Biophysical Journal Article



Data-driven prediction of $\alpha_{IIb}\beta_3$ integrin activation paths using manifold learning and deep generative modeling

Siva Dasetty,¹ Tamara C. Bidone,^{2,3} and Andrew L. Ferguson^{1,*}

¹Pritzker School of Molecular Engineering, University of Chicago, Chicago, Illinois; ²Department of Biomedical Engineering, University of Utah, Salt Lake City, Utah; and ³Scientific Computing and Imaging Institute, University of Utah, Salt Lake City, Utah

ABSTRACT The integrin heterodimer is a transmembrane protein critical for driving cellular process and is a therapeutic target in the treatment of multiple diseases linked to its malfunction. Activation of integrin involves conformational transitions between bent and extended states. Some of the conformations that are intermediate between bent and extended states of the heterodimer have been experimentally characterized, but the full activation pathways remain unresolved both experimentally due to their transient nature and computationally due to the challenges in simulating rare barrier crossing events in these large molecular systems. An understanding of the activation pathways can provide new fundamental understanding of the biophysical processes associated with the dynamic interconversions between bent and extended states and can unveil new putative therapeutic targets. In this work, we apply nonlinear manifold learning to coarse-grained molecular dynamics simulations of bent, extended, and two intermediate states of $\alpha_{IIb}\beta_3$ integrin to learn a low-dimensional embedding of the configurational phase space. We then train deep generative models to learn an inverse mapping between the low-dimensional embedding and highdimensional molecular space and use these models to interpolate the molecular configurations constituting the activation pathways between the experimentally characterized states. This work furnishes plausible predictions of integrin activation pathways and reports a generic and transferable multiscale technique to predict transition pathways for biomolecular systems.

SIGNIFICANCE Access to transition structures of the integrin heterodimer remains challenging experimentally due to their transient nature and computationally since the transitions are rare events. We address this challenge by developing a machine learning method trained on molecular simulation data within known stable states to predict activation pathways between these states. We apply this approach to the $\alpha_{IIb}\beta_3$ integrin, a transmembrane cell signaling protein whose malfunction can lead to multiple pathologies. Our results provide predictions of putative integrin activation pathways and mechanisms and can guide the identification of new metastable states as candidate therapeutic targets. Our method is generic and transferable to other large biomolecular systems.

INTRODUCTION

Integrins are transmembrane proteins that mediate the signals across the cell membrane between extracellular matrix and cytoplasm (1-4) and play a vital role in various biological processes such as cell adhesion, migration, proliferation, and differentiation, tissue development, angiogenesis, mechanosensing, and homeostasis (5-11). The two molecules constituting the integrin heterodimer are referred as

*Correspondence: andrewferguson@uchicago.edu Editor: Tamar Schlick. https://doi.org/10.1016/j.bpj.2023.12.009 © 2023 Biophysical Society. the α and β subunits that exist in a noncovalently associated complex (2). The activation mechanism of integrin involves a large-scale allosteric conformational change of the noncovalently interacting α and β subunits from an inactive (lowaffinity) to an active (high-affinity) state (2). The structures of integrin in the inactive and active ligand-binding state are respectively known as the "bent-closed" and "extendedopen" conformations (8). Fig. 1 illustrates the all-atom (AA) cartoon structures of the bent-closed and extendedopen conformations of $\alpha_{IIb}\beta_3$ integrin at near physiological conditions initially resolved using cryo-electron microscopy (12) and subsequently used to initialize AA molecular dynamics simulations (13). In the bent-closed conformation (Fig. 1 *a*), both α and β subunits are generally closely

Submitted October 14, 2023, and accepted for publication December 11, 2023.

Dasetty et al.



FIGURE 1 All-atom structures of $\alpha_{IIb}\beta_3$ integrin in the inactive bent-closed and active extended-open states along with two extended-intermediate states. The atomic coordinates of these structures were taken from a prior study by Tong et al. (13) employing all-atom molecular dynamics simulations with explicit water in the presence of lipid membrane. The initial coordinates of the four structures for the all-atom molecular dynamics simulations in Tong et al. were themselves constructed from cryo-electron microscopy data (12). (*a*) In the bent-closed state, both the α and β subunits are closely associated with a large part of the head ectodomains bent at the genu or linker regions, and the remaining cytoplasmic helical regions are extended and embedded within the cell membrane. (*b*) Structure of extended-intermediate 1 (Int 1) state exhibiting an opening at the genu or linker regions relative to the bent-closed structure. (*c*) Structure of extended-intermediate 2 (Int 2) state wherein the hybrid domain in β subunit swings out relative to the Int 1 state along with a crossing of the tail helices. (*d*) Structure of extended-open state resulting from separation of transmembrane tail helices in Int 2 state. Molecular renderings in this and all subsequent figures were constructed using VMD 1.9.4 (14). To see this figure in color, go online.

associated with a large part of their head ectodomains bent at the genu or linker regions connecting the remaining extended but associated leg regions (2). The extracellular ligand-binding site in $\alpha_{IIb}\beta_3$ lies between the headpiece $\beta_{propeller}$ domain of α subunit and β_3 -I domain of the β subunit of integrin, which in the bent-closed conformation is closer to the membrane but remains accessible to the ligand (12). The upright extended-open conformation (Fig. 1 *d*) differs from the bent-closed conformation by an extension of the ectodomain of both α and β subunits away from the membrane and a separation of leg regions of α and β subunits at the genu or linker regions by a distance of \sim 7–8 nm (2,12).

Conversion from the inactive to active state of integrin is a typically reversible process that occurs within seconds and can be triggered bidirectionally by proteins or ligands in the cytoplasm ("inside-out" activation) or in the extracellular matrix ("outside-in" activation) (15,16). In addition, integrin activation can be triggered by activating mutations and model agonists such as divalent cations, phorbol myristate acetate, calcium ionophore, and monoclonal antibodies (17–21). Experimental studies of integrin activation from the bent-closed to extended-open conformations have led to two prevailing hypothesized activation mechanisms termed the "switchblade" and "deadbolt" pathways (16,22–24). In the switchblade model, inside-out activation results in separation

of extended legs of integrin followed by an opening of the head ectodomains of both α and β subunits bent at the genu or linker regions into the straightened extended-open conformation (high-affinity ligand-binding state) in a motion similar to the opening of a pocketknife (23, 25, 26). In the deadbolt model, there is a progressive loss of interactions between the β -tail (β -T) domain (deadbolt region) and β -I domain (lock region) in the β subunit with a piston, seesaw, sliding, or rotation movement before the separation of the legs to expose the ligand-binding site (26). In contrast to switchblade model, the deadbolt model considers the bentclosed conformation to be a ligand-binding-active state, and the extended-open state may be achieved after a change in affinity or upon ligand binding (23). Alternative models of activation have also been proposed, including a cooperative activation mechanism that requires clustering of multiple integrins facilitated by cytoplasmic talin proteins (27,28) and a light switch mechanism wherein the extension of the bentclosed conformation to extended-open conformation primarily requires a change in the tilt angle of the transmembrane helix of the β -subunit potentially facilitated by the weakened interactions between the transmembrane helices of α and β subunits when embedded within the membrane (12).

Biochemical, crystallographic, and microscopy studies have shed light on the major conformational states of the integrin heterodimer, but the transient intermediate structures

and atomistic details of the activation mechanism remain much less well understood. Integrin is a major therapeutic target because its malfunction can result in various diseases such as bleeding disorders, immunodeficiency, and cancer (29). The $\alpha_{IIb}\beta_3$ variant, for example, is a major platelet integrin that plays an important role in hemostasis, thrombosis, and atherosclerosis (12,18,30). Malfunction of $\alpha_{IIb}\beta_3$ integrin activation is linked to thrombotic disorders (heart attack and stroke) and bleeding disorders (13,31,32). Molecular understanding of the integrin activation mechanism would not only provide insights into the fundamental cellular processes driven by integrin but can also identify new structural targets for integrin-targeting therapeutics (9,33-35). However, mapping transient structures and understanding molecular mechanism of integrin activation is challenging for both experiment and computation. With experimental methods, mapping high-resolution structures along the activation path is difficult because of their inherently transient nature. Currently only a few intermediate structures at high resolution are available such as those shown in Fig. 1 b and c (12). Computer simulations can shed light on the atomistic details of activation but are hindered by the large size of the molecular system: the integrin heterodimer proteins comprise ~ 1800 residues (13), making unbiased sampling of rare activation events computationally burdensome and biased simulations challenging due to the absence of good order parameters with which to drive sampling of activation.

In this work, we combine molecular dynamics (MD) simulations, nonlinear manifold learning, and deep generative modeling within a multiscale framework to predict activation pathways from simulation data collected within known intermediate metastable or stable states. We apply this approach to the activation pathways of a prototypical $\alpha_{IIb}\beta_3$ integrin heterodimer, but the technique is generically extensible to other molecular systems for which the intermediate states are known but for which the activation pathways are computationally prohibitive rare events. The remainder of this paper is structured as follows. We first provide a description of the generation of the coarse-grained MD training data within the four known intermediate states and details of each component of our computational approach to predict activation pathways: density-adaptive diffusion maps (dMaps), conditional Wasserstein generative adversarial networks (cWGAN), and targeted molecular dynamics (TMD). We then present and discuss the latent space of the model integrin $\alpha_{IIb}\beta_3$, cWGAN generated molecular structures for each intermediate state, and cWGAN predicted activation pathways between these states. We conclude with a discussion of the configurational plausibility of the predicted pathways, the possibility of using these structures to seed a targeted campaign of unbiased or enhanced sampling AA simulations, and the potential for experimental testing and validation of these predictions.

Prediction of integrin activation paths

MATERIALS AND METHODS

We present in Fig. 2 an overview of the four steps of our computational protocol: 1) calculation of the similarity matrix between coarse-grained (CG) representations of the $\alpha_{IIb}\beta_3$ integrin heterodimer harvested from AA MD simulations of the bent-closed, Int 1, Int 2, and extended-open states (Fig. 2 a), 2) application of density-adaptive dMaps to learn a nonlinear projection of the simulation data into a low-dimensional latent space (Fig. 2 b), 3) training of a cWGAN to approximate the inverse mapping from the lowdimensional latent space to the high-dimensional molecular space and the use of the trained model to predict the molecular activation pathways between metastable states by interpolation (Fig. 2 c), and 4) TMD backmapping calculations to restore AA resolution to the predicted CG configurations (Fig. 2 d). In a nutshell, the dMaps trained over simulation data in the four states furnish a low-dimensional compression of the high-dimensional simulation data into a latent space exposing the structural similarities of configurations within and between the intermediate states. The low-dimensional nature of this embedding preserves the gross structural features discriminating the states and enables the construction of interpolative pathways between the states through regions of configurational space in which no simulation data is available. The trained cWGAN learns the inverse mapping from the latent space back to the high-dimensional molecular space and serves as a deep generative model with which to predict (i.e., "hallucinate") putative molecular configurations along the activation pathways between states. For reasons of computational tractability and stability in cWGAN training and deployment, we operate these components of the pipeline using a 300-bead CG representation of the integrin heterodimer, but we perform post hoc restoration of AA detail using TMD calculations. In the following sections, we provide full details of each of the four steps in the pipeline.

All-atom training data in each intermediate state

Direct simulations of integrin activation pathways are challenging because of the high cost of simulating these large systems and the nature of the transition as a rare event. Although we are unaware of any direct unbiased simulations of activation, a limited number of studies have utilized enhanced sampling calculations to drive activation using artificial biasing forces (36). For instance, Kulke and Langel applied steered MD simulations to artificially mimic the outside-in and inside-out activation pathways of $\alpha_V \beta_3$ integrin and observed switchblade like extension mechanism from a bent-closed to an extended-open-like state in both scenarios (37). Wang and Li captured extension of $\alpha_{IIb}\beta_3$ integrin from a bent-closed state to an extended state without transmembrane tail helices separation in one of six steered MD simulations with pulling forces applied only to fibronectin ligand attached to the headpiece of integrin (38). Tong et al. (13) previously conducted unbiased MD simulations of $\alpha_{IIb}\beta_3$ integrin within the bent-closed, Int 1, Int 2, and extended-open states commencing from initial structures obtained from cryo-electron microscopy (12). Although no transitions were observed between the intermediate states in this study, we hypothesized that our computational approach could be trained over these unbiased data and used to propose biophysically plausible activation pathways without requiring the application of artificial biasing potentials. Each unbiased AA MD simulation in Tong et al. comprised the $\alpha_{IIb}\beta_3$ integrin, lipid bilayer (3:1 M ratio of 1,2-dioleoyl-sn-glycero-3-phosphocholine (DOPC):1,2-dioleoyl-sn-glycero-3-phospho-L-serine (DOPS)), and 150 mM of NaCl in water at 310 K and 1 atm (13). Snapshots were harvested at a period of 200 ps over the \sim 500 ns to collect a total of 2450 configurations for each state. Full details of the simulation protocol including system preparation using CHARMM-GUI (39) and addition of missing residues within the four experimental structures are reported in Tong et al. (13).

Although we could have conducted the dMaps nonlinear dimensionality reduction and cWGAN reconstruction in the AA representation, to reduce the computational cost and improve stability of cWGAN training and deployment, we elected to conduct these operations within a reduced-dimensional CG

Dasetty et al.

a Similarity matrix (CG resolution) b Density-adaptive diffusion maps



FIGURE 2 Schematic overview of the multiscale computational approach to predict the molecular activation pathways between intermediate states of the $\alpha_{IIb}\beta_3$ integrin heterodimer. (a) Training data is collected from ~500 ns all-atom MD simulations of the $\alpha_{IIb}\beta_3$ integrin heterodimer within a lipid bilayer in the bent-closed, Int 1, Int 2, and extended-open states. To reduce training costs and improve stability of the cWGAN, we develop 300bead coarse-grained representations of the system from the ~1750-residue AA representations and compute pairwise root mean-square deviations (RMSDs) between the translationally and rotationally aligned snapshots harvested from each trajectory. (b) The RMSD pairwise distance matrix serves as an input to the density-adaptive diffusion maps (dMaps) that learn an embedding of the simulation trajectories into a low-dimensional latent space that clusters configurations according to structural similarity within and between the four intermediate states. A schematic illustration of the latent space is presented here. (c) We train a conditional Wasserstein generative adversarial network (cWGAN) to learn the inverse mapping from the low-dimensional latent space to the high-dimensional 300-bead coarse-grained molecular space. As a deep generative model, the cWGAN can be used to interpolate/extrapolate beyond the training data. We use the trained model to predict molecular structures along pathways connecting various states within the latent space. As an example, we present the sequence of molecular configurations "hallucinated" by the trained cWGAN along the linear activation pathway between the Int 1 and bent-closed states indicated by a dashed line marked by a \star in (b). The α and β subunits are respectively colored

in blue and red. (d) We restore atomistic detail to the cWGAN predicted configurations using targeted molecular dynamics (TMD). These calculations apply biasing forces to template an all-atom model onto each coarse-grained configuration along the predicted pathway. To see this figure in color, go online.

mapping. Furthermore, the AA simulations of the bent-closed, Int 1, and Int 2 structures each have 1770 amino acid residues, whereas the extended-open structure possesses 1748 residues, so developing a CG system representation also allows us to treat all four systems on an equal footing. Following a coarse graining procedure previously developed for the construction of an essential dynamics model of $\alpha_V \beta_3$ integrin, we adopted a 300-bead CG representation for the 1770 residues in the bent-closed, Int 1, and Int 2 systems (16) (Table S1). We adapted this 300-bead mapping for the 1748-residue extended-open system by slightly adjusting the residue-to-bead mapping within the leg regions of the dimer (Table S2). We verify that the $\alpha_V \beta_3$ integrin coarse graining is also appropriate to the $\alpha_{IIb}\beta_3$ integrin studied in this work by confirming that temporal trends in the root mean-square deviation (RMSD) and positional trends in the root mean-square fluctuation observed in the AA systems are well preserved under the CG mapping (Fig. S1). The 2450 configurations harvested for each intermediate state were converted into the 300-bead CG representation and the translationally and rotationally aligned RMSD computed between each pair of frames in preparation for the application of dMaps (Fig. 2 a).

Nonlinear manifold learning of a low-dimensional latent embedding using density-adaptive dMaps

The $3 \times 300 = 900$ -dimensional CG configurational space is large, high dimensional, and sparsely populated by the simulation data, making it challenging to plot physically plausible transition pathways between the intermediate states. We employ dMaps (40–45) to learn a unified low-

imensional embedding of the gross structure relations between and within the metastable states of molecular systems (Fig. 2 b). Conceptually, this nonlinear manifold learning technique performs a spectral decomposition of a random walk over the high-dimensional data points to learn a lowdimensional projection into the leading eigenvectors of this discrete diffusion process characterizing the large-scale, slow dynamical motions of the system (42). Mathematically, Euclidean distances in the low-dimensional projection approximate diffusion distances measuring the connectivity, and therefore kinetic proximity under the random walk, of states in the high-dimensional space. This preservation of the large-scale dynamical motions within the low-dimensional projection and the imputation of an effective kinetic proximity makes dMap embeddings particularly well-suited to the interpolative construction of putative transition pathways between intermediate states within the embedding. To account for the large variability in the sampling density of training data within and outside of the intermediate states, we find it useful to employ a density-adaptive variant of dMaps to smooth out these large fluctuations and to help to learn unified global embeddings of the distinct intermediate states (46).

In this work, the high-dimensional datapoints correspond to the 2450 molecular configurations in each of the four intermediate states (Fig. 1) that are sampled via AA MD simulations (13) and featurized with a 300-bead CG model. We employ RMSD as the similarity metric, as a translationally and rotationally invariant measure of configurational similarity (Fig. 2 *a*), and adopt a diffusion kernel bandwidth of $\varepsilon = e^{3.982}$ and a scaling exponent of $\alpha = 0.1$ in the density-adaptive dMaps (46). The eigenvalue spectrum and pairwise similarity matrix identify four leading eigenvectors corresponding to the large-scale, long-time configurational dynamical relaxations of the molecular system and inform the construction of a 4D low-dimensional projection of each 300-bead CG configuration **x** into the four leading nontrivial eigenvectors $\psi_{i=2,3,4,5}$ (Fig. S2).

Inverse mapping of the low-dimensional latent space to high-dimensional molecular structure using cWGANs

A deficiency of dMaps is the absence of an explicit functional mapping between the high-dimensional ambient space and the learned low-dimensional embedding (42). The inverse mapping from the low-dimensional latent space $\psi_{i=2,3,4,5}$ to the 300-bead CG configurational space **x** is required to interpolate activation pathways between intermediate states. In this work, we approximate the inverse mapping using a conditional Wasserstein GAN (cWGAN) with gradient penalty (47-49) (Fig. 2 c). The cWGAN generator $\mathcal{G}(\mathbf{z}|\psi_{i=2,3,4,5})$ is trained over the CG-MD training data to learn a conditional distribution of the configuration x using a d-dimensional white noise vector \mathbf{z} conditioned by the projection of \mathbf{x} within the latent space $\psi_{i=2,3,4,5}$. The cWGAN critic $C(\mathbf{x})$ is co-trained with the generator to learn the Wasserstein distance between a molecular configuration x from the training data with a corresponding latent space location $\psi_{i=2,3,4,5}$ and a synthetic configuration produced by the generator $\mathcal{G}(\mathbf{z}|\psi_{i=2,3,4,5})$. The networks representing the generator and critic are co-trained in an adversarial fashion to minimize the loss function,

$$\mathcal{L}_{\text{WGAN}} = \mathbb{E}_{\mathbf{z} \sim \mathbb{P}_{\mathbf{z}}, \psi_{i}} [\mathcal{C}(\mathcal{G}(\mathbf{z} | \psi_{i = 2, 3, 4, 5}))] \\ - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{training data}}} [\mathcal{C}(\mathbf{x}, \psi_{i = 2, 3, 4, 5})] + \lambda \mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}} \\ [(\parallel \nabla_{\hat{\mathbf{x}}} \mathcal{C}(\widehat{\mathbf{x}}, \psi_{i = 2, 3, 4, 5}) \parallel_{2} - 1)^{2}].$$
(1)

The first term, $\mathbb{E}_{z \sim \mathbb{P}_x, \psi_i} [C(\mathcal{G}(\mathbf{z}|\psi_{i=2,3,4,5}))]$ is the expectation of the critic $C(\tilde{\mathbf{x}})$ with $\tilde{\mathbf{x}}$ generated from the generator $\mathcal{G}(\mathbf{z}|\psi_{i=2,3,4,5})$. The second term $\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{training data}} [C(\mathbf{x}, \psi_{i=2,3,4,5})]$ is the expectation of the critic $C(\mathbf{x}, \psi_{i=2,3,4,5})$ given a real molecular configuration \mathbf{x} from the training data with latent space coordinates $\psi_{i=2,3,4,5}$. Minimizing the loss function minimizes the Wasserstein distance between the distributions of synthetic and real molecular configurations. The third term is a regularizing penalty to stabilize cWGAN training (48), the strength of which is controlled by the hyperparameter λ and which enforces the *L*2-norm of the gradient of the critic $(|| \nabla_{\tilde{\mathbf{x}}} C(\tilde{\mathbf{x}}, \psi_{i=2,3,4,5})||_2)$ to unity for input molecular configurations $\hat{\mathbf{x}}$. These input molecular configurations $\hat{\mathbf{x}} \sim \mathbb{P}(\hat{\mathbf{x}})$ are drawn uniformly at random along straight lines connecting a training configuration \mathbf{x} and a generated configuration $\tilde{\mathbf{x}}$ with latent space coordinates $\psi_{2,3,4,5}$.

In this work, we employ fully connected, feedforward neural networks for both the generator and the critic. The generator is modeled by a 132-256-256-256-900 network employing a 128D Gaussian noise vector $z\sim$ $\mathbb{P}_{\mathbf{z}}(\mathbf{z}) = \mathcal{N}(0,1)$ and a 4D conditioning vector $\psi_{i=2,3,4,5}$. The critic is modeled by a 904-256-256-256-1 network, where the input contains the 900D flattened coordinates of the 300-bead CG configuration translationally and rotationally aligned to a reference structure corresponding to Int 2 state observed at $t=20\ \text{ns}$ in AA MD simulations (13) and a 4D conditioning vector $\psi_{i=2,3,4,5}$. Inputs to both networks were normalized to [-1,1] (50). Sigmoid-weighted linear unit activations or swish functions (51) were employed in all hidden layers, and a tanh activation function was applied to the output of the generator. Batch normalization (52) was applied to each hidden layer of the generator. The generator and critic comprise, respectively, 393,000 and 363,000 trainable parameters. Training was performed using RMSprop (53) with a learning rate of 0.00005, momentum of 0.1, batch size of 100, and regularization parameter of $\lambda = 10$. Training was conducted for 1000 epochs over all $4 \times 2450 = 9800$ CG snapshots harvested over the four intermediate states, and the critic was updated five times for each generator update to balance training of the two networks. Models were constructed and trained using PyTorch lightning (54,55). Training curves for the cWGAN are illustrated in Fig. S3 illustrating convergence to an equilibrium at which the critic is unable to distinguish the synthetic molecular configurations produced by the generator from the training data. At this point, the generator has learned an excellent approximation for the inverse mapping from the 4D latent space to the 300-bead CG configurational space over the training data collected in the intermediate states.

CG to AA backmapping using TMD

The 300-bead CG configurations produced by the trained cWGAN were upgraded to AA resolution using TMD (56) (Fig. 2 d). Given a CG configuration produced by the cWGAN, we began from a candidate AA configuration, typically the Int 1 or extended-open state taken from Tong et al. (13), and we applied a moving harmonic biasing potential to minimize the RMSD between the AA configuration and the CG target (56,57). The RMSD between the configurations was computed under the CG mapping detailed in Tables S1 and S2 using the central residue for both this calculation and the application of the biasing forces. The potential bias can be expressed as $V = (\kappa(t)/2)(\text{RMSD} - \text{RMSD}_{\text{center}}(t))^2$, where $\kappa(t)$ refers to the harmonic potential force constant at time t, and the RMSD_{center}(t) indicates the reference RMSD at time t. In this work, we linearly scale $\kappa(t)$ from 0 to 20,000 kJ/mol/nm² and RMSD_{center}(t) from the initial RMSD to ~0 nm in a total time of 2.5 ns. We find that applying identical moving harmonic potentials to the α and β subunits separately performs better than applying a single moving harmonic potential for the entire $\alpha_{IIb}\beta_3$ integrin. An example application of the TMD procedure is illustrated in Fig. S4.

The AA simulations were conducted using GROMACS 2021.6 (58,59) patched with the PLUMED 2.8.1 plugin (60). Integrin was modeled using the CHARMM36m (61) force field with its default CHARMM-modified TIP3P water model (62). In contrast to Tong et al. (13), a lipid bilayer was not added to the system. Cubic boxes of dimensions 12.7 imes 12.7 imes26.0 nm³ for the Int 1 state and $14.8 \times 15.0 \times 42.5$ nm³ for the extendedopen state were employed. Sodium counterions were added to neutralize the net charge. In total, the Int 1 and extended-open AA systems contained, respectively, 132,794 and 310,258 water molecules and 61 and 59 sodium counterions. Real space cutoffs of 1.2 nm were employed for both van der Waals and Coulomb interactions (61). The force-switch modifier was employed to smoothly switch van der Waals forces between a cutoff of 1 nm and 1.2 nm. Long range electrostatic interactions were computed using the particle mesh Ewald method (63). Energy minimization and equilibration were performed as described in Tong et al. (13). In brief, steepest descent energy minimization was followed by NVT equilibration to 300 K and then NPT equilibration to 300 K and 1 bar using a Berendsen thermostat and barostat (64). TMD production runs were performed in the NPT ensemble at 300 K and 1 bar for 2.5 ns using the same leapfrog algorithm (65) utilized in equilibration runs with a time step of 2 fs. Hydrogen atom positions were constrained using the LINCS algorithm (66). A velocityrescale thermostat (67) with a time constant of 1 ps and Parinello-Rahman barostat (68) with time constant of 5 ps and compressibility of 4.5×10^{-5} bar⁻¹ were employed to regulate temperature and pressure.

RESULTS AND DISCUSSION

Nonlinear manifold learning of a low-dimensional latent space embedding of $\alpha_{IIb}\beta_3$ intermediate states

As the first step in our pipeline, we compute a low-dimensional embedding of the four states of the $\alpha_{IIb}\beta_3$ integrin heterodimer. In Fig. 3 *a*, we illustrate the pairwise RMSD matrix between the 300-bead CG representations of the





FIGURE 3 Nonlinear manifold learning of a lowdimensional embedding of the intermediate states of $\alpha_{IIb}\beta_3$ CG integrin using density-adaptive dMaps. (a) Pairwise RMSD matrix within and between the 300-bead CG representations of the 2450 configurations harvested from AA MD simulations in each of the four states-bent-closed, Int 1, Int 2, and extended-open. RMSD values are denoted by color corresponding to the scale bar. (b) Representative 300-bead CG configurations of each intermediate state. The α and β subunits of the $\alpha_{IIb}\beta_3$ integrin heterodimer are colored blue and red, respectively. (c) Learned 4D embedding of the $4 \times 2450 = 9800$ CG snapshots into the four leading nontrivial diffusion map eigenvectors $(\psi_2, \psi_3, \psi_4, \psi_5)$ visualized in all 2D projections. The snapshots corresponding to the bent-closed, Int 1, Int 2, and extended-open states are colored in gray, orange, green, and black, respectively (see legend), and exhibit a clear clustering within the low-dimensional projection. To see this figure in color, go online.

2450 configurations harvested from the AA MD simulations in the bent-closed, Int 1, Int 2, and extended-open states in Tong et al. (13). Representative CG structures from each intermediate state are presented in Fig. 3 *b*. The pairwise distance matrix exhibits a clear block structure evincing much closer configurational similarity of configurations within each intermediate state relative to between the states. The mean within-state pairwise RMSD across the four states is 0.92 nm compared with 4.4 nm between the states. This matrix structure is a consequence of the absence of any transitions between the four states in the simulation data wherein configurations within each state dynamically interconvert, whereas those between states do not. This suggests that we should anticipate gaps between the four states in our low-dimensional dMap embedding.

In Fig. 3 *c*, we present the projection of the $4 \times 2450 =$ 9800 snapshots into the 4D latent space learned using density-adaptive diffusion maps and spanned by the leading four nontrivial eigenvectors ($\psi_2, \psi_3, \psi_4, \psi_5$). The embedding demonstrates a clear clustering of the four states reflecting the closer structural similarity of configurations within each state relative to between states. We hypothesize, however, that the dMaps have learned a structurally meaningful embedding between the four states that will allow us to predict putative activation pathways between the states. Specifically, the ψ_2 projections of the embedding illustrate a clear separation of the bent-closed state from the other three states (Fig. 3 *c* (i–iii)), implying that the collective variable ψ_2 captures primarily the differences in the bent and extended con-

figurations. The Int 1, Int 2, and extended-open states are themselves separated out along the ψ_3 axis (Fig. 3 c (i, iv, and v)), and ψ_4 separates the Int 2 state from the other three (Fig. 3 c (ii, iv, and vi)). ψ_5 appears to offer little interstate separation but provides an axis over which the clusters spread out to accentuate intrastate differences between the constituent configurations (Fig. 3 c (iii, v, and vi)).

We gain further insight into the physical interpretation of some of the leading dMap eigenvectors by coloring the latent space with candidate physical observables. We find that ψ_2 is strongly correlated with the distance $d_{\beta_{propeller}} - \beta_{TD}$ between the $\beta_{propeller}$ and β_T domains (Fig. S5) and the distance $d_{head-tail}$ between the ectodomain headpiece and cytoplasmic tail helices (Fig. S6). This is consistent with our previous observation that ψ_2 characterizes differences between the bent-closed state and the more extended states. We observe a moderate correlation between ψ_3 and the distance between the cytoplasmic tail helices $d_{\alpha_{helix}} - \beta_{helix}$ (Fig. S7) and between the α - β subunits $d_{\alpha_{unit}} - \beta_{unit}$ (Fig. S8) exposing its role in characterizing the separation of the α and β subunits of $\alpha_{IIb}\beta_3$ integrin.

Generative decoding of CG configurations of $\alpha_{\rm IIb}\beta_3$ intermediate states from the learned latent space

Having learned a low-dimensional projection of the four states, we now train a cWGAN to learn the inverse mapping from the 4D latent space to the 300-bead CG configurational

Prediction of integrin activation paths

space. We train the cWGAN over all $4 \times 2450 = 9800$ CG snapshots harvested from the four states, and in Fig. 4 *a*, we illustrate its performance within the training set by comparing the pairwise RMSDs between each MD configuration and a corresponding synthetic cWGAN configuration generated by conditioning the trained cWGAN on the 4D latent space coordinates ($\psi_2, \psi_3, \psi_4, \psi_5$) of the training sample. The generation procedure is stochastic and only conditioned on the CG information contained in the four di-

mensions of the dMap embedding, so while we should not expect the cWGAN to generate exact copies of the training points, a well-trained generator should be capable of producing synthetic 300-bead configurations with the same gross structural configurations as the training data and low pairwise RMSDs. The four histograms of the pairwise RMSDs pertaining to each state demonstrate that the cWGAN has learned to approximate the inverse mapping and can accurately recapitulate the training data with a



FIGURE 4 Generative decoding of $\alpha_{IIb}\beta_3$ integrin heterodimer 300-bead CG structures conditioned on location within the 4D latent space embedding learned by density-adaptive dMaps. (*a*) Histograms of the pairwise RMSD between the 2,450 300-bead CG training configurations in each state—(i) bent-closed, (ii) Int 1, (iii) Int 2, and (iv) extended-open—and synthetic configurations generated by the trained cWGAN conditioned on the latent space coordinates of each training configuration. The low pairwise RMSD scores with means of the distribution centered at < 1 nm indicate that the cWGAN can accurately reproduce configurations in the training data for all four intermediate states. (*b*) Comparison of selected representative 300-bead CG training configurations from the MD simulations (*gray*) conducted in each intermediate states with the corresponding cWGAN synthetic configuration (*cyan*) generated using the latent space coordinates of the training sample. (*c*) Projection of the (i) 2450 MD configurations and (ii) 2450 cWGAN reconstructions from each of the four states into a space spanned by two physical collective variables quantifying the distance between tail cytoplasmic helices $d_{\alpha_{heftx}} - \beta_{hoftx}$ and the distance between the $\beta_{propeller}$ and β_T domain $d_{\beta_{propeller}} - \beta_{TD}$. (iii) Illustration of the structural regions corresponding to the two intramolecular distances $d_{\alpha_{heftx}} - \beta_{hoftx}$ and $d_{\beta_{propeller}} - \beta_{TD}$ on a molecular snapshot of bent-closed state. (*d*) Comparison of the distribution of the pairwise RMSD between each pair of configurations in the 2450 training configurations harvested from the MD simulations in each intermediate state and an identical number of synthetic configurations generated by conditioning the cWGAN on latent space coordinates generated from a 4D Gaussian mixture model with two components and mean and bandwidth fitted to the projection of the training data for each intermediate state into the latent space (Fig. S9). The similarity

Dasetty et al.



FIGURE 5 Interpolative paths between the four intermediate states of the $\alpha_{IIb}\beta_3$ integrin heterodimer within the 4D dMaps latent space. Ten intermediate points are generated using linear interpolation from which to seed conditional generation of synthetic molecular configurations using the trained cWGAN. The bent-closed, Int 1, Int 2, and extended-open state are colored gray, orange, green, and black, respectively. The image presents the paths within the $\psi_2 - \psi_3$ projection of the latent space; projections into all other 2D latent space elevations are presented in Fig. S10. To see this figure in color, go online.

mean pairwise RMSD of 0.86 nm. Moreover, the cWGAN performs well in all four states, indicating that the model has not overfit to any one structural ensemble and suggesting that it may be capable of predictively interpolating the structural ensembles between the states. In Fig. 4 b, we present a comparison of representative MD training configurations and the corresponding synthetic cWGAN structures along with the pairwise RMSD scores. In Fig. 4 c, we present a projection of the 2450 MD configurations and 2450 cWGAN reconstructions in each state into a space spanned by two physical collective variables quantifying the distance between tail cytoplasmic helices $d_{\alpha_{helix} - \beta_{helix}}$ and the distance between the $\beta_{propeller}$ and β_T domain $d_{\beta_{propeller} - \beta_{TD}}$. The good agreement between the two projections into these two key intramolecular distances that separate the four states provides further support that the cWGAN produces physically plausible molecular configurations.

We now assess the out-of-training generative capacity of the trained cWGAN by generating novel synthetic configurations over the latent space. In Fig. 4 *d*, we compare for each of the four states the distribution of pairwise RMSD scores between each pair of configurations in the 2450 training configurations harvested from the MD simulations in each state and an identical number of synthetic configurations generated by conditioning the cWGAN on latent space generated from a 4D Gaussian mixture model with two components and mean and bandwidth fitted to the projection of the training data for each state into the latent space (Fig. S9). The similarity of the two distributions in each case demonstrates that the trained cWGAN is capable of generating configurational ensembles of synthetic configurations within each state that are in good accord with those harvested from the MD simulation trajectories. Finally, we observe that generation of the cWGAN synthetic configurational ensembles is also computationally cheap, requiring only ~ 0.01 CPU-seconds to produce the 2450 novel configurations.

Prediction of integrin activation pathways between intermediate states

Having learned a low-dimensional latent space embedding of the four states and having trained a cWGAN to generate physically plausible 300-bead CG configurations within each state, we now seek to use these two models to chart pathways between the four states in the latent space and to generate configurations along these putative activation pathways that are challenging to access by both conventional molecular simulations and experiment. To do so, we first generate 10 intermediate points in the 4D latent space between every pair of states using linear interpolation. Fig. 5 shows the $\binom{4}{2} = 6$ interpolation pathways within the $\psi_2 - \psi_3$ projection of the latent space. Projections into all other 2D latent space projections are presented in Fig. S10.

Next, we use the intermediate points along the 4D interpolative pathways in the dMaps latent space to condition the generation of 300-bead CG configurations using the trained cWGAN. We illustrate the resulting interpolated configurations between each pair of the four states in Fig. 6. These interstate interpolations represent putative activation pathways between the states that can be examined to infer possible activation mechanisms. We pause here to caution against over-interpretation of these pathways: in the absence of the underlying free energy landscape, the actual (possibly nonlinear) minimum free energy route taken by each path through the latent space is not known, the preferred forward and backward pathways may not be coincident, and the accuracy of the cWGAN far from the training data cannot be assessed without collecting additional simulation data. We return to these issues and how the proposed paths may be used to seed additional more detailed, and more computationally burdensome, calculations in the conclusion. With these caveats in mind, we view the paths largely as model-guided interpolations between the states and examine them as putative activation pathways that can shed light on possible transition mechanisms without the application of artificial biasing forces, can make contact with the hypothesized switchblade and deadbolt mechanisms, and can be used to seed more detailed future studies.

Fig. 6 a illustrates a putative activation pathway between the bent-closed and Int 1 state that evinces a smooth

Prediction of integrin activation paths



FIGURE 6 Synthetic 300-bead CG configurations generated at 10 evenly spaced intervals along interpolative paths connecting each pair of the four states in the 4D dMaps latent space. (α -f) The symbols annotating each row of synthetic configurations correspond to those marking the interpolative pathways in Fig. 5. The α and β integrin subunits are illustrated in blue and red, respectively. To see this figure in color, go online.

extension of the ectodomain headpiece while the α and β subunits remain associated. Similarly, the pathway between the bent-closed and Int 2 state shown in Fig. 6 *b* exhibits extension of the ectodomain headpiece but now accompanied by a separation of the α and β subunits to ultimately form a cross shape. The Int 1 and Int 2 states are largely differentiated by the separation of the ectodomain headpiece that can be quantified by the distance between the $\beta_{propeller}$ and β_T domain $d_{\beta_{propeller}} - \beta_{TD}$ (cf. Fig. 4 *c*) and a crossing of the cytoplasmic tail helices. The putative activation pathway between Int 1 and Int 2 is illustrated in Fig. 6 *c*. In Fig. 6 *d*, we present the pathway between the Int 1 and extended-open state. These structures show that the extended-open state can be realized from Int 1 by simultaneous changes in the separation of the headpiece and tail helices of the α

and β subunits. On the other hand, the Int 2 state already contains a separated headpiece and, as illustrated in Fig. 6 *e*, transitions to the extended-open state along its putative activation pathway by a gradual separation of tail helices. Although the activation pathway between the bent-closed and extended-open states is anticipated to pass through intermediate metastable states, for the sake of completeness, we illustrate the putative direct pathway between these two states in Fig. 6*f*. Despite the large structural differences, our approach generates what appear to be physically plausible intermediate configurations that involve an extension of the headpiece followed by separation of both headpiece and tail helices of the α and β subunits. We compare each of the intermediate configurations along this direct pathway to the Int 1 and Int 2 states to assess the degree to which this

Dasetty et al.

path may pass through these previously characterized intermediates. As illustrated in Fig. S11, none of the intermediate configurations possess an RMSD to the Int 1 and Int 2 configurations lower than, respectively, 3.64 nm and 2.96 nm. This analysis indicates that the Int 1 and Int 2 states do not appear to be intermediates on the direct bent-closed to extended-open pathway reported in Fig. 6 f.

To make contact with the proposed switchblade and deadbolt activation models (25,23), we now project the training configurations and those generated by the cWGAN into the 2D space spanned by the physical distances $d_{\beta_{propeller} - \beta_{TD}}$ and $d_{\alpha_{helix}-\beta_{helix}}$ (16) (Fig. 7). In the switchblade model, insideout activation of the extended-open conformation is realized from the bent-closed conformation first by separation of extended legs followed by opening of the head ectodomains of both α and β subunits bent at the genu or linker regions. This mechanism, similar to opening of a pocketknife, can be traced as first a horizontal transition in $d_{\alpha_{helix} - \beta_{helix}}$ followed by vertical transition in $d_{\beta_{propeller} - \beta_{TD}}$. In contrast, the deadbolt model requires progressive loss of interactions between the β -T domain (the deadbolt region) and the β -I domain (lock region) in the β subunit before the separation of the legs. This can be traced as first a vertical transition in $d_{\beta_{propeller} - \beta_{TD}}$ and then a horizontal transition in $d_{\alpha_{helix} - \beta_{helix}}$. This deadbolt path passes through both the extended-intermediate structures.



(As described in Hanein and Volkmann (28), these intermediates could also be part of the light switch mechanism involving subtler rearrangements of the cytoplasmic tail helix of the β subunit by the weakened cytoplasmic tail helix interactions when inserted into the membrane.) Using the trained cWGAN, we use the synthetic configurations generated over the 4D $(\psi_2, \psi_3, \psi_4, \psi_5)$ latent space and projected into the 2D $d_{\alpha_{helix}-\beta_{helix}}$ - $d_{\beta_{propeller}-\beta_{TD}}$ space to trace structures along the switchblade and deadbolt paths and create putative activation pathways (Fig. 7). These pathways are generated by constructing cubic Bézier curves following the general trends of the switchblade and deadbolt paths within this 2D physical space (16). Visual analysis of these paths exposes plausible structures for the switchblade and deadbolt models. In particular, the generated structures along the switchblade path illustrate a gradual separation of the cytoplasmic tails followed by the opening of the head ectodomains of both α and β subunits from the leg regions and separation of α and β subunits. This partly follows the plausible intermediate structures generated between the bent-closed and extended-open states (Fig. 6 f), but the separation of legs precedes the extension and separation of headpieces of α and β subunits. On the other hand, the generated structures along the deadbolt path effectively combine the plausible intermediate structures generated between bent-closed and Int 1 state (Fig. 6 a), Int 1 and Int 2

> FIGURE 7 Synthetic 300-bead CG configurations generated along interpolative paths in the $d_{\alpha_{helix} - \beta_{helix}} - d_{\beta_{propeller} - \beta_{TD}}$ space following the proposed switchblade and deadbolt activation pathways. The switchblade path describes inside-out activation from the bent-closed to extended-open state following first a horizontal excursion along $d_{\alpha_{helix} - \beta_{helix}}$ (separation of legs) followed by a vertical transition in $d_{\beta_{propeller} - \beta_{TD}}$ (ectodomain opening). The deadbolt path, in contrast, first follows a vertical path in $d_{\beta_{\text{propeller}} - \beta_{TD}}$ (unlocking of deadbolt region), passes through the extended-intermediate structures Int 1 and Int 2, and then moves horizontally in $d_{\alpha_{helix} - \beta_{helix}}$ (separation of legs). Gray, orange, green, and black colored markers correspond to training configurations in the bent-closed, Int 1, Int 2, and extended-open states. The light purple colored markers correspond to the projection into the $d_{\alpha_{helix} - \beta_{helix}} - d_{\beta_{propeller} - \beta_{TD}}$ space of 90,000 randomly sampled locations in the learned 4D dMaps latent space, each of which may be used to condition the cWGAN to generatively decode a 300-bead CG structure. Since the sampling of these points was performed randomly, their distribution in this space cannot be interpreted as exposing thermodynamically preferred regions of space. Snapshots are rendered at the blue markers in the intermediate space along the arrows highlighting potential switchblade and deadbolt paths and are arranged around the plot in order. The blue markers are selected as the closest light purple point and its corresponding

300-bead decoded CG structure at 10 evenly spaced intervals along three cubic Bézier polynomials trained using four points selected within this 2D physical space. One Bézier curve corresponds to the switchblade mechanism passing along the bent-closed state to extended-open state (10 blue points). Two Bézier curves correspond to the deadbolt mechanism passing first along the bent-closed state to the Int 1 state followed by the Int 1 state to the extended-open state (17 blue points). The α and β integrin subunits in the snapshots are illustrated in blue and red, respectively. To see this figure in color, go online.

Prediction of integrin activation paths

state (Fig. 6 *c*), and Int 2 and extended-open state (Fig. 6 *e*). This involves first a gradual separation of head ectodomain headpiece from the leg regions, a minor separation of α and β subunits with a simultaneous change in the arrangement of a tail helices to a cross shape, and finally the separation of tail helices to result in the extended-open state.

Finally, we backmap the 300-bead CG pathways generated in Fig. 6 to AA resolution using TMD simulations to anneal an AA model to each CG target configuration. As is typically the case with all backmapping approaches, the AA configurations generated by this procedure are unlikely to be fully energetically relaxed, but they do serve two valuable purposes in providing higher structural resolution of the putative activation pathways and in furnishing initial AA structures for subsequent molecular simulations. We present in Fig. 8 the backmapped AA structures corresponding to each of the intermediate 300-bead CG structures shown in Fig. 6. The relatively high resolution of the CG model means that the additional structural insights obtained from the AA backmapping are largely limited to the relative arrangements of the different domains in the α and β subunits of integrin. For example, the backmapped AA structures between bentclosed and Int 1 state (Fig. 8 *a*) provide insights into the progressive loss of contacts between the headpiece (β -I and hybrid domains) and lower leg regions (linker and β -T domains) in the β subunit and a gradual extension of the headpiece (β -propeller and thigh domain) in the α subunit indicating their plausible role in bent-closed and Int 1 state



FIGURE 8 Prediction of intermediate AA structures of integrin generated with the trained cWGAN between every pair of the four intermediate states in the latent space that are shown in Fig. 6. (a-f) The symbols annotating each row of synthetic configurations correspond to those marking the interpolative pathways in Fig. 5. The α and β integrin subunits are illustrated in blue and red, respectively. To see this figure in color, go online.

Dasetty et al.

transition mechanism. Similarly, the gradual changes in the curvature of the ectodomain (headpiece and lower leg) of α subunit, loss of contacts between linker region in β subunit and genu or thigh domain of α subunit, and rearrangement of cytoplasmic tail helical domains are revealed in the back-mapped AA structures between Int 1 and Int 2 state (Fig. 8 *c*). Although we do not do so here due to the high computational expense, these AA structures can be utilized to launch new atomistic simulations to perform high-resolution exploration of the thermodynamics and dynamics of $\alpha_{IIb}\beta_3$ integrin activation using unbiased or biased sampling. To facilitate this and enable further research into integrin activation mechanism, we make both the trained cWGAN model as well as the predicted CG and AA intermediate structures available at GitHub: https://github.com/Ferg-Lab/integrin_molgen.git.

CONCLUSION

In this work, we have developed a computational method to interpolate between conformational states of the $\alpha_{IIb}\beta_3$ integrin heterodimer as a large multimolecular system for which direct molecular simulation of the transition process is a computationally prohibitive rare event and experimental measurement is challenging due to the transient nature of the activated states. The method integrates MD, nonlinear manifold learning, and deep generative models to learn a unified lowdimensional embedding capturing the gross structural relationships between the states and interpolate putative activation pathways between them. We train the models over unbiased simulation data, generate these pathways at both 300-bead CG and AA resolution, develop mechanistic insight into the putative activation paths, and relate them to the hypothesized switchblade and deadbolt activation mechanisms. This work reports, to the best of our knowledge, the first proposed continuous activation pathways of the integrin dimer without the use of artificial biasing forces that can deform and distort the structures along the putative activation path. The degree to which the proposed pathways follow the preferred minimum free energy routes through configurational space can, however, only be determined with additional simulations, and it is a natural, although computationally very demanding, extension of this work to seed new CG or AA simulations to refine the pathways. In particular, we propose that the string method is very well-suited to relaxing these initial pathways into the free energy minimum channel (69,70). We anticipate that the structures along the putative pathways may also offer useful guidance to experimental work seeking to isolate metastable states of the integrin dimer along its activation course. Finally, we demonstrated our approach for dynamical interconversion of various metastable states of the $\alpha_{IIb}\beta_3$ integrin heterodimer, but we anticipate that our approach offers a generic and transferable multiscale technique to predict transition pathways for other biomolecular systems for which intermediate states are well characterized but transition regions are not.

SUPPORTING MATERIAL

Supporting material can be found online at https://doi.org/10.1016/j.bpj. 2023.12.009.

AUTHOR CONTRIBUTIONS

A.L.F., T.C.B., and S.D. designed the research. S.D. performed the research. S.D. and A.L.F. analyzed the data and wrote the article. S.D., T.C.B., and A.L.F. edited the manuscript.

ACKNOWLEDGMENTS

This work was supported by the US Department of Energy, Office of Science, Basic Energy Sciences, under award #DE-SC0023318. We gratefully acknowledge computing time provided by the University of Chicago Research Computing Center (https://rcc.uchicago.edu), the University of Chicago high-performance GPU-based cyberinfrastructure supported by the National Science Foundation under grant no. DMR-1828629, and the Texas Advanced Computing Center (TACC) at The University of Texas at Austin (http://www.tacc.utexas.edu). We thank Dr. Gregory Voth for helpful discussions and insights.

DECLARATION OF INTERESTS

A.L.F. is a co-founder and consultant of Evozyne, Inc., and a co-author of US Patent Applications 16/887,710 and 17/642,582, US Provisional Patent Applications 62/853,919, 62/900,420, 63/314,898, and 63/479,378, and International Patent Applications PCT/US2020/035,206 and PCT/US2020/050,466.

REFERENCES

- Giancotti, F. G., and E. Ruoslahti. 1999. Integrin signaling. Science. 285:1028–1032.
- Hynes, R. O. 2002. Integrins: Bidirectional, Allosteric Signaling Machines. *Cell*. 110:673–687. https://www.sciencedirect.com/science/ article/pii/S0092867402009716.
- Geiger, B., A. Bershadsky, ..., K. M. Yamada. 2001. Transmembrane crosstalk between the extracellular matrix and the cytoskeleton. *Nat. Rev. Mol. Cell Biol.* 2:793–805.
- Kechagia, J. Z., J. Ivaska, and P. Roca-Cusachs. 2019. Integrins as biomechanical sensors of the microenvironment. *Nat. Rev. Mol. Cell Biol.* 20:457–473.
- Barczyk, M., S. Carracedo, and D. Gullberg. 2010. Integrins. *Cell Tissue Res.* 339:269–280.
- Kim, C., F. Ye, and M. H. Ginsberg. 2011. Regulation of integrin activation. *Annu. Rev. Cell Dev. Biol.* 27:321–345.
- 7. Kumar, C. C. 1998. Signaling by integrin receptors. *Oncogene*. 17:1365–1373.
- Shattil, S. J., C. Kim, and M. H. Ginsberg. 2010. The final steps of integrin activation: the end game. *Nat. Rev. Mol. Cell Biol.* 11:288–300.
- Giancotti, F. G. 2007. Targeting integrin β4 for cancer and anti-angiogenic therapy. *Trends Pharmacol. Sci.* 28:506–611.
- Paavolainen, O., and E. Peuhu. 2021. Integrin-mediated adhesion and mechanosensing in the mammary gland. *Semin. Cell Dev. Biol.* 114:113–125. https://www.sciencedirect.com/science/article/pii/S1084952120301671.
- Humphrey, J. D., E. R. Dufresne, and M. A. Schwartz. 2014. Mechanotransduction and extracellular matrix homeostasis. *Nat. Rev. Mol. Cell Biol.* 15:802–812.

Prediction of integrin activation paths

- Xu, X.-P., E. Kim, ..., D. Hanein. 2016. Three-dimensional structures of full-length, membrane-embedded human αIIbβ3 integrin complexes. *Biophys. J.* 110:798–809.
- Tong, D., N. Soley, ..., T. C. Bidone. 2023. Integrin αIIbβ3 intermediates: From molecular dynamics to adhesion assembly. *Biophys. J.* 122:533–543.
- Humphrey, W., A. Dalke, and K. Schulten. 1996. VMD Visual Molecular Dynamics. J. Mol. Graph. 14:33–38.
- 15. Cai, T.-Q., S. K. Law, ..., S. D. Wright. 1995. Reversible Inactivation of Purified Leukocyte Integrin CR3 (CD11b/CD18, β m β 2) by Removal of Divalent Cations from a Cryptic Site. *Cell Adhes. Commun.* 3:399–406.
- Bidone, T. C., A. Polley, ..., G. A. Voth. 2019. Coarse-grained simulation of full-length integrin activation. *Biophys. J.* 116:1000–1010.
- Smith, J. W., R. S. Piotrowicz, and D. Mathis. 1994. A mechanism for divalent cation regulation of beta 3-integrins. J. Biol. Chem. 269:960–967.
- Mehrbod, M., S. Trisno, and M. R. K. Mofrad. 2013. On the activation of integrin αIIbβ3: outside-in and inside-out pathways. *Biophys. J.* 105:1304–1315.
- Bledzka, K., J. Qin, and E. F. Plow. 2019. Integrin αIIbβ3. Platelets 227–241.
- Kashiwagi, H., M. Shiraga, ..., Y. Tomiyama. 2004. Activation of integrin αIIbβ3 in the glycoprotein Ib-high population of a megakaryocytic cell line, CMK, by inside-out signaling. *J. Thromb. Haemostasis*. 2:177–186.
- Ye, F., C. Kim, and M. H. Ginsberg. 2012. Reconstruction of integrin activation. *Blood, The Journal of the American Society of Hematology*. 119:26–33.
- Luo, B.-H., and T. A. Springer. 2006. Integrin structures and conformational signaling. *Curr. Opin. Cell Biol.* 18:579–586.
- 23. Ye, F., J. Liu, ..., K. A. Taylor. 2008. Integrin αIIbβ3 in a membrane environment remains the same height after Mn2+ activation when observed by cryoelectron tomography. J. Mol. Biol. 378:976–986.
- 24. Takagi, J., B. M. Petre, ..., T. A. Springer. 2002. Global conformational rearrangements in integrin extracellular domains in outside-in and inside-out signaling. *Cell*. 110:599–611.
- Zhu, J., B. Boylan, ..., T. A. Springer. 2007. Tests of the extension and dead bolt models of integrin activation. *J. Biol. Chem.* 282:11914– 11920.
- Xiong, J.-P., T. Stehle, ..., M. A. Arnaout. 2003. New insights into the structural basis of integrin activation. *Blood*. 102:1155–1159.
- Bunch, T. A. 2010. Integrin αIIbβ3 Activation in Chinese Hamster Ovary Cells and Platelets Increases Clustering Rather than Affinity. *J. Biol. Chem.* 285:1841–1849.
- 28. Hanein, D., and N. Volkmann. 2018. Conformational Equilibrium of Human Platelet Integrin Investigated by Three-Dimensional Electron Cryo-Microscopy. *In* Membrane Protein Complexes: Structure and Function. J. R. Harris and E. J. Boekema, eds Springer Singapore, Singapore, pp. 353–363.
- 29. Niu, G., and X. Chen. 2011. Why integrin as a primary target for imaging and therapy. *Theranostics*. 1:30–47.
- Chen, Y., L. A. Ju, ..., C. Zhu. 2019. An integrin αIIbβ3 intermediate affinity state mediates biomechanical platelet aggregation. *Nat. Mater.* 18:760–769.
- Nurden, A. T., X. Pillois, ..., P. Nurden. 2011. Glanzmann thrombasthenia-like syndromes associated with macrothrombocytopenias and mutations in the genes encoding the αIIbβ3 integrin. Semin. Thromb. Hemost. 37:698–706.
- 32. Kaneva, V. N., A. A. Martyanov, ..., A. N. Sveshnikova. 2019. Platelet integrin αIIbβ3: mechanisms of activation and clustering; involvement into the formation of the thrombus heterogeneous structure. *Biochem.* (*Moscow*), Suppl. Ser. 13:97–110.
- Staunton, D. E., M. L. Lupher, ..., W. M. Gallatin. 2006. Targeting integrin structure and function in disease. *Adv. Immunol.* 91:111–157.

- Gu, Y., B. Dong, ..., Y. Cui. 2023. The challenges and opportunities of αvβ3-based therapeutics in cancer: from bench to clinical trials. *Pharmacol. Res.* 106694
- Cox, D., M. Brennan, and N. Moran. 2010. Integrins as therapeutic targets: lessons and opportunities. *Nat. Rev. Drug Discov.* 9:804–820.
- Tvaroška, I., S. Kozmon, and J. Kóňa. 2023. Molecular modeling insights into the structure and behavior of integrins: a review. *Cells*. 12:324.
- 37. Kulke, M., and W. Langel. 2020. Molecular dynamics simulations to the bidirectional adhesion signaling pathway of integrin $\alpha V\beta 3$. *Proteins*. 88:679–688.
- Wang, K., and Z. Li. 2023. Steered molecular dynamics simulation of force triggering the integrin αIIbβ3 extension via its ligand. *Eur. Phys. J. Spec. Top.* 1–9.
- **39.** Jo, S., J. B. Lim, ..., W. Im. 2009. CHARMM-GUI Membrane Builder for mixed bilayers and its application to yeast membranes. *Biophys. J.* 97:50–58.
- Nadler, B., S. Lafon, ..., R. Coifman. 2005. Diffusion maps, spectral clustering and eigenfunctions of Fokker-Planck operators. *Adv. Neural Inf. Process. Syst.* 18
- Ferguson, A. L., A. Z. Panagiotopoulos, ..., I. G. Kevrekidis. 2010. Systematic determination of order parameters for chain dynamics using diffusion maps. *Proc. Natl. Acad. Sci. USA*. 107:13597–13602.
- Ferguson, A. L., A. Z. Panagiotopoulos, ..., P. G. Debenedetti. 2011. Nonlinear dimensionality reduction in molecular simulation: The diffusion map approach. *Chem. Phys. Lett.* 509:1–11.
- 43. Coifman, R. R., S. Lafon, ..., S. W. Zucker. 2005. Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. *Proc. Natl. Acad. Sci. USA*. 102:7426–7431.
- 44. Coifman, R. R., and S. Lafon. 2006. Diffusion maps. Appl. Comput. Harmon. Anal. 21:5–30.
- Belkin, M., and P. Niyogi. 2003. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Comput.* 15:1373–1396.
- 46. Wang, J., M. A. Gayatri, and A. L. Ferguson. 2017. Mesoscale simulation and machine learning of asphaltene aggregation phase behavior and molecular assembly landscapes. J. Phys. Chem. B. 121:4923–4944.
- Sidky, H., W. Chen, and A. L. Ferguson. 2020. Molecular latent space simulators. *Chem. Sci.* 11:9459–9467.
- Gulrajani, I., F. Ahmed, ..., A. C. Courville. 2017. Improved Training of Wasserstein GANs. *In* Advances in Neural Information Processing Systems, *Vol.30*. I. Guyon, U. V. Luxburg, and ..., R. Garnetteds. Curran Associates, Inc.. https://proceedings.neurips.cc/paper_files/ paper/2017/file/892c3b1c6dccd52936e27cbd0ff683d6-Paper.pdf.
- Arjovsky, M., S. Chintala, and L. Bottou. 2017. Wasserstein Generative Adversarial Networks. *In* Proceedings of the 34th International Conference on Machine Learning. PMLR, Volume 70 of *Proceedings of Machine Learning Research*. D. Precup and Y. W. Teh, eds, pp. 214–223. https://proceedings.mlr.press/v70/arjovsky17a.html.
- Scikit-Learn. sklearn.preprocessing.MinMaxScaler. https://scikit-learn. org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler. html#sklearn.preprocessing.MinMaxScaler. (Accessed July 2023), accessed.
- Ramachandran, P., B. Zoph, and Q. V. Le. 2017. Searching for activation functions. Preprint at arXiv. https://doi.org/10.48550/arXiv.1710. 05941.
- 52. Ioffe, S., and C. Szegedy. 2015. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *In* Proceedings of the 32nd International Conference on International Conference on Machine Learning JMLR.org, pp. 448–456, ICML'15.
- Hinton, G., N. Srivastava, and K. Swersky. Neural Networks for Machine Learning (Lecture 6a): Overview of Mini-Batch Gradient Descent. http://www.cs.toronto.edu/~hinton/coursera/lecture6/lec6. pdf. (Accessed July 2023), accessed.
- 54. Falcon, W. 2019. PyTorch Lightning.

Dasetty et al.

- Shmilovich, K., and A. L. Ferguson. 2023. Generative models for conditional molecular structure generation. https://github.com/Ferg-Lab/ molgen.
- Schlitter, J., M. Engels, and P. Krüger. 1994. Targeted molecular dynamics: a new approach for searching pathways of conformational transitions. J. Mol. Graph. 12:84–89. https://www.sciencedirect.com/ science/article/pii/0263785594800723.
- Hénin, J., T. Lelièvre, ..., L. Delemotte. 2022. Enhanced Sampling Methods for Molecular Dynamics Simulations [Article v1. 0]. *Living J. Comput. Mol. Sci.* 4:1583. https://livecomsjournal.org/index.php/ livecoms/article/view/v4i1e1583.
- Lindahl, E., M. J. Abraham, ..., D. van der Spoel. 2022. GROMACS 2021.6 Source Code.
- Abraham, M. J., T. Murtola, ..., E. Lindahl. 2015. GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX*. 1–2:19–25. https://www. sciencedirect.com/science/article/pii/S2352711015000059.
- Tribello, G. A., M. Bonomi, ..., G. Bussi. 2014. PLUMED 2: New feathers for an old bird. *Comput. Phys. Commun.* 185:604–613. https://www.sciencedirect.com/science/article/pii/ S0010465513003196.
- Huang, J., S. Rauscher, ..., A. D. MacKerell, Jr. 2017. CHARMM36m: an improved force field for folded and intrinsically disordered proteins. *Nat. Methods.* 14:71–73.
- Neria, E., S. Fischer, and M. Karplus. 1996. Simulation of activation free energies in molecular systems. J. Chem. Phys. 105:1902–1921.

- Darden, T., D. York, and L. Pedersen. 1993. Particle mesh Ewald: An N.log(N) method for Ewald sums in large systems. J. Chem. Phys. 98:10089–10092.
- Berendsen, H. J. C., J. P. M. Postma, ..., J. R. Haak. 1984. Molecular dynamics with coupling to an external bath. J. Chem. Phys. 81:3684–3690.
- Hockney, R., S. Goel, and J. Eastwood. 1974. Quiet high-resolution computer models of a plasma. J. Comput. Phys. 14:148–158. https:// www.sciencedirect.com/science/article/pii/0021999174900102.
- Hess, B., H. Bekker, ..., J. G. E. M. Fraaije. 1997. LINCS: A linear constraint solver for molecular simulations. *J. Comput. Chem.* 18:1463–1472.
- Bussi, G., D. Donadio, and M. Parrinello. 2007. Canonical sampling through velocity rescaling. J. Chem. Phys. 126, 014101.
- Parrinello, M., and A. Rahman. 1981. Polymorphic transitions in single crystals: A new molecular dynamics method. *J. Appl. Phys.* 52:7182–7190.
- 69. E, W., W. Ren, and E. Vanden-Eijnden. 2007. Simplified and improved string method for computing the minimum energy paths in barriercrossing events. J. Chem. Phys. 126, 164103.
- Vanden-Eijnden, E., and M. Venturoli. 2009. Revisiting the finite temperature string method for the calculation of reaction tubes and free energies. J. Chem. Phys. 130, 194103.