# Al-Driven Multiscale Simulations Illuminate Mechanisms of SARS-CoV-2 Spike Dynamics

Lorenzo Casalino<sup>1†</sup>, Abigail Dommer<sup>1†</sup>, Zied Gaieb<sup>1†</sup>, Emilia P. Barros<sup>1</sup>, Terra Sztain<sup>1</sup>, Surl-Hee Ahn<sup>1</sup>, Anda Trifan<sup>2,3</sup>, Alexander Brace<sup>2</sup>, Anthony Bogetti<sup>4</sup>, Heng Ma<sup>2</sup>, Hyungro Lee<sup>5</sup>, Matteo Turilli<sup>5</sup>, Syma Khalid<sup>6</sup>, Lillian Chong<sup>4</sup>, Carlos Simmerling<sup>7</sup>, David J. Hardy<sup>3</sup>, Julio D. C. Maia<sup>3</sup>, James C. Phillips<sup>3</sup>, Thorsten Kurth<sup>8</sup>, Abraham Stern<sup>8</sup>, Lei Huang<sup>9</sup>, John McCalpin<sup>9</sup>, Mahidhar Tatineni<sup>10</sup>, Tom Gibbs<sup>8</sup>, John E. Stone<sup>3</sup>, Shantenu Jha<sup>5</sup>, Arvind Ramanathan<sup>2\*</sup>, Rommie E. Amaro<sup>1\*</sup> <sup>1</sup>University of California San Diego, <sup>2</sup>Argonne National Lab, <sup>3</sup>University of Illinois at Urbana-Champaign, <sup>4</sup>University of Pittsburgh, <sup>5</sup>Rutgers University & Brookhaven National Lab, <sup>6</sup>University of Southampton, <sup>7</sup> Stony Brook University, <sup>8</sup>NVIDIA Corporation, <sup>9</sup>Texas Advanced Computing Center, <sup>10</sup>San Diego Supercomputing Center, <sup>†</sup>Joint first authors, \*Contact authors: ramaro@ucsd.edu, ramanathana@anl.gov

# **ABSTRACT**

We develop a generalizable AI-driven workflow that leverages heterogeneous HPC resources to explore the time-dependent dynamics of molecular systems. We use this workflow to investigate the mechanisms of infectivity of the SARS-CoV-2 spike protein, the main viral infection machinery. Our workflow enables more efficient investigation of spike dynamics in a variety of complex environments, including within a complete SARS-CoV-2 viral envelope simulation, which contains 305 million atoms and shows strong scaling on ORNL Summit using NAMD. We present several novel scientific discoveries, including the elucidation of the spike's full glycan shield, the role of spike glycans in modulating the infectivity of the virus, and the characterization of the flexible interactions between the spike and the human ACE2 receptor. We also demonstrate how AI can accelerate conformational sampling across different systems and pave the way for the future application of such methods to additional studies in SARS-CoV-2 and other molecular systems.

### **KEYWORDS**

molecular dynamics, deep learning, multiscale simulation, weighted ensemble, computational virology, SARS-CoV-2, COVID19, HPC, GPU, AI

#### **ACM Reference Format:**

Lorenzo Casalino<sup>1†</sup>, Abigail Dommer<sup>1†</sup>, Zied Gaieb<sup>1†</sup>, Emilia P. Barros<sup>1</sup>, Terra Sztain<sup>1</sup>, Surl-Hee Ahn<sup>1</sup>, Anda Trifan<sup>2,3</sup>, Alexander Brace<sup>2</sup>, Anthony Bogetti<sup>4</sup>, Heng Ma<sup>2</sup>, Hyungro Lee<sup>5</sup>, Matteo Turilli<sup>5</sup>, Syma Khalid<sup>6</sup>, Lillian Chong<sup>4</sup>, Carlos Simmerling<sup>7</sup>, David J. Hardy<sup>3</sup>, Julio D. C. Maia<sup>3</sup>, James C. Phillips<sup>3</sup>, Thorsten Kurth<sup>8</sup>, Abraham Stern<sup>8</sup>, Lei Huang<sup>9</sup>, John McCalpin<sup>9</sup>, Mahidhar Tatineni<sup>10</sup>, Tom Gibbs<sup>8</sup>, John E. Stone<sup>3</sup>, Shantenu Jha<sup>5</sup>, Arvind Ramanathan<sup>2\*</sup>, Rommie E. Amaro<sup>1\*</sup>. 2020. AI-Driven Multiscale Simulations Illuminate Mechanisms of SARS-CoV-2 Spike Dynamics. In *Supercomputing* 

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Supercomputing '20, November 16–19, 2020, Virtual © 2020 Association for Computing Machinery. ACM ISBN ISBN...\$15.00 https://doi.org/finalDOI '20: International Conference for High Performance Computing, Networking, Storage, and Analysis. ACM, New York, NY, USA, 14 pages. https://doi.org/finalDOI

# 1 JUSTIFICATION

We:

- develop an AI-driven multiscale simulation framework to interrogate SARS-CoV-2 spike dynamics,
- reveal the spike's full glycan shield and discover that glycans play an active role in infection, and
- achieve new high watermarks for classical MD simulation of viruses (305 million atoms) and the weighted ensemble method (600,000 atoms).

# 2 PERFORMANCE ATTRIBUTES

Our Submission
calability, Time-to-solution
Explicit, Deep Learning
nole application including I/O
Mixed Precision
Measured on full system
dware performance counters,
Application timers,
Performance Modeling

# 3 OVERVIEW OF THE PROBLEM

The SARS-CoV-2 virus is the causative agent of COVID19, a world-wide pandemic that has infected over 35 million people and killed over one million. As such it is the subject of intense scientific investigations. Researchers are interested in understanding the structure and function of the proteins that constitute the virus, as this knowledge aids in the understanding of transmission, infectivity, and potential therapeutics.

A number of experimental methods, including x-ray crystallography, cryoelectron (cryo-EM) microscopy, and cryo-EM tomography are able to inform on the structure of viral proteins and the other (e.g., host cell) proteins with which the virus interacts. Such structural information is vital to our understanding of these molecular machines, however, there are limits to what experiments can tell us.

Casalino et al.

For example, achieving high resolution structures typically comes at the expense of dynamics: flexible parts of the proteins (e.g., loops) are often not resolved, or frequently not even included in the experimental construct. Glycans, the sugar-like structures that decorate viral surface proteins, are particularly flexible and thus experimental techniques are currently unable to provide detailed views into their structure and function beyond a few basic units. Additionally, these experiments can resolve static snapshots, perhaps catching different states of the protein, but they are unable to elucidate the thermodynamic and kinetic relationships between such states.

In addition to the rich structural datasets, researchers have used a variety of proteomic, glycomic, and other methods to determine detailed information about particular aspects of the virus. In one example, deep sequencing methods have informed on the functional implications of mutations in a key part of the viral spike protein [57]. In others, mass spectrometry approaches have provided information about the particular composition of the glycans at particular sites on the viral protein [54, 69]. These data are each valuable in their own right but exist as disparate islands of knowledge. Thus there is a need to integrate these datasets into cohesive models, such that the fluctuations of the viral particle and its components that cause its infectivity can be understood.

In this work, we used all-atom molecular dynamics (MD) simulations to combine, augment, and extend available experimental datasets in order to interrogate the structure, dynamics, and function of the SARS-CoV-2 spike protein (Fig. 1). The spike protein is considered the main infection machinery of the virus because it is the only glycoprotein on the surface of the virus and it is the molecular machine that interacts with the human host cell receptor, ACE2, at the initial step of infection. We have developed MD simulations of the spike protein at three distinct scales, where each system (and scale) is informative, extensive, and scientifically valuable in its own right (as will be discussed). This includes the construction and simulation of the SARS-CoV-2 viral envelope that contains 305 million atoms, and is thus among one of the largest and most complex biological systems ever simulated (Fig. 1A). We employ both conventional MD as well as the weighted ensemble enhanced sampling approach (which again breaks new ground in terms of applicable system size). We then collectively couple these breakthrough simulations with artificial intelligence (AI) based methods as part of an integrated workflow that transfers knowledge gained at one scale to 'drive' (enhance) sampling at another.

An additional significant challenge faced in bringing this work to fruition is that it pushes the boundaries of several fields simultaneously, including biology, physics, chemistry, mathematics, and computer science. It is intersectional in nature, and requires the collective work of and effective communication among experts in each of these fields to construct, simulate, and analyze such systems - all while optimizing code performance to accelerate scientific discovery against SARS-CoV-2.

Our work has brought HPC to bear to provide unprecedented detail and atomic-level understanding of virus particles and how they infect human cells. Our efforts shed light on many aspects of the spike dynamics and function that are currently inaccessible with experiment, and have provided a number of experimentally testable hypotheses - some of which have already been experimentally validated. By doing so, we provide new understandings for

vaccine and therapeutic development, inform on basic mechanisms of viral infection, push technological and methodological limits for molecular simulation, and bring supercomputing to the forefront in the fight against COVID19.

#### 3.1 Methods

Full-length, fully-glycosylated spike protein. In this work, we built two full-length glycosylated all-atom models of the SARS-CoV-2 S protein in both closed and open states, fully detailed in Casalino et al [10]. The two all-atom models were built starting from the cryo-EM structures of the spike in the open state (PDB ID: 6VSB [70]), where one receptor binding domain (RBD) is in the "up" conformation, and in the closed state, bearing instead three RBDs in the "down" conformation (PDB ID: 6VXX [66]). Given that the experimental cryo-EM structures were incomplete, the remaining parts, namely (i) the missing loops within the head (residues 16-1141), (ii) the stalk (residues 1141-1234) and (iii) the cytosolic tail (residues 1235-1273), were modelled using MODELLER [51] and I-TASSER [79]. The resulting full-length all-atom constructs were subsequently N-/O-glycosylated using the Glycan Reader & Modeler tool [24] integrated into Glycan Reader [25] in CHARMM-GUI [38]. Importantly, an asymmetric glycoprofile was generated (e.g., not specular across monomers) taking into account the N-/Oglycans heterogeneity as described in the available glycoanalytic data [54, 69]. The two glycosylated systems were embedded into their physiological environment composed of an ERGIC-like lipid bilayer [11, 65] built using CHARMM-GUI [24, 72], explicit TIP3P water molecules [26], and neutralizing chloride and sodium ions at 150 mM concentration, generating two final systems each tallying ~1.7 million atoms. Using CHARMM36 all-atom additive force fields [19, 21] and NAMD 2.14 [42], the systems were initially relaxed through a series of minimization, melting (for the membrane), and equilibration cycles. The equilibrated systems were then subjected to multiple replicas of all-atom MD simulation production runs of the open (6x) and closed (3x) systems on the NSF Frontera computing system at the Texas Advanced Computing Center (TACC). A cumulative extensive sampling of  $\sim$ 4.2 and  $\sim$ 1.7  $\mu s$  was attained for the open and closed systems, respectively. Additionally, a third, mutant system bearing N165A and N234A mutations was built from the open system in order to delete the N-linked glycans and delineate their structural role in the RBD dynamics. This system was also simulated for  $\sim$ 4.2 µs in 6 replicas [10].

ACE2-RBD complex MD simulations. The model of the ACE2-RBD complex was based on cryo-EM structure trapping ACE2 as a homodimer co-complexed with two RBDs and B0AT1 transporter (PDB ID 6M17 [73]). Upon removal of B0AT1, ACE2 missing residues at the C terminal end were modeled using I-TASSER [79], whereas those missing at the N terminal end were taken from 6M0J and properly positioned upon alignment of the N terminal helix. Zinc sites including the ions and the coordinating residues were copied from 1R42. The construct was fully N-/O-glycosylated using CHARMM-GUI tools [24, 25, 38] for glycan modeling, reproducing the glycan heterogeneity for ACE2 and RBD reported in the available glycoanalytic data [53, 62, 81]. Similarly, the apo ACE2 homo-dimer was also built upon removal of the RBDs from the holo construct. The glycosylated models were embedded into separate lipid patches

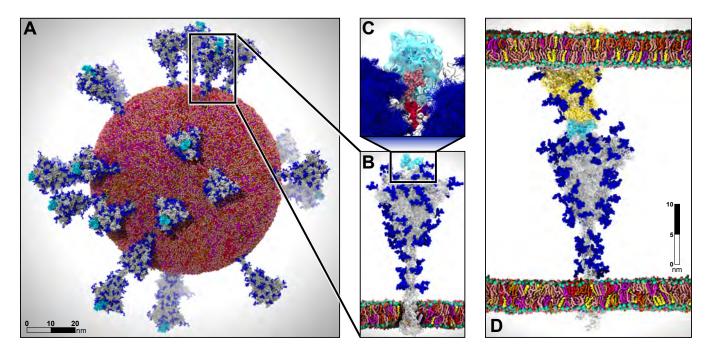


Figure 1: Multiscale modeling of SARS-CoV-2. A) All-atom model of the SARS-CoV-2 viral envelope (305 M atoms), including 24 spike proteins (colored in gray) in both the open (16) and closed states (8). The RBDs in the "up" state are highlighted in cyan) N-/O-Glycans are shown in blue. Water molecules and ions have been omitted for clarity. B) Full-length model of the glycosylated SARS-CoV-2 spike protein (gray surface) embedded into an ERGIC-like lipid bilayer (1.7 M atoms). RBD in the "up" state is highlighted in cyan. C) The glycan shield is shown by overlaying multiple conformations for each glycan collected at subsequent timesteps along the dynamics (blue bushlike representation). Highlighted in pink and red are two N-glycans (linked to N165 and N234, respectively) responsible for the modulation of the RBD dynamics, thus priming the virus for infection. The RBD "up" is depicted with a cyan surface. D) Two-parallel-membrane system of the spike-ACE2 complex (8.5 M atoms). The spike protein, embedded into an ERGIC-like membrane, is depicted with a gray transparent surface, whereas ACE2 is shown with a yellow transparent surface and it is embedded into a lipid bilayer mimicking the composition of mammalian cell membranes. Glycans are shown in blue, whereas water has been omitted for clarity. Visualizations were created in VMD using its custom GPU-accelerated ray tracing engine [23, 58–61].

with a composition mimicking that of mammalian cellular membranes [11, 65] and simulated in explicit water molecules at 150 mM ion concentration, affording two final systems of ~800,000 atoms each. MD simulations were performed using CHARMM36 all-atom additive force fields [19, 21] along with NAMD 2.14 [42]. The MD protocol was identical to that adopted for the simulation of the full-length spike and it is fully described in Casalino et al [10]. This work is fully detailed in Barros et al [5].

Weighted ensemble simulations of spike opening. The spike must undergo a large conformational change for activation and binding to ACE2 receptors, where the receptor binding domain transitions from the "down", or closed state to the "up," or open state [71]. Such conformational changes occur on biological timescales generally not accessible by classical molecular dynamics simulations [37]. To simulate the full unbiased path at atomic resolution, we used the weighted ensemble (WE) enhanced sampling method [22, 82]. Instead of running one single long simulation, the WE method runs many short simulations in parallel along the chosen reaction coordinates. The trajectories that rarely sample high energy regions are replicated, while the trajectories that frequently sample low energy

regions are merged, which makes sampling rare events computationally tractable and gives enhanced sampling. The trajectories also carry probabilities or weights, which are continuously updated, and there is no statistical bias added to the system. Hence, we are able to directly obtain both thermodynamic and kinetic properties from the WE simulations [78].

For this study, the closed model of the glycosylated spike from Casalino et al. [10], was used as the initial structure by only keeping the head domain. The WE simulations were run using the highly scalable WESTPA software [83], with the Amber GPU accelerated molecular dynamics engine [20, 52], version 18. Chamber [13] was used to convert CHARMM36 [19, 21] force fields and parameters from the system developed by Casalino et al. [10] into an Amber readable format. A TIP3P [27] water box with at least 10 Šbetween protein and box edges was used with 150 mM NaCl, leading the total number of atoms to 548,881. Amber minimization was carried out in two stages. First the solvent was minimized for 10,000 cycles with sugars and proteins restrained with a weight of 100 kcal/mol Ų, followed by unrestrained minimization for 100,000 cycles. Next the system was incrementally heated to 300 K over 300 ps. Equilibration

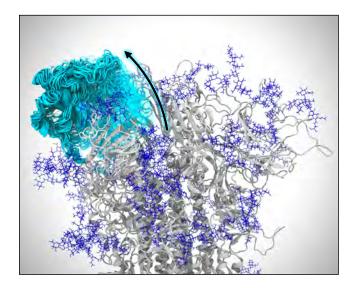


Figure 2: Opening of the spike protein. VMD visualization of weighted ensemble simulations shows the transition of the spike's RBD from the closed state to the open state. Many conformations of the RBD along its opening pathway are represented at the same time using cyan cartoons and a transparency gradient. Glycans appear as dark blue.

and production were carried out in 2 fs timesteps with SHAKE [49] constraints on non-polar hydrogens and NPT ensemble. Pressure and temperature were controlled with Monte Carlo barostat and Langevin thermostat with 1 ps-1 collision frequency. The particlemesh Ewald (PME) method was used with 10 Å cutoff for non-bonded interactions. The system was first equilibrated for 21 ns of conventional MD. The RMSD of the alpha carbons began to level off around 16 ns, and 24 structures were taken at regular intervals between 16 and 21 ns to use as equally weighted basis states for the WE simulation.

For each WE, tau was set to 100 ps of MD production followed by progress coordinate evaluation, and splitting / merging of walkers and updating weights, with a maximum of 8 walkers per bin. A two dimensional progress coordinate was defined by (i): the distance between the center of mass (COM) of the alpha carbons in the structured region of the spike helical core, and the alpha carbons in the four main beta sheets of the RBD (refers to RBD from chain A unless otherwise specified) and (ii): the RMSD of the alpha carbons in the four main beta sheets of the RBD to the initial structure (obtained from 1 ns equilibration). This simulation was run for 8.77 days on 80 P100 GPUs on Comet at SDSC collecting a comprehensive sampling of ~7.5  $\mu$ s, with bin spacing continuously monitored and adjusted to maximize sampling.

After extensive sampling of the RBD closed state, the second progress coordinate was changed to the RMSD of the alpha carbons in the four main beta sheets of the RBD compared to the final open structure, obtained from system 1, after 1 ns of equilibration carried out with identical methods as the closed structure described above, which was initially calculated as 11.5 Å. This allowed more efficient sampling of the transition to the open state by focusing

sampling on states which are closer in rotational or translational space to the final state, rather than sampling all conformations that are distinctly different from the closed state. Bin spacing was continuously monitored and adjusted to maximize traversing the RMSD coordinate. The full transition was confirmed when the RMSD coordinate reached below 6 Å and the RBD COM coordinate reached above 8.5 Å (Fig. 2). The simulation was stopped for analysis after 1099 iterations, upon running for 26.74 days on 100 V100 GPUs on Longhorn at TACC and harvesting  $\sim\!\!70.0~\mu s$ .

A second, independent WE simulation was conducted to determine if the findings of the initial simulation were reproducible, and to use the information on the free energy landscape of the successful transition in the first WE to inform bin spacing and target state definition to run an unsupervised simulation. After 19.64 days on 100 V100 GPUs on TACC Longhorn and  $\sim\!51.5~\mu s$  of comprehensive sampling, successful transitions to the open state were observed, as well as further open states, in which the RBD was observed to be peeling off of the spike core.

Two-parallel-membrane system of the spike-ACE2 complex. The SARS-CoV-2 virus gains entry into the host cell through a membrane fusion process taking place upon the recognition of the ACE2 receptors exposed on the host cell. This binding event triggers several, dramatic conformational changes within the spike protein, which becomes primed to pull the two membranes together for fusion, allowing the virus to pour the viral RNA into the host cell. In order to disentangle the mechanistic intricacies underlying this key process, we exploited the wealth of information obtained from the individual simulations described above to assemble an all-atom complex between the full-length spike and the ACE2 dimer. As a first step, equilibrated structures of the spike in the open state and of the ACE2-RBD complex were extracted from their respective individual simulations [5, 9]. Subsequently, the spike protein was superimposed onto the ACE2-RBD complex by aligning the spikes's RBD "up" with the RBD of the ACE2-RBD complex, allowing for a fairly vertical arrangement of the new construct. In order to preserve the best possible binding interface, the RBD of the spike was discarded, whereas the RBD from the ACE2-RBD complex was retained and linked to the rest of the spike. The spike-ACE2 complex was embedded into a double membrane system: the spike's transmembrane domain was inserted into a 330 Å × 330 Å ERGIC-like lipid bilayer, whereas for ACE2 a mammalian cellular membrane of the same dimension was used [11, 65]. The two membranes were kept parallel to each other, allowing the use of an orthorhombic box. In order to facilitate the water and ion exchange between the internal and external compartment, an outer-membrane-protein-G (OmpG) porin folded into a beta barrel was embedded into each membrane. The OmpG equilibrated model was obtained from Chen et al [12]. The generated two-membrane construct was solvated with explicit TIP3P water molecules, with the total height of the external water compartment matching the internal one exhibiting a value of 380 Å. Sodium and chloride ions were added at a concentration of 150 mM to neutralize the charge and reshuffled to balance the charge between the two compartments.

The composite system, counting 8,562,698 atoms with an orthorhombic box of 330 Å  $\times$  330 Å  $\times$  850 Å, was subjected to allatom MD simulation on the Summit computing system at ORNL

AI-Driven Multiscale Simulations Illuminate Mechanisms of SARS-CoV-2 Spike Dynamics

Supercomputing '20, November 16-19, 2020, Virtual

using NAMD 2.14 [42] and CHARMM36 all-atom additive force fields [19, 21]. Two cycles of conjugate gradient energy minimization and NPT pre-equilibration were conducted using a 2 fs timestep for a total of ~3 ns. During this phase, the ACE2 and spike proteins and the glycans were harmonically restrained at 5 kcal/mol, allowing for the relaxation of the two lipid bilayers, the OmpG porins, water molecules and ions within the context of the double membrane system. We remark that the two lipid patches were previously equilibrated, therefore not requiring a melting phase at this stage. The dimension of the cell in the xy plane was maintained constant while allowing fluctuation along the z axis. Upon this initial preequilibration phase, a ~17 ns NPT equilibration was performed by releasing all the restraints, preparing the system for production run. From this point, three replicas were run or a total of ~522 ns comprehensive simulation time. By using the trained AI learning model, three conformations were extracted from this set of simulations, each of them representing a starting point of a new replica with re-initialized velocities. A total of three additional simulations were therefore performed, collecting ~180 ns and bringing the total simulation time to  $\sim$ 702 ns.

SARS-CoV-2 viral envelope. The full-scale viral envelope was constructed using the LipidWrapper program (v1.2) previously developed and described by Durrant et al. [14]. A 350 Å  $\times$  350 Å lipid bilayer patch used as the pdb input was generated using CHARMM-GUI with an ERGIC-like lipid composition and an estimated area per lipid of 63 Å. An icospherical mesh with a 42.5 nm radius, in accordance with experimentally-observed CoV-2 radii, was exported as a collada file from Blender (v2.79b) and used as the surface file [31]. LipidWrapper was run in a Python 2.7 conda environment with lipid headgroup parameters "\_P,CHL1\_O3", a lipid clash cut-off of 1.0 Å, and filling holes enabled. The final bilayer pdb was solvated in a 110 nm cubic box using explicit TIP3P water molecules and neutralized with sodium and chloride ions to a concentration of 150 mM. The final system contained 76,134,149 atoms.

Since the LipidWrapper program operates via tessellation, lipid clash removal, and a subsequent lipid patching algorithm, the bilayer output attains a lower surface pressure than that of a bilayer of the same lipid composition at equilibrium [9]. Due to this artifact, as the bilayer equilibrates, the lipids undergo lateral compression resulting in the unwanted formation of pores. Thus, the envelope was subjected to multiple rounds of minimization, heating, equilibration, and patching until the appropriate equilibrium surface pressure was reached.

All-atom MD simulations were performed using NAMD 2.14 and CHARMM36 all-atom additive force fields. The conjugate-gradient energy minimization procedure included two phases in which the lipid headgroups were restrained with 100 and 10 kcal/mol weights, respectively, at 310K for 15,000 cycles each. The membrane was then melted by incremental heating from 25 K to 310 K over 300 ps prior to NPT equilibration. The equilibration sequentially released the harmonic restraints on the lipid headgroups from 100 to 0 kcal/mol over 0.5 ns. Following this sequence, the structure was visually evaluated to determine whether to continue equilibration or to proceed with pore patching. Most structures continued with unrestrained

equilibration for 4–26 ns prior to patching, with longer unrestrained equilibrations attributed to later, more stable envelopes.

Patching of the envelope was done by overlapping the initial LipidWrapper bilayer output with the newly-equilibrated envelope. All superimposed lipids within 2.0 Å of the equilibrated lipids were removed to eliminate clashes. Superimposed lipids within 4.0 Å of an equilibrated cholesterol molecule were also removed to eliminate ring penetrations. The patched system, with new lipids occupying the pores, was then re-solvated, neutralized, and subjected to the next round of minimization, heating, and equilibration.

After ten rounds of equilibration and patching, 24 spike proteins with glycans, 8 in the closed and 16 in the open state, were inserted randomly on the envelope using a house tcl script. A random placement algorithm was used in accordance with experimental microscopy imaging which has suggested that there is no obvious clustering of the spikes and no correlation between RBD state and location on the spike surface [31]. The number of spikes was selected based on experimental evidence reporting a concentration of 1000 spikes/nm<sup>2</sup> on the envelope [31]. The new structure containing spikes was re-solvated, neutralized, and processed to remove clashing lipids prior to further simulation. The resulting cubic solvent box was 146 nm per side and contained 304,780,149 atoms. The spike-inclusive envelope was then subjected to three more equilibration and patching sequences. The final virion used for all-atom MD production runs had a lipid envelope of 75 nm in diameter with a full virion diameter of 120 nm. The complete equilibration of the viral envelope totaled 41 ns on the TACC Frontera system and 75 ns on ORNL Summit. Full-scale viral envelope production simulations were performed on Summit for a total of 84 ns in an NPT ensemble at 310 K, with a PME cutoff of 12 Å for non-bonded interactions.

# 4 CURRENT STATE OF THE ART

### 4.1 Parallel molecular dynamics

NAMD [41] has been developed over more than two decades, with the goal of harnessing parallel computing to create a computational microscope [34, 55] enabling scientists to study the structure and function of large biomolecular complexes relevant to human health. NAMD uses adaptive, asynchronous, message-driven execution based on Charm++[28, 29]. It was one of the first scientific applications to make use of heterogeneous computing with GPUs [43], and it implements a wide variety of advanced features supporting state-of-the-art simulation methodologies. Continuing NAMD and Charm++ developments have brought improved work decomposition and distribution approaches and support for low overhead hardware-specific messaging layers, enabling NAMD to achieve greater scalability on larger parallel systems [32, 44]. NAMD incorporates a collective variables module supporting advanced biasing methods and a variety of in-situ analytical operations [16]. Simulation preparation, visualization, and post-hoc analysis are performed using both interactive and offline parallel VMD jobs [23, 59-61]. NAMD has previously been used to study viruses and large photosynthetic complexes on large capability-oriented and leadership class supercomputing platforms, enabling the high-fidelity determination of the HIV-1 capsid structure [80], the characterization of

https://www.blender.org/

<sup>2</sup>https://docs.anaconda.com/

Casalino et al.

substrate binding in influenza [15], and the structure and kinetics of light harvesting bacterial organelles [56].

# 4.2 Weighted Ensemble MD simulations

The weighted ensemble (WE) method is an enhanced sampling method for MD simulations that can be orders of magnitude more efficient than standard simulations in generating pathways and rate constants for rare-event processes. WE runs many short simulations in parallel, instead of one long simulation, and directly gives both thermodynamic and kinetic properties, which most enhanced sampling methods cannot do. The simulations go through "resampling" where simulations are merged for over-sampled regions and replicated for rare regions so that regions are continuously sampled regardless of energy barriers. The simulations also carry probabilities or "weights" that are continuously updated and no statistical bias is added to the system, so we are able to directly obtain both thermodynamic (e.g., free energy landscape) and kinetic (e.g., rates and pathways) properties from the simulation. In addition, the WE method is one of the few methods that can obtain continuous unbiased pathways between states, so this was the most suitable method for us to obtain and observe the closed to open transition for the spike system. Before the WE method was applied to the spike system under investigation here (about 600,000 atoms), the largest system used for the WE method was the barnase-barnstar complex (100,000 atoms)[50].

#### 4.3 AI-driven multiscale MD simulations

A number of approaches, including deep learning methods, have been developed for analysis of long timescale MD simulations [36]. These linear, non-linear, and hybrid ML approaches cluster the simulation data along a small number of latent dimensions to identify conformational transitions between states [6, 46]. Our group developed a deep learning approach, namely the variational autoencoder that uses convolutional filters on contact maps (from MD simulations) to analyze long time-scale simulation datasets and organize them into a small number of conformational states along biophysically relevant reaction coordinates [7]. We have used this approach to characterize protein conformational landscapes [48]. However, with the spike protein, the intrinsic size of the simulation posed a tremendous challenge in scaling our deep learning approaches to elucidate conformational states relevant to its function.

Recently, we extended our approach to adaptively run MD simulation ensembles to fold small proteins. This approach, called DeepDriveMD [35], successively learns which parts of the conformational landscape have been sampled sufficiently and initiates simulations from undersampled regions of the conformational landscape (that also constitute "interesting" features from a structural perspective of the protein). While a number of adaptive sampling techniques exist [2, 8, 30, 33, 47, 67, 68], including based on reinforcement learning methods [39], these techniques have been demonstrated on prototypical systems. In this paper, we utilize the deep learning framework to suggest additional points for sampling and do not necessarily use it in an adaptive manner to run MD simulations (mainly due to the limitations posed by the size of the system). However, extensions to our framework for enabling support of such large-scale systems are straightforward and further work will examine such large-scale simulations.

#### 5 INNOVATIONS REALIZED

# 5.1 Parallel molecular dynamics

Significant algorithmic improvements and performance optimizations have been required for NAMD to achieve high performance on the GPU-dense Summit architecture [1, 42, 58]. New CUDA kernels for computing the short-range non-bonded forces were developed that implement a "tile list" algorithm for decomposing the workload into lists of finer grained tiles that more fully and equitably distribute work across the larger SM (streaming multiprocessor) counts in modern NVIDIA GPUs. This new decomposition uses the symmetry in Newton's Third Law to eliminate redundant calculation without incurring additional warp-level synchronization [58]. CUDA kernels also were added to offload the calculation of the bonded force terms and non-bonded exclusions [1]. Although these terms account for a much smaller percentage of the work per step than that of the short-range non-bonded forces, NAMD performance on Summit benefits from further reduction of CPU workload. NAMD also benefits from the portable high-performance communication layer in Charm++ that communicates using the IBM PAMI (Parallel Active Messaging Interface) library, which improves performance by up to 20% over an MPI-based implementation [1, 32].

Additional improvements have benefited NAMD performance on Frontera. Recent developments in Charm++ now include support for the UCX (Unified Communication X) library which improves performance and scaling for Infiniband-based networks. Following the release of NAMD 2.14, a port of the CUDA tile list algorithm to Intel AVX-512 intrinsics was introduced, providing a 1.8× performance gain over the "Sky Lake" (SKX) builds of NAMD.

A significant innovation in NAMD and VMD has been the development of support for simulation of much larger system sizes, up to two billion atoms. Support for larger systems was developed and tested through all-atom modeling and simulation of the protocell as part of the ORNL CAAR (Center for Accelerated Application Readiness) program that provided early science access to the Summit system [42]. This work has greatly improved the performance and scalability of internal algorithms and data structures of NAMD and VMD to allow modeling of biomolecular systems beyond the previous practical limitation on the order of 250 million atoms. This work has redefined the practical simulation size limits in both NAMD and VMD and their associated file formats, added new analysis methods specifically oriented toward virology [17], and facilitates modeling of cell-scaled billion-atom assemblies, while making smaller modeling projects significantly more performant and streamlined than before [1, 17, 42, 58, 60].

# 5.2 Multiscale molecular dynamics simulations

Often referred to as "computational microscopy," MD simulations are a powerful class of methods that enable the exploration of complex biological systems, and their time-dependent dynamics, at the atomic level. The systems studied here push state of the art in both their size and complexity. The system containing a full-length, fully-glycosylated spike protein, embedded in a realistic viral membrane (with composition that mimics the endoplasmic reticulum) contains essentially all of the biological complexity known about the SARS-CoV-2 spike protein. The composite system contains  $\sim\!1.7$  million atoms and combines data from multiple cryoEM, glycomics, and

AI-Driven Multiscale Simulations Illuminate Mechanisms of SARS-CoV-2 Spike Dynamics

Supercomputing '20, November 16-19, 2020, Virtual

lipidomics datasets. The system was simulated with conventional MD out to microseconds in length, and several mutant systems were simulated and validated with independent experiments.

A related set of experiments utilizing the weighted ensemble method, an enhanced sampling technique, explored a truncated version of the spike protein ( $\sim$ 600,000 atoms with explicit solvent) in order to simulate an unbiased spike protein conformational transition from the closed to open state. This is the largest system, by an order of magnitude, that has been simulated using the WE method (biggest system until now was  $\sim$ 60,000 atoms). Using calculations optimized to efficiently make use of extensive GPU resources, we obtained several full, unbiased paths of the glycosylated spike receptor binding domain activation mechanism.

The second system increases the complexity by an order of magnitude by combining the spike system described above with a full-length, fully-glycosylated model of the ACE2 receptor bound into a host cell plasma membrane. This system represents the encounter complex between the spike and the ACE2 receptor, contains two parallel membranes of differing composition, has both the spike and ACE2 fully glycosylated, and forming a productive binding event at their interface. The composite system contains ~8.5 Million atoms with explicit water molecules and provides unseen views into the critical handshake that must occur between the spike protein and the ACE2 receptor to begin the infection cascade.

Our final system is of the SARS-CoV-2 viral envelope. This system incorporates 24 full-length, fully-glycosylated spike proteins into a viral membrane envelope of realistic (ER-like) composition, where the diameter of the viral membrane is  $\sim\!\!80\mathrm{nm}$  and the diameter of the virion, inclusive of spikes, is 146 nm. Until now, the largest system disclosed in a scientific publication was the influenza virus, which contained  $\sim\!160$  million atoms. The SARS-CoV-2 viral envelope simulation developed here contains a composite 305 million atoms, and thus breaks new ground for MD simulations of viruses in terms of particle count, size, and complexity.

Moreover, typical state of the art simulations are run in isolation, presenting each as a self-contained story. While we also do that for each of the systems presented here, we advance on state of the art by using an AI-driven workflow that drives simulation at one scale, with knowledge gained from a disparate scale. In this way, we are able to explore relevant phase space of the spike protein more efficiently and in environments of increasing complexity.

#### 5.3 Using AI for driving multiscale simulations

Using deep learning to characterize conformational states sampled in the SARS-CoV-2 spike simulations. MD simulations such as the ones described above generate tremendous amounts of data. For e.g., the simulations of the WE sampling of the spike protein's closed-to-open state generated over 100 terabytes of data. This imposes a heavy burden in terms of understanding the intrinsic latent dimensions along which large-scale conformational transitions can be characterized. A key challenge then is to use the raw simulation datasets (either coordinates, contact matrices, or other data collected as part of a standard MD runs) to cluster conformational states that have been currently sampled, to identify biologically relevant transitions between such states (e.g., open/closed states of spike), and suggest conformational states that may not be fully sampled to characterize these transitions [46].

To deal with the size and complexity of these simulation datasets, approaches that analyze 3D point clouds are more appropriate. Indeed, such approaches are becoming more commonly utilized for characterizing protein binding pockets and protein-ligand interactions. We posited that such representations based on the  $C^{\alpha}$  representation of protein structures could be viable to characterize largescale conformational changes within MD simulation trajectories. We leverage the 3D PointNet based [45] adversarial autoencoder (3D-AAE) developed by Zamorski and colleagues [77] to analyze the spike protein trajectories. In this work, we employ the chamfer distance based reconstruction loss and a Wasserstein [4] adversarial loss with gradient penalty [18] to stabilize training. The original PointNet backbone treats the point cloud as unordered, which is true for general point clouds. In our case however, the protein is essentially a 1D embedding into a 3D space. This allows us to define a canonical order of points, i.e. the order in which they appear in the chain of atoms. For that reason, we increase the size-1 1D convolutional encoder kernels from the original PointNet approach to larger kernels up to size 5. This allows the network to not only learn features solely based on distance, but also based on neighborhood in terms of position of each atom in the chain. We found that a 4layer encoder network with kernel sizes [5, 3, 3, 1, 1] and filter sizes [64, 128, 256, 256, 512] performs well for most tasks. A final dense layer maps the vectors into latent space with dimensionality 64. For the generator, we only use unit size kernels with filter dimensions [64, 128, 512, 1024, 3] respectively (the output filter size is always the dimensionality of the problem). The discriminator is a 5 layer fully connected network with layer widths [512, 512, 128, 64, 1].

The trajectories from the WE simulations were used to build a combined data set consisting of 130,880 examples. The point cloud data, representing the coordinates of the 3,375 backbone  $C^{\alpha}$  atoms of the protein, was randomly split into training (80%) and validation input (20%) and was used to train the 3D-AAE model for 100 epochs using a batch size of 32. The data was projected onto a latent space of 64 dimensions constrained by a gaussian prior distribution with a standard deviation of 0.2. The loss optimization was performed with the Adam optimizer, a variant of stochastic gradient descent, using a learning rate of 0.0001. We also added hyperparameters to scale individual components of the loss. The reconstruction loss was scaled by 0.5 and the gradient penalty by a factor of 10.

The embedding learned from the 3D-AAE model summarizes a latent space that is similar to variational autoencoders, except that 3D-AAEs tend to be more robust to outliers within the simulation data. The embeddings learned from the simulations allow us to cluster the conformations (in an unsupervised manner) based on their similarity in overall structure, which can be typically measured using quantities such as root-mean squared deviations (RMSD).

We trained the model using several combinations of hyperparameters, mainly varying learning rate, batch size and latent dimension. For visualizing and assessing the quality of the model in terms latent space structure, we computed t-SNE [64] dimensionality reductions on the high-dimensional embeddings from the validation set. A good model should generate clusters with respect to relevant biophysical observables not used in the training process. Therefore, we painted the t-SNE plot with the root mean squared deviation (RMSD) of each structure to the starting conformation and observed intelligible clustering of RMSD values. We tested this model on a

Casalino et al.

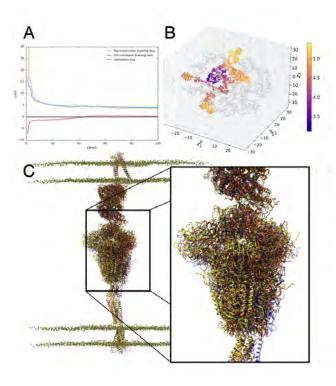


Figure 3: 3D-AAE training and test results. A) The loss progression for reconstruction, discriminator and validation loss over 100 epochs. B) The t-SNE plot visualization of the reduced latent space, with training embeddings represented in grey and test examples represented in color over the range of RMSD values. Outliers identified in the outlier detection stage are represented with an outlined diamond. C) VMD visualization of outlier structures (yellow, orange, dark orange) aligned and compared to the starting structure (blue).

set of trajectories from the full scale spike-ACE2 system, using the same atom selection (3,375  $C^{\alpha}$  atoms) as the corresponding WE spike protein. We subsequently performed outlier detection using the local outlier factor (LOF) algorithm, which uses distance from neighboring points to identify anomalous data. The goal of the outlier detection step is to identify conformations of the protein that are most distinct from the starting structure, in order to story board important events in the transition of the protein from an open to closed conformation. Although the number of outlier conformations detected can be a parameter that the end-user can specify, we selected 20 outlier conformations, based on the extreme LOF scores. These conformations were visualized in VMD [23, 58], and further analyzed using tilt angles of the stalk and the RBD. The final selection included 3 structures which were used as the starting conformations for the next set of simulations. These 'outlier' conformers are cycled through additional MD simulations that are driven by the ML-methods.

Table 1: NAMD AVX-512 FP operation breakdown.

FP Instr.	Ops	% total	FP Instr.	Ops	% total
DblScalar	4.99e16	26.9%	SglScalar	2.09e15	1.1%
Dbl128b	6.86e15	3.7%	Sgl128b	3.61e15	1.9%
Dbl256b	1.06e17	57.1%	Sgl256b	1.18e16	6.3%
Dbl512b	4.96e15	2.7%	Sgl512b	3.43e14	0.2%

# 6 HOW PERFORMANCE WAS MEASURED

#### 6.1 3D-AAE

Since this application dominantly utilizes the GPU, we do not need to profile CPU FLOPs. Instead, we measure FLOPs for all precisions using the methodology explained in [74] with the NVIDIA NSight Compute 2020 GPU profiling tool. We collect floating point instructions of relevant flavors (i.e. adds, mults, fmas (fused multiply adds) and tensor core operations for FP16, FP32 and FP64) and multiply those with weighting factors of {1, 1, 2, 512} respectively in order to transform those into FLOP counts. The sum of all these values for all precisions will yield our overall mixed precision FLOP count. To exclude FLOPs occuring during initialization and shutdown, we wrap the training iteration loop into start/stop profiler hooks provided by the NVIDIA CuPy Python package.<sup>3</sup>

#### 6.2 NAMD

NAMD performance metrics were collected on TACC Frontera, using the Intel msr-tools utilities, with NAMD 2.14 with added Intel AVX-512 support. FLOP counts were measured for each NAMD simulation with runs of two different step counts. The results of the two simulation lengths were subtracted to eliminate NAMD startup operations, yielding an accurate estimate of the marginal FLOPs per step for a continuing simulation [40].

FLOP counts were obtained by reading the hardware performance counters on all CPU cores on all nodes, using the rdmsr utility from msr-tools. At the beginning of each job, the "TACC stats" system programs the core performance counters to count the 8 sub-events of the Intel FP\_ARITH\_INST\_RETIRED. Counter values are summed among the 56 cores in each node, and ultimately among each node. Each node-summed counter value is scaled by the nominal SIMD-width of the floating point instruction being counted and the 8 classes are added together to provide the total FLOP count per node. The hardware counters do not take masked SIMD instructions into account. SIMD lanes that are masked-out still contribute to the total FLOPs, however static analysis of the AVX-512-enabled NAMD executable showed that only 3.7% of FMA instructions were masked.

A breakdown of floating point instruction execution frequency for the AVX-512 build of NAMD across 2048 nodes is shown in Table 1. For CPU versions of NAMD, arithmetic is performed in double precision, except for single-precision PME long-range electrostatics calculations and associated FFTs. In the GPU-accelerated NAMD on Summit, single-precision arithmetic is used for both PME and

<sup>3</sup>https://cupy.dev/

<sup>4</sup>https://github.com/intel/msr-tools

 $<sup>^5</sup>$ https://github.com/TACC/tacc\_stats

AI-Driven Multiscale Simulations Illuminate Mechanisms of SARS-CoV-2 Spike Dynamics

Supercomputing '20, November 16-19, 2020, Virtual

Table 2: 3D-AAE training performance on one V100 GPU.

Latent Dimensions	Peak TFLOP/s	Sustained TFLOP/s
32	2.96	0.97
64	3.16	2.28
128	3.13	0.91

also for short-range non-bonded force calculations, significantly increasing the fraction of single-precision instructions, at the cost of requiring a mixed-precision patch-center-based atomic coordinate representation to maintain full force calculation precision [42, 58].

# 7 PERFORMANCE RESULTS

# 7.1 3D-AAE training performance

We used the aforementioned recipe for GPU profiling to determine the performance for the 3D-AAE training. We measure the FLOP counts individually for 2 training and 1 validation steps for a batch size of 32. The latent dimension of the model is a free hyperparameter and affects the FLOP count. We trained three models with latent dimensions [32, 64, 128] in order to determine an optimal model for the task and thus we profile and report numbers for all of those. All models were trained for 100 epochs with batch size 32 on a single V100 GPU each. As mentioned above, the train/valdiation dataset split is 80%/20% and we do one validation pass after each training epoch. Thus, we can assume that this fraction translates directly into the FLOP counts for these alternating two stages. Our sustained performance numbers are computed using this weighted FLOP count average and the total run time. In order to determine peak performance, we compute the instantaneous FLOP rate for the fastest batch during training. Note that the 3D-AAE does exclusively use float (FP32) precision. The performance results are summarized in table 2. Although the model is dense linear algebra heavy, it is also rather lightweight so it cannot utilize the full GPU and thus only delivering 20% of theoretical peak performance.

As expected, the peak performance is very consistent between the runs. The big difference in sustained performance between latent dim 64 and the other two models is that the frequency for computing the t-SNE was significantly reduced, i.e. from every epoch to every 5th. The t-SNE computation and plotting happens after each validation in a background thread on the CPU, but the training epochs can be much shorter than the t-SNE time. In that case, the training will stall until the previous t-SNE has completed. Evidently, decreasing the t-SNE frequency reduces that overhead significantly. We expect that the other models would perform similarly if we would have enabled this optimization for those runs as well. The remaining difference in peak vs. sustained performance can be explained by other overhead, e.g. storing embedding vectors, model checkpoints and the initial scaffolding phase. Furthermore, it includes the less FLOP-intensive validation phase whereas the peak estimate is obtained from the FLOP-heavy training phase.

#### 7.2 NAMD simulation performance

Low-level NAMD performance measurements were made on the TACC Frontera system, to establish baseline counts of FLOPs per timestep for the four different biomolecular systems simulated as

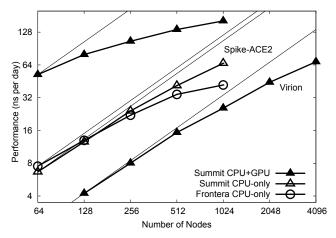


Figure 4: NAMD scaling on Summit and Frontera for 8.5Matom spike-ACE2 complex (upper lines) and 305M-atom virion (lower line). Thin lines indicate linear scaling.

Table 3: NAMD simulation floating point ops per timestep.

NAMD Simulation	Atoms	FLOPS/step
ACE2-RBD complex	800k	21.57 GFLOPS/step
Single Spike	1.7M	47.96 GFLOPS/step
Spike-ACE2 complex	8.5M	243.7 GFLOPS/step
SARS-CoV-2 virion	305M	8.3511 TFLOPS/step

Table 4: NAMD performance: 8.5M-atom Spike-ACE2.

Nodes	Frontera	Summit	Summit
	CPU-only	CPU-only	CPU + GPU
64	7.52 ns/day	6.67 ns/day	52.15 ns/day
128	13.00 ns/day	12.59 ns/day	79.68 ns/day
256	22.09 ns/day	24.19 ns/day	105.54 ns/day
512	34.32 ns/day	41.31 ns/day	135.31 ns/day
1024	41.88 ns/day	66.31 ns/day	162.22 ns/day

Table 5: NAMD performance: 305M-atom virion.

Nodes	Summit	Speedup	Efficiency
	CPU + GPU		
128	4.23 ns/day	~1.0×	~100%
256	8.02 ns/day	1.9×	95%
512	15.32 ns/day	3.6×	91%
1024	25.66 ns/day	6.1×	75%
2048	44.27 ns/day	10.5×	65%
4096	68.36 ns/day	16.2×	51%

part of this work, summarized in Table 3, with the breakdown of CPU FLOPs described in Table 1. Sustained NAMD performance measurements were obtained using built-in application timers over long production science runs of several hours, including all I/O, and reported in units of nanoseconds per day of simulation time. NAMD sustained simulation performance for the spike-ACE2 complex is summarized for the TACC Frontera and ORNL Summit systems

Table 6: Peak NAMD FLOP rates, ORNL Summit

NAMD Simulation	Atoms	Nodes	Sim rate	Performance
Spike-ACE2 complex	8.5M	1024	162 ns/day	229 TFLOP/s
SARS-CoV-2 virion	305M	4096	68 ns/day	3.06 PFLOP/s

in Table 4 and Fig. 4. NAMD sustained simulation performance, parallel speedup, and scaling efficiency are reported for the full SARS-CoV-2 virion in Table 5. Peak NAMD mixed-precision FLOP rates on ORNL Summit are estimated in Table 6 by combining sustained performance measurements with FLOPs/timestep measurements

# 8 IMPLICATIONS

Our major scientific achievements are:

- (1) We characterize for the first time the glycan shield of the full-length SARS-CoV-2 spike protein (including the stalk), and find that two N-glycans linked to N165 and N234 have a functional role in modulating the dynamics of the spike's RBD. This unprecedented finding establishes a major new role of glycans in this system as playing an active role in infection, beyond shielding (Fig. 1C) [10].
- (2) We discover that the human ACE2 receptor has a flexible hinge in the linker region near the membrane that enables it to undergo exceptionally large angular motions relative to the plane of the membrane. We predict this flexibility will aid forming productive complexes with the spike protein and may serve as a mechanical energy source during the cell fusion process [5].
- (3) We openly share our models, methods, and data, making them freely available to the scientific community. We are committed to the shared set of principles outlined in Ref. [3]: depositing findings as preprints in advance of formal peer review, making available our models at the time of deposition into a preprint server [5], and releasing the full datasets upon peer review [10]. By doing so, the reproducibility and robustness of our findings and methods are enhanced, and the scientific findings from our simulations are amplified through reuse by others.
- (4) We describe for the first time unbiased pathways for the full closed-to-open transition of the spike's RBD (Fig. 2), where knowledge of this pathway has the potential to inform on mechanisms of viral infection as well as potentially aid in the discovery of novel druggable pockets within the spike. Our work set a new milestone for the use of the weighted ensemble method in biomolecular simulation, increasing applicable system size by an order or magnitude over current state of the art.
- (5) We characterize the spike's flexibility in the context of ACE2 binding. One of the most important properties of the spike protein is its intrinsic flexibility, a key feature that facilitates the interaction with the ACE2 receptors exposed on the host cell. CryoEM and cryoET structural data revealing the architecture of the SARS-CoV-2 viral particle showed that the spike can tilt up to 60° with respect to the perpendicular to the membrane [31, 75]. Behind

Casalino et al.

this flexibility is the structural organization of the extra-virion portion of the spike, composed of two major domains, the stalk and the head, that are connected through a flexible junction that has been referred to as "hip" (Fig. 5A) [10, 63]. Moreover, the stalk can be further divided into an upper and a lower leg, which correspond to the extra-virion alpha-helices of the coil-coiled trimeric bundle, and the transmembrane domain, which can be intended as the foot of this organizational scaffold. The stalk's upper leg, lower leg and the foot are interspersed by highly flexible loops defined as "knee" and "ankle" junctions (Fig. 5A) [63].

We then harnessed DeepDriveMD to perform adaptive MD on the Spike-ACE2 8.5 million atoms system. Following this workflow, we extracted three conformations from the first set of Spike-ACE2 MD simulations (replicas 1-3) and subsequently used them as starting points for a new round of MD (replicas 4-6). We then calculated the distribution of the overall spike tilting with respect to the perpendicular to the membrane (Fig. 5E) and of other three angles involving the stalk, namely the "hip" angle between the stalk's upper leg and the head (Fig. 5B), the "knee" angle between the stalk's lower and upper legs (Fig. 5C), and the "ankle" angle between the perpendicular to the membrane and the stalk's lower leg (Fig. 5D).

The AI-driven adaptive MD approach expanded the conformational space explored, especially for the knee and hip angles, showing average values of  $18.5^{\circ} \pm 7.7^{\circ}$  and  $13.8^{\circ} \pm 7.6^{\circ}$  for replicas 1-3, shifted to  $30.4^{\circ} \pm 5.1^{\circ}$  and  $18.8^{\circ} \pm 6.0^{\circ}$  for the subsequent set of MD (replicas 4-6), respectively. The population shift is less pronounced for the ankle, exhibiting an average angle of  $21.8^{\circ} \pm 2.7^{\circ}$ . These results, in agreement with the data from Turonova et al. [63] that however did not consider the spike in complex with ACE2, reveal large hinge motions throughout the stalk and between the stalk and the head that accommodate the interaction between the spike's RBD and the ACE2 receptor, preventing the disruption of the binding interface. This is further highlighted by the overall tilting of the spike that remains well defined around  $7.3^{\circ} \pm 2.0^{\circ}$  (Fig. 5E), showing that the stalk's inner hinge motions prevent a larger scale bending that could potentially disrupt the RBD-ACE2 interaction.

(6) Our approach points to the very near term ability to accelerate the sampling of dynamical configurations of the complicated viral infection machinery within in the context of its full biological complexity using AI. The enormous amount of data arising from MD and WE simulations of the single spike served to build and train an AI model using the variational autoencoder deep learning approach, which we demonstrate to accelerate dynamical sampling of the spike in a larger, more complex system (i.e., the two parallel membrane spike-ACE2 complex). Thus, the combination of the AI-driven workflows together with the groundbreaking simulations opens the possibility to overcome a current major bottleneck in the development and use of such ultra-large scale MD simulations, which relates to the efficient and effective sampling of the conformational dynamics of a system with so many degrees of freedom. The scientific implications of such a technological advance, in terms of understanding of the basic science of molecular mechanisms of infection as well as the development of novel therapeutics, are vast.

Al-Driven Multisc; 6–19, 2020, Virtual

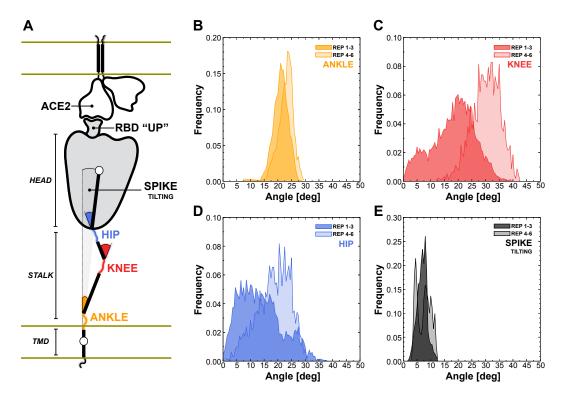


Figure 5: Flexibility of the spike bound to the ACE2 receptor. A) Schematic representation of the two-parallel-membrane system of the spike-ACE2 complex. (B-E) Distributions of the ankle, knee, hip and spike-tilting angles resulting from MD replicas 1-3 (darker color) and 4-6 (lighter color). Starting points for replicas 4-6 have been selected using DeepDriveMD.

- (7) We establish a new high watermark for the atomic-level simulation of viruses with the simulation of the SARS-CoV-2 viral envelope, tallying 305 million atoms including explicit water molecules, and exhibiting a strong scaling on Summit . The virion has a realistic ERGIC-like membrane, contains 24 fully glycosylated full-length spikes (in both the open and closed states) and replicates the spatial patterning and density of viral proteins as determined from cryoelectron tomography experiments [31]. These groundbreaking simulations, just now in the process of being fully analyzed, set the stage for future work on SARS-CoV-2 that will be unprecedented in terms of their ability to more closely mimic realistic biological conditions. This includes, for example, the ability to explore the interactions of the virus with multiple receptors on the host cell, or multiple antibodies. It will allow researchers to explore the correlated dynamics of the molecular pieceparts on the surface of the virus and the host cell, and the effects of curvature on such behavior. It will be used as the ground-truth in the development of other simulation approaches, including coarse grained simulation methods, which are under development [76]. It will aid in the development of methods related to the construction of complicated biological membranes [17]. And the list goes on.
- (8) We developed an AI-driven workflow as a generalizable framework for multiscale simulation. Though we focus here on advances made relevant to COVID19, the methods and workflow established here will be broadly applicable to the multiscale simulation of molecular systems.

#### **ACKNOWLEDGMENTS**

The authors thank D. Maxwell, B. Messer, J. Vermaas, and the Oak Ridge Leadership Computing Facility at Oak Ridge National Laboratory supported by the DOE under Contract DE-AC05-00OR22725. We also thank the Texas Advanced Computing Center Frontera team, especially D. Stanzione and T. Cockerill, and for compute time made available through a Director's Discretionary Allocation (NSF OAC-1818253). We thank the Argonne Leadership Computing Facility supported by the DOE under DE-AC02-06CH11357. NAMD and VMD are funded by NIH P41-GM104601. The NAMD team thanks Intel and M. Brown for contributing the AVX-512 tile list kernels. Anda Trifan acknowledges support from a DOE CSGF (DE-SC0019323). This work was supported by NIH GM132826, NSF RAPID MCB-2032054, an award from the RCSA Research Corp., a UC San Diego Moore's Cancer Center 2020 SARS-COV-2 seed grant, to R.E.A. This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the US DOE Office of Science and the National Nuclear Security Administration. Research was supported by the DOE through the National Virtual Biotechnology Laboratory, a consortium of DOE national laboratories focused on response to COVID-19, with funding from the Coronavirus CARES Act. This work used resources, services, and support from the COVID-19 HPC Consortium (https://covid19hpc-consortium.org/), a private-public effort uniting government, industry, and academic leaders who are volunteering free compute time and resources in support of COVID-19 research. We dedicate this contribution to the memory of Klaus Schulten.

Casalino et al.

#### **REFERENCES**

- Bilge Acun, David J. Hardy, Laxmikant Kale, Ke Li, James C. Phillips, and John E. Stone. 2019. Scalable Molecular Dynamics with NAMD on the Summit System. IBM J. Res. Dev. 62 (2019), 4:1–4:9.
- [2] Jane R. Allison. 2020. Computational methods for exploring protein conformations. Biochemical Society Transactions 48, 4 (08 2020), 1707–1724. https://doi.org/ 10.1042/BST20200193 arXiv:https://portlandpress.com/biochemsoctrans/articlepdf/48/4/1707/891950/bst-2020-0193c.pdf
- [3] Rommie E. Amaro and Adrian J. Mulholland. 2020. A Community Letter Regarding Sharing Biomolecular Simulation Data for COVID-19. Journal of chemical information and modeling (2020), 0-6. https://doi.org/10.1021/acs.jcim.0c00319
- [4] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein GAN. arXiv:stat.ML/1701.07875
- [5] E P Barros, L Casalino, Z Gaieb, A C Dommer, Y Wang, L Fallon, L Raguette, K Belfon, C Simmerling, and R E Amaro. 2020. The flexibility of ACE2 in the context of SARS-CoV-2 infection. bioRxiv (2020). https://doi.org/10.1101/2020. 09.16.300459
- [6] Mattia Bernetti, Martina Bertazzo, and Matteo Masetti. 2020. Data-Driven Molecular Dynamics: A Multifaceted Challenge. *Pharmaceuticals* 13, 9 (2020). https://doi.org/10.3390/ph13090253
- [7] Debsindhu Bhowmik, Shang Gao, Michael T. Young, and Arvind Ramanathan. 2018. Deep clustering of protein folding simulations. BMC Bioinformatics 19, 18 (2018), 484. https://doi.org/10.1186/s12859-018-2507-5
- [8] Luigi Bonati, Yue-Yu Zhang, and Michele Parrinello. 2019. Neural networks-based variationally enhanced sampling. Proceedings of the National Academy of Sciences 116, 36 (2019), 17641–17647. https://doi.org/10.1073/pnas.1907975116 arXiv:https://www.pnas.org/content/116/36/17641.full.pdf
- [9] Lorenzo Casalino, Zied Gaieb, Abigail C Dommer, Aoife M Harbison, Carl A Fogarty, Emilia P Barros, Bryn C Taylor, Elisa Fadda, and Rommie E Amaro. 2020. Shielding and Beyond: The Roles of Glycans in SARS-CoV-2 Spike Protein. bioRxiv (jan 2020), 2020.06.11.146522. https://doi.org/10.1101/2020.06.11.146522
- [10] Lorenzo Casalino, Zied Gaieb, Jory A. Goldsmith, Christy K. Hjorth, Abigail C. Dommer, Aoife M. Harbison, Carl A. Fogarty, Emilia P. Barros, Bryn C. Taylor, Jason S. McLellan, Elisa Fadda, and Rommie E. Amaro. 2020. Beyond Shielding: The Roles of Glycans in the SARS-CoV-2 Spike Protein. ACS Central Science (sep 2020). https://doi.org/10.1021/acscentsci.0c01056
- [11] Doralicia Casares, Pablo V. Escribá, and Catalina Ana Rosselló. 2019. Membrane lipid composition: Effect on membrane and organelle structure, function and compartmentalization and therapeutic avenues. *International Journal of Molecular Sciences* 20, 9 (may 2019). https://doi.org/10.3390/ijms20092167
- [12] Min Chen, Syma Khalid, Mark S P Sansom, and Hagan Bayley. 2008. Outer membrane protein G: Engineering a quiet pore for biosensing. Proceedings of the National Academy of Sciences 105, 17 (2008), 6272–6277. https://doi.org/10.1073/ pnas.0711561105
- [13] Michael F. Crowley, Mark J. Williamson, and Ross C. Walker. 2009. CHAMBER: Comprehensive support for CHARMM force fields within the AMBER software. International Journal of Quantum Chemistry 109, 15 (2009), 3767–3772. https://doi.org/10.1002/qua.22372
- [14] Jacob D. Durrant and Rommie E. Amaro. 2014. LipidWrapper: An Algorithm for Generating Large-Scale Membrane Models of Arbitrary Geometry. PLoS Computational Biology 10, 7 (2014). https://doi.org/10.1371/journal.pcbi.1003720
- [15] Jacob D Durrant, Sarah E Kochanek, Lorenzo Casalino, Pek U Ieong, Abigail C Dommer, and Rommie E Amaro. 2020. Mesoscale all-atom influenza virus simulations suggest new substrate binding mechanism. ACS Cent. Sci. 6, 2 (2020), 189–196.
- [16] G. Fiorin, M. L. Klein, and J. Hénin. 2013. Using collective variables to drive molecular dynamics simulations. Mol. Phys. 111, 22-23 (2013), 3345–3362.
- [17] Fabio Gonzalez-Arias, Tyler Reddy, John Stone, Jodi Hadden-Perilla, and Juan Perilla. 2020. Scalable Analysis of Authentic Viral Envelopes on FRONTERA. Computing in Science and Engineering (2020). https://doi.org/10.1109/MCSE.2020. 3020508
- [18] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. 2017. Improved Training of Wasserstein GANs. arXiv:cs.LG/1704.00028
- [19] Olgun Guvench, Elizabeth Hatcher, Richard M. Venable, Richard W. Pastor, and Alexander D. MacKerell. 2009. CHARMM additive all-atom force field for glycosidic linkages between hexopyranoses. *Journal of Chemical Theory and Compu*tation 5, 9 (sep 2009), 2353–2370. https://doi.org/10.1021/ct900242e
- [20] Andreas W. Götz, Mark J. Williamson, Dong Xu, Duncan Poole, Scott Le Grand, and Ross C. Walker. 2012. Routine Microsecond Molecular Dynamics Simulations with AMBER on GPUs. 1. Generalized Born. *Journal of Chemical Theory and Computation* 8, 5 (May 2012), 1542–1555. https://doi.org/10.1021/ct200909j
- [21] Jing Huang and Alexander D. Mackerell. 2013. CHARMM36 all-atom additive protein force field: Validation based on comparison to NMR data. *Journal of Computational Chemistry* (2013). https://doi.org/10.1002/jcc.23354
- [22] G A Huber and S Kim. 1996. Weighted-ensemble Brownian dynamics simulations for protein association reactions. *Biophysical Journal* 70, 1 (Jan. 1996), 97–110. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1224912/

- [23] William Humphrey, Andrew Dalke, and Klaus Schulten. 1996. VMD Visual Molecular Dynamics. J. Mol. Graphics 14, 1 (1996), 33–38. https://doi.org/10. 1016/0263-7855(96)00018-5
- [24] Sunhwan Jo, Taehoon Kim, Vidyashankara G. Iyer, and Wonpil Im. 2008. CHARMM-GUI: A web-based graphical user interface for CHARMM. Journal of Computational Chemistry 29, 11 (aug 2008), 1859–1865. https://doi.org/10.1002/jcc.20945
- [25] Sunhwan Jo, Kevin C. Song, Heather Desaire, Alexander D. MacKerell, and Wonpil Im. 2011. Glycan reader: Automated sugar identification and simulation preparation for carbohydrates and glycoproteins. *Journal of Computational Chemistry* (2011). https://doi.org/10.1002/jcc.21886
- [26] William L. Jorgensen, Jayaraman Chandrasekhar, Jeffry D. Madura, Roger W. Impey, and Michael L. Klein. 1983. Comparison of simple potential functions for simulating liquid water. *The Journal of Chemical Physics* 79, 2 (jul 1983), 926–935. https://doi.org/10.1063/1.445869
- [27] William L. Jorgensen, Jayaraman Chandrasekhar, Jeffry D. Madura, Roger W. Impey, and Michael L. Klein. 1983. Comparison of simple potential functions for simulating liquid water. *The Journal of Chemical Physics* 79, 2 (July 1983), 926–935. https://doi.org/10.1063/1.445869
- Laxmikant Kalé, Bilge Acun, Seonmyeong Bak, Aaron Becker, Milind Bhandarkar, Nitin Bhat, Abhinav Bhatele, Eric Bohm, Cyril Bordage, Robert Brunner, Ronak Buch, Sayantan Chakravorty, Kavitha Chandrasekar, Jaemin Choi, Michael Denardo, Jayant DeSouza, Matthias Diener, Harshit Dokania, Isaac Dooley, Wayne Fenton, Juan Galvez, Fillipo Gioachin, Abhishek Gupta, Gagan Gupta, Manish Gupta, Attila Gursoy, Vipul Harsh, Fang Hu, Chao Huang, Narain Jagathesan, Nikhil Jain, Pritish Jetley, Prateek Jindal, Raghavendra Kanakagiri, Greg Koenig, Sanjeev Krishnan, Sameer Kumar, David Kunzman, Michael Lang, Akhil Langer, Orion Lawlor, Chee Wai Lee, Jonathan Lifflander, Karthik Mahesh, Celso Mendes, Harshitha Menon, Chao Mei, Esteban Meneses, Eric Mikida, Phil Miller, Ryan Mokos, Venkatasubrahmanian Narayanan, Xiang Ni, Kevin Nomura, Sameer Paranjpye, Parthasarathy Ramachandran, Balkrishna Ramkumar, Evan Ramos, Michael Robson, Neelam Saboo, Vikram Saletore, Osman Sarood, Karthik Senthil, Nimish Shah, Wennie Shu, Amitabh B. Sinha, Yanhua Sun, Zehra Sura, Ehsan Totoni, Krishnan Varadarajan, Ramprasad Venkataraman, Jackie Wang, Lukasz Wesolowski, Sam White, Terry Wilmarth, Jeff Wright, Joshua Yelon, and Gengbin Zheng. 2019. The Charm++ Parallel Programming System. https://doi.org/10.5281/zenodo.3370873
- [29] Laxmikant V. Kalé and Gengbin Zheng. 2013. Chapter 1: The Charm++ Programming Model. In Parallel Science and Engineering Applications: The Charm++ Approach (1st ed.), Laxmikant V. Kale and Abhinav Bhatele (Eds.). CRC Press, Inc., Boca Raton, FL, USA, Chapter 1, 1–16. https://doi.org/10.1201/b16251
- 30] Peter M Kasson and Shantenu Jha. 2018. Adaptive ensemble simulations of biomolecules. Current Opinion in Structural Biology 52 (2018), 87 – 94. https: //doi.org/10.1016/j.sbi.2018.09.005 Cryo electron microscopy: the impact of the cryo-EM revolution in biology • Biophysical and computational methods - Part A
- [31] Zunlong Ke, Joaquin Oton, Kun Qu, Mirko Cortese, Vojtech Zila, Lesley McKeane, Takanori Nakane, Jasenko Zivanov, Christopher J. Neufeldt, Berati Cerikan, John M. Lu, Julia Peukes, Xiaoli Xiong, Hans Georg Kräusslich, Sjors H.W. Scheres, Ralf Bartenschlager, and John A.G. Briggs. 2020. Structures and distributions of SARS-CoV-2 spike proteins on intact virions. *Nature* (aug 2020), 1–7. https: //doi.org/10.1038/s41586-020-2665-2
- [32] S. Kumar, A. R. Mamidala, D. A. Faraj, B. Smith, M. Blocksome, B. Cernohous, D. Miller, J. Parker, J. Ratterman, P. Heidelberger, D. Chen, and B. Steinmacher-Burrow. 2012. PAMI: A Parallel Active Message Interface for the Blue Gene/Supercomputer. In 2012 IEEE 26th International Parallel and Distributed Processing Symposium. 763–773. https://doi.org/10.1109/IPDPS.2012.73
- [33] João Marcelo Lamim Ribeiro and Pratyush Tiwary. 2019. Toward Achieving Efficient and Accurate Ligand-Protein Unbinding with Deep Learning and Molecular Dynamics through RAVE. Journal of Chemical Theory and Computation 15, 1 (01 2019), 708-719. https://doi.org/10.1021/acs.jctc.8b00869
- [34] E. H. Lee, J. Hsin, M. Sotomayor, G. Comellas, and K. Schulten. 2009. Discovery through the computational microscope. Structure 17 (2009), 1295–1306.
- [35] H. Lee, M. Turilli, S. Jha, D. Bhowmik, H. Ma, and A. Ramanathan. 2019. Deep-DriveMD: Deep-Learning Driven Adaptive Molecular Simulations for Protein Folding. In 2019 IEEE/ACM Third Workshop on Deep Learning on Supercomputers (DLS). 12–19.
- [36] Frank Noé. 2020. Machine Learning for Molecular Dynamics on Long Timescales. Springer International Publishing, Cham, 331–372. https://doi.org/10.1007/978-3-030-40245-7\_16
- [37] José Nelson Onuchic, Zaida Luthey-Schulten, and Peter G. Wolynes. 1997. THE-ORY OF PROTEIN FOLDING: The Energy Landscape Perspective. Annual Review of Physical Chemistry 48, 1 (1997), 545–600. https://doi.org/10.1146/annurev.physchem.48.1.545 arXiv:https://doi.org/10.1146/annurev.physchem.48.1.545 PMID: 9348663.
- [38] Sang Jun Park, Jumin Lee, Yifei Qi, Nathan R. Kern, Hui Sun Lee, Sunhwan Jo, Insuk Joung, Keehyung Joo, Jooyoung Lee, and Wonpil Im. 2019. CHARMM-GUI

- Glycan Modeler for modeling and simulation of carbohydrates and glycoconjugates. *Glycobiology* (2019). https://doi.org/10.1093/glycob/cwz003
- [39] Adrià Pérez, Pablo Herrera-Nieto, Stefan Doerr, and Gianni De Fabritiis. 2020. AdaptiveBandit: A Multi-armed Bandit Framework for Adaptive Sampling in Molecular Simulations. Journal of Chemical Theory and Computation 16, 7 (07 2020), 4685–4693. https://doi.org/10.1021/acs.jetc.0c00205
- [40] James Phillips, Gengbin Zheng, Sameer Kumar, and Laxmikant Kale. 2002. NAMD: Biomolecular Simulation on Thousands of Processors. In Proceedings of the IEEE/ACM SC2002 Conference, Technical Paper 277. IEEE Press, Baltimore, Maryland, 1–18.
- [41] James C. Phillips, Rosemary Braun, Wei Wang, James Gumbart, Emad Tajkhorshid, Elizabeth Villa, Christophe Chipot, Robert D. Skeel, Laxmikant Kale, and Klaus Schulten. 2005. Scalable molecular dynamics with NAMD. J. Comput. Chem. 26 (2005), 1781–1802.
- [42] James C. Phillips, David J. Hardy, Julio D. C. Maia, John E. Stone, João V. Ribeiro, Rafael C. Bernardi, Ronak Buch, Giacomo Fiorin, Jérôme Hénin, Wei Jiang, Ryan McGreevy, Marcelo C. R. Melo, Brian Radak, Robert D. Skeel, Abhishek Singharoy, Yi Wang, Benoît Roux, Aleksei Aksimentiev, Zaida Luthey-Schulten, Laxmikant V. Kalé, Klaus Schulten, Christophe Chipot, and Emad Tajkhorshid. 2020. Scalable molecular dynamics on CPU and GPU architectures with NAMD. J. Chem. Phys. 153 (2020), 044130. https://doi.org/10.1063/5.0014475
- [43] James C. Phillips, John E. Stone, and Klaus Schulten. 2008. Adapting a Message-Driven Parallel Application to GPU-Accelerated Clusters. In SC '08: Proceedings of the 2008 ACM/IEEE Conference on Supercomputing. IEEE Press, Piscataway, NJ, USA, 1–9. (9 pages).
- [44] James C. Phillips, Yanhua Sun, Nikhil Jain, Eric J. Bohm, and Laximant V. Kalé. 2014. Mapping to Irregular Torus Topologies and Other Techniques for Petascale Biomolecular Simulation. In Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC '14). IEEE Press, 81–91. https://doi.org/10.1109/SC.2014.12
- [45] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. 2017. Point-Net: Deep Learning on Point Sets for 3D Classification and Segmentation. arXiv:cs.CV/1612.00593
- [46] A. Ramanathan, Andrej J. Savol, Virginia M. Burger, S. Quinn, P. Agarwal, and C. Chennubhotla. 2012. Statistical Inference for Big Data Problems in Molecular Biophysics.
- [47] João Marcelo Lamim Ribeiro, Pablo Bravo, Yihang Wang, and Pratyush Tiwary. 2018. Reweighted autoencoded variational Bayes for enhanced sampling (RAVE). The Journal of Chemical Physics 149, 7 (2018), 072301. https://doi.org/10.1063/1. 5025487 arXiv:https://doi.org/10.1063/1.5025487
- [48] Raquel Romero, Arvind Ramanathan, Tony Yuen, Debsindhu Bhowmik, Mehr Mathew, Lubna Bashir Munshi, Seher Javaid, Madison Bloch, Daria Lizneva, Alina Rahimova, Ayesha Khan, Charit Taneja, Se-Min Kim, Li Sun, Maria I. New, Shozeb Haider, and Mone Zaidi. 2019. Mechanism of glucocerebrosidase activation and dysfunction in Gaucher disease unraveled by molecular dynamics and deep learning. Proceedings of the National Academy of Sciences 116, 11 (2019), 5086–5095. https://doi.org/10.1073/pnas.1818411116 arXiv:https://www.pnas.org/content/116/11/5086.full.pdf
- [49] Jean-Paul Ryckaert, Giovanni Ciccotti, and Herman J. C Berendsen. 1977. Numerical integration of the cartesian equations of motion of a system with constraints: molecular dynamics of n-alkanes. J. Comput. Phys. 23, 3 (March 1977), 327–341. https://doi.org/10.1016/0021-9991(77)9008-5
- [50] Ali S Saglam and Lillian T Chong. 2019. Protein-protein binding pathways and calculations of rate constants using fully-continuous, explicit-solvent simulations. *Chemical science* 10, 8 (2019), 2360–2372.
- [51] Andrej Šali and Tom L. Blundell. 1993. Comparative protein modelling by satisfaction of spatial restraints. *Journal of Molecular Biology* 234, 3 (dec 1993), 779–815. https://doi.org/10.1006/jmbi.1993.1626
- [52] Romelia Salomon-Ferrer, Andreas W. Götz, Duncan Poole, Scott Le Grand, and Ross C. Walker. 2013. Routine Microsecond Molecular Dynamics Simulations with AMBER on GPUs. 2. Explicit Solvent Particle Mesh Ewald. *Journal of Chemical Theory and Computation* 9, 9 (Sept. 2013), 3878–3888. https://doi.org/ 10.1021/ct400314v
- [53] Asif Shajahan, Stephanie Archer-Hartmann, Nitin T. Supekar, Anne S. Gleinich, Christian Heiss, Parastoo Azadi, and Ka Sheraton. 2020. Comprehensive characterization of N- and O- glycosylation of SARS-CoV-2 human receptor angiotensin converting enzyme 2. bioRxiv (aug 2020), 2020.05.01.071688. https: //doi.org/10.1101/2020.05.01.071688
- [54] Asif Shajahan, Nitin T Supekar, Anne S Gleinich, and Parastoo Azadi. [n.d.]. Deducing the N-and O-glycosylation profile of the spike protein of novel coronavirus SARS-CoV-2. Glycobiology 2020 ([n.d.]), 1–8. https://doi.org/10.1093/glycob/cwaa042
- [55] D.E. Shaw, M.M. Deneroff, R.O. Dror, J.S. Kuskin, R.H. Larson, J.K. Salmon, C. Young, B. Batson, K.J. Bowers, J.C. Chao, M.P. Eastwood, J. Gagliardo, J.P. Grossman, C.R. Ho, D.J. Ierardi, I. Kolossváry, J.L. Klepeis, T. Layman, C. McLeavey, M.A. Moraes, R. Mueller, E.C. Priest, Y. Shan, J. Spengler, M. Theobald, B. Towles, and S.C. Wang, 2007. Anton, a special-purpose machine for molecular dynamics

- simulation. SIGARCH Comput. Archit. News 35 (2007), 1-12.
- [56] A. Singharoy, C. Maffeo, K. Delgardo, D. J. K. Swainsbury, M. Sener, U. Kleinekathöfer, B. Isralewitz, I. Teo, D. Chandler, J. Stone, J. Phillips, T. Pogorelov, M. I. Mallus, C. Chipot, Z. Luthey-Schulten, P. Tieleman, C. N. Hunter, Emad Tajkhorshid, A. Aksimentiev, and K. Schulten. 2019. Atoms to Phenotypes: Molecular Design Principles of Cellular Energy Metabolism. Cell 179 (2019), 1098–1111.
- [57] Tyler N Starr, Allison J Greaney, Sarah K Hilton, Daniel Ellis, Katharine H D Crawford, Adam S Dingens, Mary Jane Navarro, John E Bowen, M Alejandra Tortorici, Alexandra C Walls, Neil P King, David Veesler, and Jesse D Bloom. 2020. Deep Mutational Scanning of SARS-CoV-2 Receptor Binding Domain Reveals Constraints on Folding and ACE2 Binding. Cell 182, 5 (sep 2020), 1295–1310.e20. https://doi.org/10.1016/j.cell.2020.08.012
- [58] John E. Stone, Antti-Pekka Hynninen, James C. Phillips, and Klaus Schulten. 2016. Early Experiences Porting the NAMD and VMD Molecular Simulation and Analysis Software to GPU-Accelerated OpenPOWER Platforms. *International Workshop on OpenPOWER for HPC (IWOPH'16)* (2016), 188–206.
- [59] John E. Stone, Barry Isralewitz, and Klaus Schulten. 2013. Early Experiences Scaling VMD Molecular Visualization and Analysis Jobs on Blue Waters. In Extreme Scaling Workshop (XSW), 2013. 43–50. https://doi.org/10.1109/XSW.2013. 10
- [60] John E. Stone, Melih Sener, Kirby L. Vandivort, Angela Barragan, Abhishek Singharoy, Ivan Teo, Joao V. Ribeiro, Barry Isralewitz, Bo Liu, Boon Chong Goh, James C. Phillips, Craig MacGregor-Chatwin, Matthew P. Johnson, Lena F. Kourkoutis, C. Neil Hunter, and Klaus Schulten. 2016. Atomic Detail Visualization of Photosynthetic Membranes with GPU-Accelerated Ray Tracing. Parallel Comput. 55 (2016), 17–27. https://doi.org/10.1016/j.parco.2015.10.015
- [61] John E. Stone, Kirby L. Vandivort, and Klaus Schulten. 2013. GPU-Accelerated Molecular Visualization on Petascale Supercomputing Platforms. In Proceedings of the 8th International Workshop on Ultrascale Visualization (UltraVis '13). ACM, New York, NY, USA, Article 6, 8 pages.
- [62] Zeyu Sun, Keyi Ren, Xing Zhang, Jinghua Chen, Zhengyi Jiang, Jing Jiang, Feiyang Ji, Xiaoxi Ouyang, and Lanjuan Li. 2020. Mass Spectrometry Analysis of Newly Emerging Coronavirus HCoV-19 Spike Protein and Human ACE2 Reveals Camouflaging Glycans and Unique Post-Translational Modifications. Engineering (2020). https://doi.org/10.1016/j.eng.2020.07.014
- [63] Beata Turoňová, Mateusz Sikora, Christoph Schürmann, Wim J. H. Hagen, Sonja Welsch, Florian E. C. Blanc, Sören von Bülow, Michael Gecht, Katrin Bagola, Cindy Hörner, Ger van Zandbergen, Jonathan Landry, Nayara Trevisan Doimo de Azevedo, Shyamal Mosalaganti, Andre Schwarz, Roberto Covino, Michael D. Mühlebach, Gerhard Hummer, Jacomine Krijnse Locker, and Martin Beck. 2020. In situ structural analysis of SARS-CoV-2 spike reveals flexibility mediated by three hinges. Science (aug 2020), eabd5223. https://doi.org/10.1126/science.abd5223
- [64] L.J.P. van der Maaten and G.E. Hinton. 2008. Visualizing High-Dimensional Data Using t-SNE. Journal of Machine Learning Research 9, nov (2008), 2579–2605. Pagination: 27.
- [65] Gerrit Van Meer, Dennis R. Voelker, and Gerald W. Feigenson. 2008. Membrane lipids: Where they are and how they behave. , 112–124 pages. https://doi.org/ 10.1038/nrm2330
- [66] Alexandra C. Walls, Young Jun Park, M. Alejandra Tortorici, Abigail Wall, Andrew T. McGuire, and David Veesler. 2020. Structure, Function, and Antigenicity of the SARS-CoV-2 Spike Glycoprotein. Cell 181, 2 (apr 2020), 281–292.e6. https://doi.org/10.1016/j.cell.2020.02.058
- [67] Yihang Wang, João Marcelo Lamim Ribeiro, and Pratyush Tiwary. 2019. Past–future information bottleneck for sampling molecular reaction coordinate si-multaneously with thermodynamics and kinetics. *Nature Communications* 10, 1 (2019), 3573. https://doi.org/10.1038/s41467-019-11405-4
- [68] Yihang Wang, Joao Marcelo Lamim Ribeiro, and Pratyush Tiwary. 2020. Machine learning approaches for analyzing and enhancing molecular dynamics simulations. Current Opinion in Structural Biology 61 (2020), 139–145.
- [69] Yasunori Watanabe, Joel D Allen, Daniel Wrapp, Jason S McLellan, and Max Crispin. 2020. Site-specific glycan analysis of the SARS-CoV-2 spike. Science (New York, N.Y.) (may 2020). https://doi.org/10.1126/science.abb9983
- [70] Daniel Wrapp, Nianshuang Wang, Kizzmekia S Corbett, Jory A Goldsmith, Ching-Lin Hsieh, Olubukola Abiona, Barney S Graham, and Jason S McLellan. 2020. Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. Science (New York, N.Y.) 1263, March (2020), 1260–1263. https://doi.org/10.1126/science. abb2507
- [71] Daniel Wrapp, Nianshuang Wang, Kizzmekia S. Corbett, Jory A. Goldsmith, Ching-Lin Hsieh, Olubukola Abiona, Barney S. Graham, and Jason S. McLellan. 2020. Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* 367, 6483 (March 2020), 1260–1263. https://doi.org/10.1126/science. abb2507
- [72] Emilia L. Wu, Xi Cheng, Sunhwan Jo, Huan Rui, Kevin C. Song, Eder M. Dávila-Contreras, Yifei Qi, Jumin Lee, Viviana Monje-Galvan, Richard M. Venable, Jeffery B. Klauda, and Wonpil Im. 2014. CHARMM-GUI membrane builder toward realistic biological membrane simulations. https://doi.org/10.1002/jcc.23702

Casalino et al.

- [73] Renhong Yan, Yuanyuan Zhang, Yaning Li, Lu Xia, Yingying Guo, and Qiang Zhou. 2020. Structural basis for the recognition of SARS-CoV-2 by full-length human ACE2. Science 367, 6485 (mar 2020), 1444–1448. https://doi.org/10.1126/ science.abb2762
- [74] Charlene Yang. 2020. Hierarchical Roofline Analysis: How to Collect Data using Performance Tools on Intel CPUs and NVIDIA GPUs. arXiv:cs.DC/2009.02449
- [75] Hangping Yao, Yutong Song, Yong Chen, Yigong Shi, Lanjuan Li, Sai Li Correspondence, Nanping Wu, Jialu Xu, Chujie Sun, Jiaxing Zhang, Tianhao Weng, Zheyuan Zhang, Zhigang Wu, Linfang Cheng, Danrong Shi, Xiangyun Lu, Jianlin Lei, Max Crispin, and Sai Li. 2020. Molecular Architecture of the SARS-CoV-2 Virus. Cell 183 (2020). https://doi.org/10.1016/j.cell.2020.09.018
- [76] Alvin Yu, Alexander J Pak, Peng He, Viviana Monje-Galvan, Lorenzo Casalino, Zied Gaieb, Abigail C Dommer, Rommie E Amaro, and Gregory A Voth. 2020. A Multiscale Coarse-grained Model of the SARS-CoV-2 Virion. bioRxiv (oct 2020), 2020.10.02.323915. https://doi.org/10.1101/2020.10.02.323915
- [77] Maciej Zamorski, Maciej Zięba, Piotr Klukowski, Rafał Nowak, Karol Kurach, Wojciech Stokowiec, and Tomasz Trzciński. 2020. Adversarial autoencoders for compact representations of 3D point clouds. Computer Vision and Image Understanding 193 (2020), 102921. https://doi.org/10.1016/j.cviu.2020.102921
- [78] Bin W. Zhang, David Jasnow, and Daniel M. Zuckerman. 2010. The "weighted ensemble" path sampling method is statistically exact for a broad class of stochastic processes and binning procedures. The Journal of Chemical Physics 132, 5 (2010), 054107. https://doi.org/10.1063/1.3306345

- arXiv:https://doi.org/10.1063/1.3306345
- [79] Yang Zhang. 2008. I-TASSER server for protein 3D structure prediction. BMC Bioinformatics 9, 1 (jan 2008), 1–8. https://doi.org/10.1186/1471-2105-9-40
- [80] Gongpu Zhao, Juan R. Perilla, Ernest L. Yufenyuy, Xin Meng, Bo Chen, Jiying Ning, Jinwoo Ahn, Angela M. Gronenborn, Klaus Schulten, Christopher Aiken, and Peijun Zhang. 2013. Mature HIV-1 Capsid structure by Cryo-electron microscopy and All-Atom Molecular Dynamics. *Nature* 497 (2013), 643–646. https://doi.org/ 10.1038/nature12162
- [81] Peng Zhao, Jeremy L. Praissman, Oliver C. Grant, Yongfei Cai, Tianshu Xiao, Katelyn E. Rosenbalm, Kazuhiro Aoki, Benjamin P. Kellman, Robert Bridger, Dan H. Barouch, Melinda A. Brindley, Nathan E. Lewis, Michael Tiemeyer, Bing Chen, Robert J. Woods, and Lance Wells. 2020. Virus-Receptor Interactions of Glycosylated SARS-CoV-2 Spike and Human ACE2 Receptor. Cell Host and Microbe 28, 4 (oct 2020), 586–601.e6. https://doi.org/10.1016/j.chom.2020.08.004
- [82] Daniel M. Zuckerman and Lillian T. Chong. 2017. Weighted Ensemble Simulation: Review of Methodology, Applications, and Software. Annual Review of Biophysics 46 (2017), 43–57. https://doi.org/10.1146/annurev-biophys-070816-033834
- [83] Matthew C. Zwier, Joshua L. Adelman, Joseph W. Kaus, Adam J. Pratt, Kim F. Wong, Nicholas B. Rego, Ernesto Suárez, Steven Lettieri, David W. Wang, Michael Grabe, Daniel M. Zuckerman, and Lillian T. Chong. 2015. WESTPA: An Interoperable, Highly Scalable Software Package for Weighted Ensemble Simulation and Analysis. Journal of Chemical Theory and Computation 11, 2 (Feb. 2015), 800–809. https://doi.org/10.1021/ct5010615 Publisher: American Chemical Society.