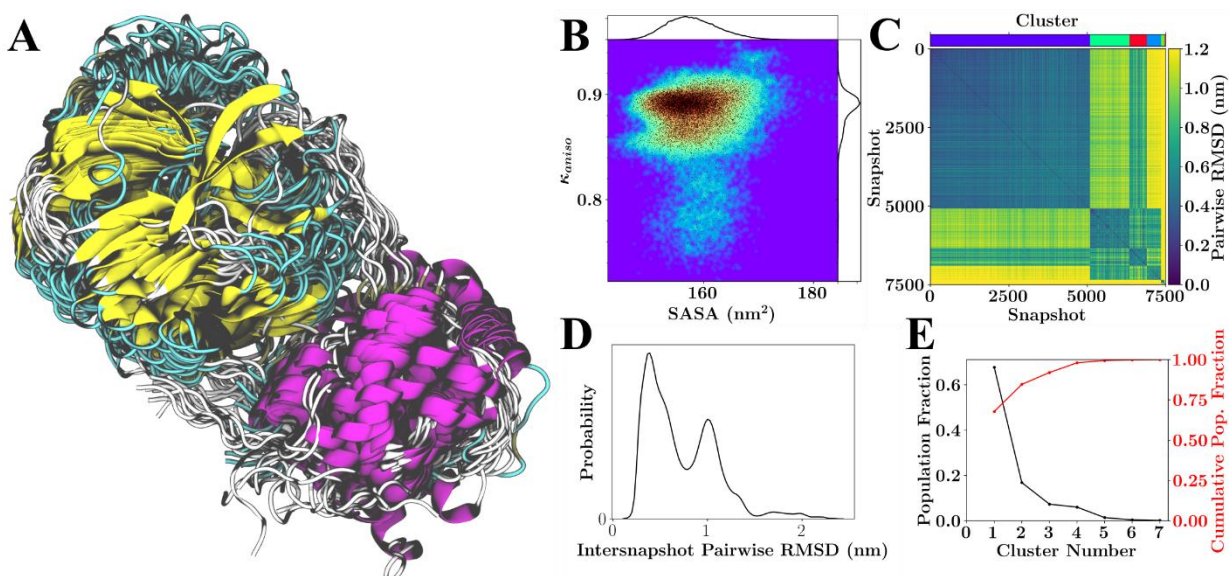
SI Figure 15. General results for simulations of NSP3 Phosphatase Domain asymmetric unit. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



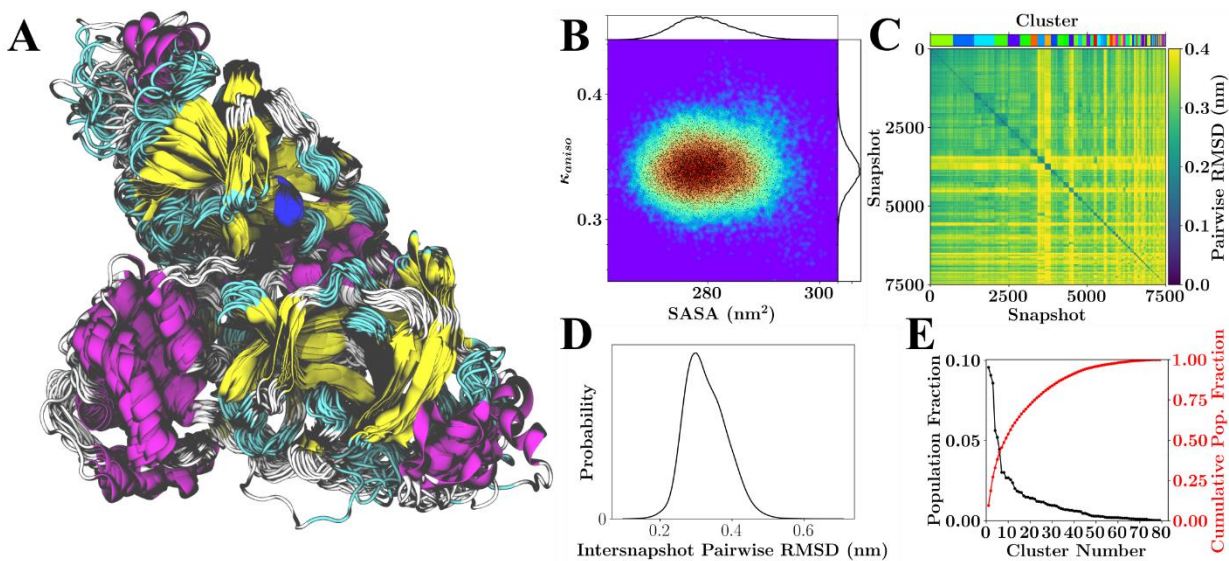
SI Figure 16. General results for simulations of MPro monomer with all HIS treated as HID (Default Output from CHARMM-GUI). (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



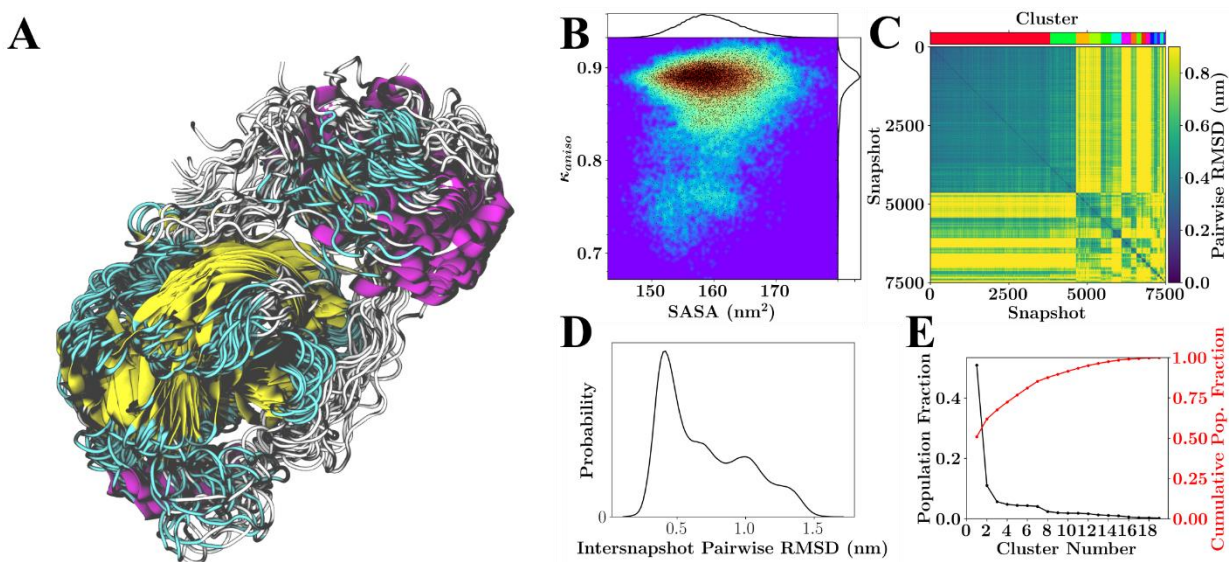
SI Figure 17. General results for simulations of MPro Dimer with all HIS treated as HID (Default Output from CHARMM-GUI). (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



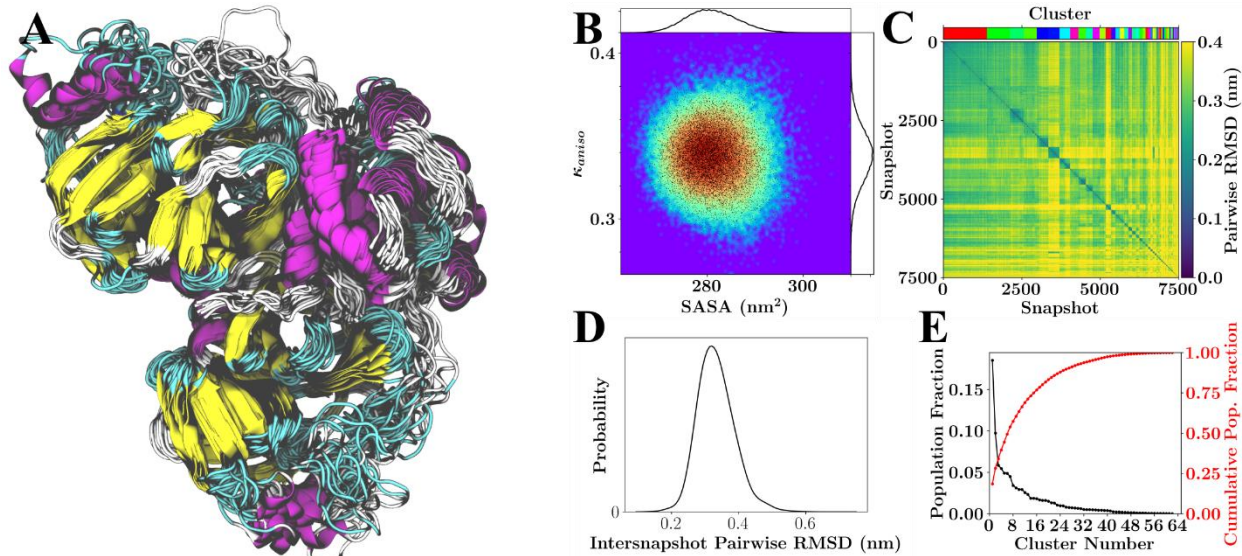
SI Figure 18. General results for simulations of MPro monomer with the following protonation states HSD41, HSD64, HSD80, HSE163, HSE164, HSE172, & HSE246. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



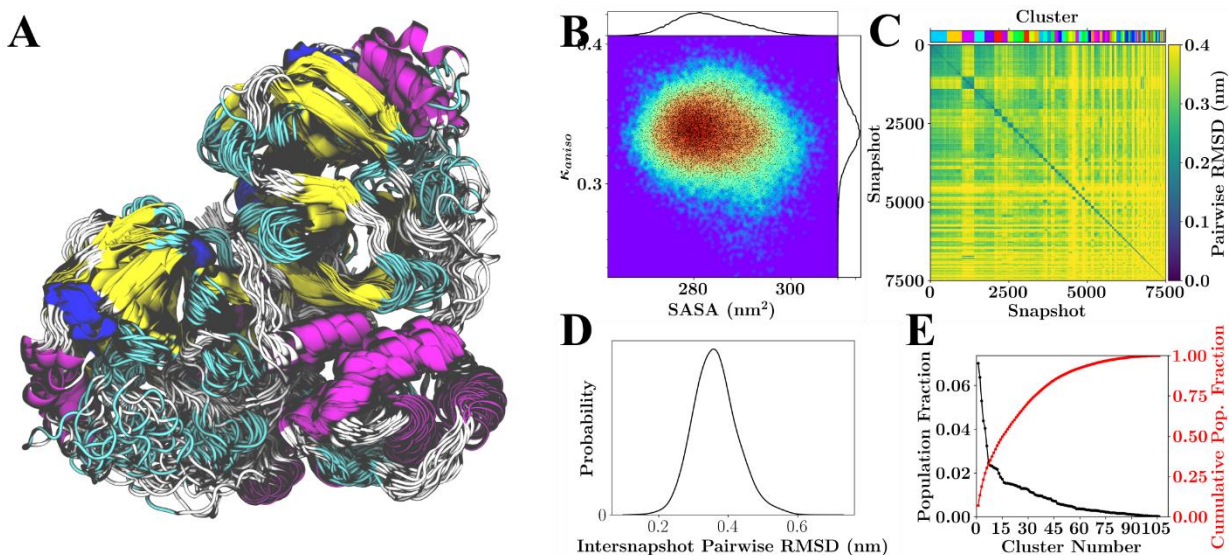
SI Figure 19. General results for simulations of MPro dimer with the following protonation states HSD41, HSD64, HSD80, HSE163, HSE164, HSE172, & HSE246. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



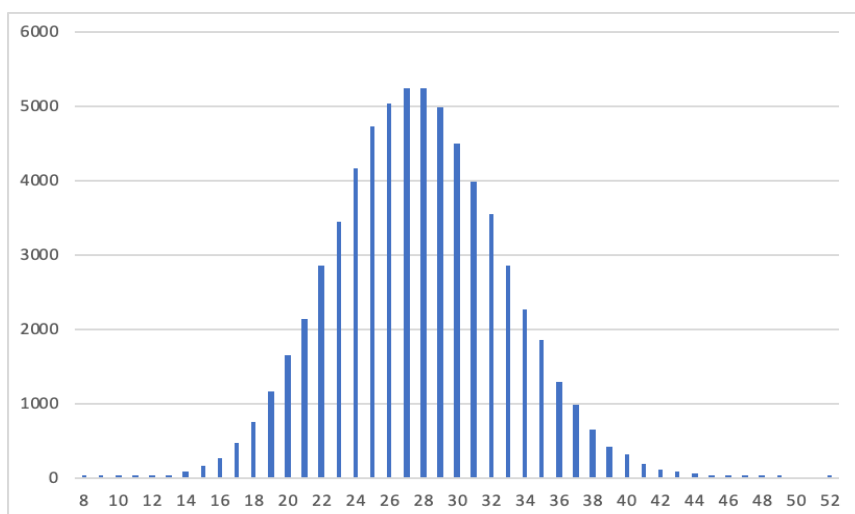
SI Figure 20. General results for simulations of MPro monomer with the following protonation states HSE41, HSD64, HSD80, HSE163, HSE164, HSE172, & HSE246. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



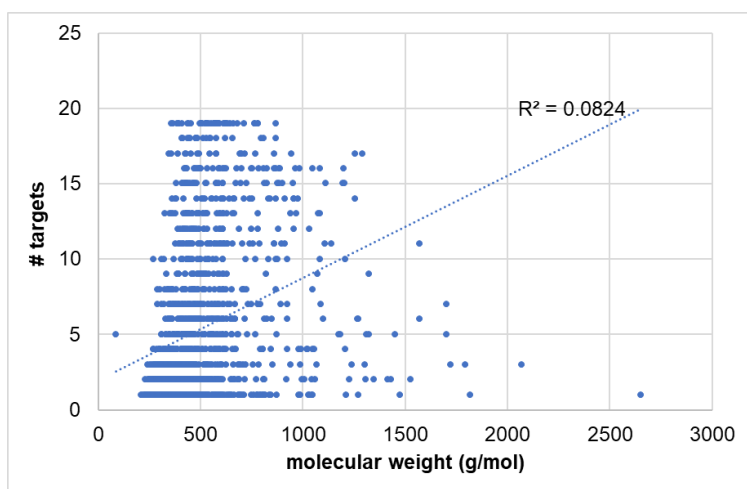
SI Figure 21. General results for simulations of MPro monomer with the following protonation states HSE41, HSD64, HSD80, HSE163, HSE164, HSE172, & HSE246. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



SI Figure 22. General results for simulations of the 'charged' MPro dimer with the following protonation states HSD41, HSP64, HSP80, HSP163, HSP164, HSE172, HSE246. (A) Overlay of 26 RMSD aligned structures from the lowest temperature replicate spanning the 750 ns of sampling. (B) Population distribution for shape anisotropy (κ) and solvent accessible surface area (SASA), with redder colors indicating greater occupancy of these kappa-SASA combinations. The distributions are also reflected by one-dimensional histograms above and to the right of the plot, and black dots within the population distribution, which represent position information for 10% of the total snapshots considered. (C) Pairwise RMSD clustering for the lowest temperature replica, with the snapshots ordered according to their cluster. The clusters in this instance were defined using a cutoff of half the maximum RMSD observed within the simulation and are labeled according to color with a color-bar for reference located above the plot. (D) Pairwise RMSD distribution across all snapshots. (E) Population statistics for the clusters introduced in (C).



SI Figure 23. Distribution of expected numbers of identical compounds for two random hits from 500 compounds out of the complete smaller database.



SI Figure 24. Correlation of ligand molecular weight and the number of targets a particular ligand was in the top 500 ranked compounds.