

# Scalable computation of anisotropic vibrations for large macromolecular assemblies

---

Received: 31 August 2023

Accepted: 2 April 2024

Published online: 24 April 2024

 Check for updates

---

Jordy Homing Lam <sup>1,2,3</sup>, Aiichiro Nakano <sup>1,4,5</sup>  & Vsevolod Katritch <sup>1,2,3,6</sup> 


---

The Normal Mode Analysis (NMA) is a standard approach to elucidate the anisotropic vibrations of macromolecules at their folded states, where low-frequency collective motions can reveal rearrangements of domains and changes in the exposed surface of macromolecules. Recent advances in structural biology have enabled the resolution of megascale macromolecules with millions of atoms. However, the calculation of their vibrational modes remains elusive due to the prohibitive cost associated with constructing and diagonalizing the underlying eigenproblem and the current approaches to NMA are not readily adaptable for efficient parallel computing on graphic processing unit (GPU). Here, we present eigenproblem construction and diagonalization approach that implements level-structure bandwidth-reducing algorithms to transform the sparse computation in NMA to a globally-sparse-yet-locally-dense computation, allowing batched tensor products to be most efficiently executed on GPU. We map, optimize, and compare several low-complexity Krylov-subspace eigensolvers, supplemented by techniques such as Chebyshev filtering, sum decomposition, external explicit deflation and shift-and-inverse, to allow fast GPU-resident calculations. The method allows accurate calculation of the first 1000 vibrational modes of some largest structures in PDB (> 2.4 million atoms) at least 250 times faster than existing methods.

The Normal Mode Analysis (NMA) is a standard approach to derive motions from static snapshots of a macromolecular structure. As a computational probe to shape-changing motions, the analysis has found wide applicability in the refinement and fitting of macromolecular structures<sup>1,2</sup> and in the simulations of functional motions when coupled with enhanced sampling techniques<sup>3,4</sup>. Many of these predicted modes of motion align consistently with available experimental data in kinases<sup>5</sup>, ion channels<sup>6</sup> and transporters<sup>7</sup>. Such successful applications have led to an increasing appreciation of the

interdependence among biological structure, dynamics, and function<sup>8</sup>. In NMA, the input macromolecular structure is assumed to be in a conformational minimum of its potential energy. Spatial arrangements of atoms (or atom groups like residues) and forces among them are then assimilated into the Hessian matrix and the equation of motion is analyzed under the harmonic approximation. The outcome of NMA are eigenvectors, representing sets of optimal displacements for each atom in the system, ranked by their ascending eigenvalues that reflect increasing strain along those directions. Early attempts of NMA on

---

<sup>1</sup>Department of Quantitative and Computational Biology, University of Southern California, Los Angeles, CA, USA. <sup>2</sup>Bridge Institute and Michelson Center for Convergent Biosciences, University of Southern California, Los Angeles, CA, USA. <sup>3</sup>Center for New Technologies in Drug Discovery and Development, University of Southern California, Los Angeles, CA, USA. <sup>4</sup>Department of Physics and Astronomy, University of Southern California, Los Angeles, CA, USA. <sup>5</sup>Department of Computer Science, University of Southern California, Los Angeles, CA, USA. <sup>6</sup>Department of Chemistry, University of Southern California, Los Angeles, CA, USA.  e-mail: [anakano@usc.edu](mailto:anakano@usc.edu); [katritch@usc.edu](mailto:katritch@usc.edu)

macromolecules<sup>9–12</sup> were derived from all-atom potentials and full dense Hessian matrices. This approach to NMA is still the golden standard in analyzing the shear and hinge motions of proteins with sizes typically less than 5000 atoms<sup>13</sup>. However, applying NMA in this form to larger macromolecules is difficult as their energy minimizations are prone to overstepping, rendering unphysical modes<sup>14,15</sup>. In a pioneering work by Tirion<sup>15</sup>, this requirement was relaxed by replacing the detailed all-atom potentials with an elastic network of pairwise Hookean potential, connecting atoms within a neighborhood boundary. The NMA analysis then proceeds by assuming that the provided structure was minimized under experimental conditions. As such, crystal coordinates, which lack hydrogens and/or sidechains, can be analyzed directly in absence of a forcefield. This practical form of NMA is now commonly known as the elastic network model (ENM). Despite its simplicity, as demonstrated convincingly by Hinsen<sup>16,17</sup> and Sanejouand<sup>18</sup>, the low-frequency ENM modes were shown to agree with the fluctuations observed in X-ray crystallography as well as with NMAs performed under standard all-atom potentials, as long as only low-frequency modes are concerned. The ENM also found major applications in Cryogenic Electron Microscopy (Cryo-EM) and Cryogenic Electron Tomography (Cryo-ET) for the flexible refinement of the composite density maps<sup>19–21</sup> and the flexible fitting of atomic models<sup>22–24</sup> when EM volume or atomic templates are available. These applications of NMA are especially compelling for large macromolecular complexes, when other methods are computationally expensive. However, computing the NMA, even in this practical ENM form, remains challenging for large macromolecular systems of more than 1,000,000 particles. Namely, the storage, the construction and the diagonalization of the Hessian matrix all create imminent difficulties in terms of time and memory complexities as the size of the matrix increases quadratically with the number of atoms  $N$ . While ideally the memory complexity of NMA only increases linearly, the linear factors due to packing density of atoms as well as the cubic time complexity to diagonalize the Hessian matrix present a major bottleneck in calculations (See Methods). As such, many innovative approaches, including Anisotropic Network Model (ANM)<sup>15,25</sup>, Rotation-Translation Block (RTB) method<sup>26</sup> and Block Normal Mode (BNM)<sup>27</sup>, were undertaken to eliminate the degrees of freedom (DOF) to be differentiated in these larger systems, hence reducing  $N$  for 8–10 times. Agreement of these methods with experimental data<sup>18,25</sup> are satisfactory, though the displacement information of the discounted atoms were lost and the appropriate level of granularity, especially in presence of elongated (e.g., lipids) or planar (e.g., aromatics) chemical moieties, are hard to be determined. To apply NMA without excessive coarse-graining, the development of faster numerical recipes is necessary<sup>28,29</sup>. More recently, the choices on diagonalization algorithms were examined on a hyperthreaded machine in the work of Koehl<sup>30</sup>. This ultimately allowed the first 100 ENM modes of a ZIKV virus (PDBID: 5IZ7), with around 800 thousand atoms and around 300 million nonzero entries in its Hessian, to be calculated within an hour when the Jacobi-Davidson Method (JDM)<sup>31,32</sup> were filtered with an 80-degree low-pass Chebyshev polynomial of the Hessian<sup>30</sup>. However, the current state-of-the-art is still far from handling megascale systems with more than a million atoms and billions of non-zero entries in the Hessian.

In this work, we designed several synergistic algorithms to fully engage dense graphic processing unit (GPU) kernels in both construction and diagonalization of the Hessian in atomic NMAs. Specifically, we developed a level-structure algorithm to minimize the bandwidth of the Hessian, prior to its construction, and thereby converting its sparse computation into a globally-sparse-yet-locally-dense computation, which enables the batched execution of tensor products on GPUs. We also mapped, optimized and compared several low-complexity Krylov-subspace eigensolvers, supplemented by techniques such as Chebyshev filtering<sup>33–35</sup>, sum decomposition, external explicit deflation<sup>36</sup> and shift-and-inverse, to allow fast calculations on a

GPU device. The level-structure algorithm produces an isomorph of the original elastic network, for which its unpermuted eigenpairs can be recovered in linear complexity using the bijection mapping. The implementation of this INCHING (“Isomorphic Nma Calculations Harnessing 1 Necessary Gpu”) algorithm presented here was benchmarked on macromolecules from the Protein Data Bank<sup>37</sup> with sizes up to 2.4 million atoms. Compared to other existing methods, 250–370 times speedup in throughput was achieved while maintaining residual error under  $10^{-12}$ . The utility of the INCHING method was also demonstrated through examples, including the largest experimentally resolved atomistic structure of a mature HIV-1 capsid<sup>38</sup> (PDBID:3J3Q) with 2.4 million atoms and 1.6 billion non-zero entries in its Hessian. Using our INCHING program, we were able to resolve its first 64 atomic normal modes within 44 min of wall-clock time on a single NVIDIA® A100 Tensor Core GPU and its first 1000 atomic normal modes within 63 h. Fast, accurate, and highly scalable GPU-optimized implementation of NMA approach will find practical applications in conformational analysis of large and dynamic macromolecular complexes, and facilitate their refinement from cryo-EM/cryo-ET data.

## Results

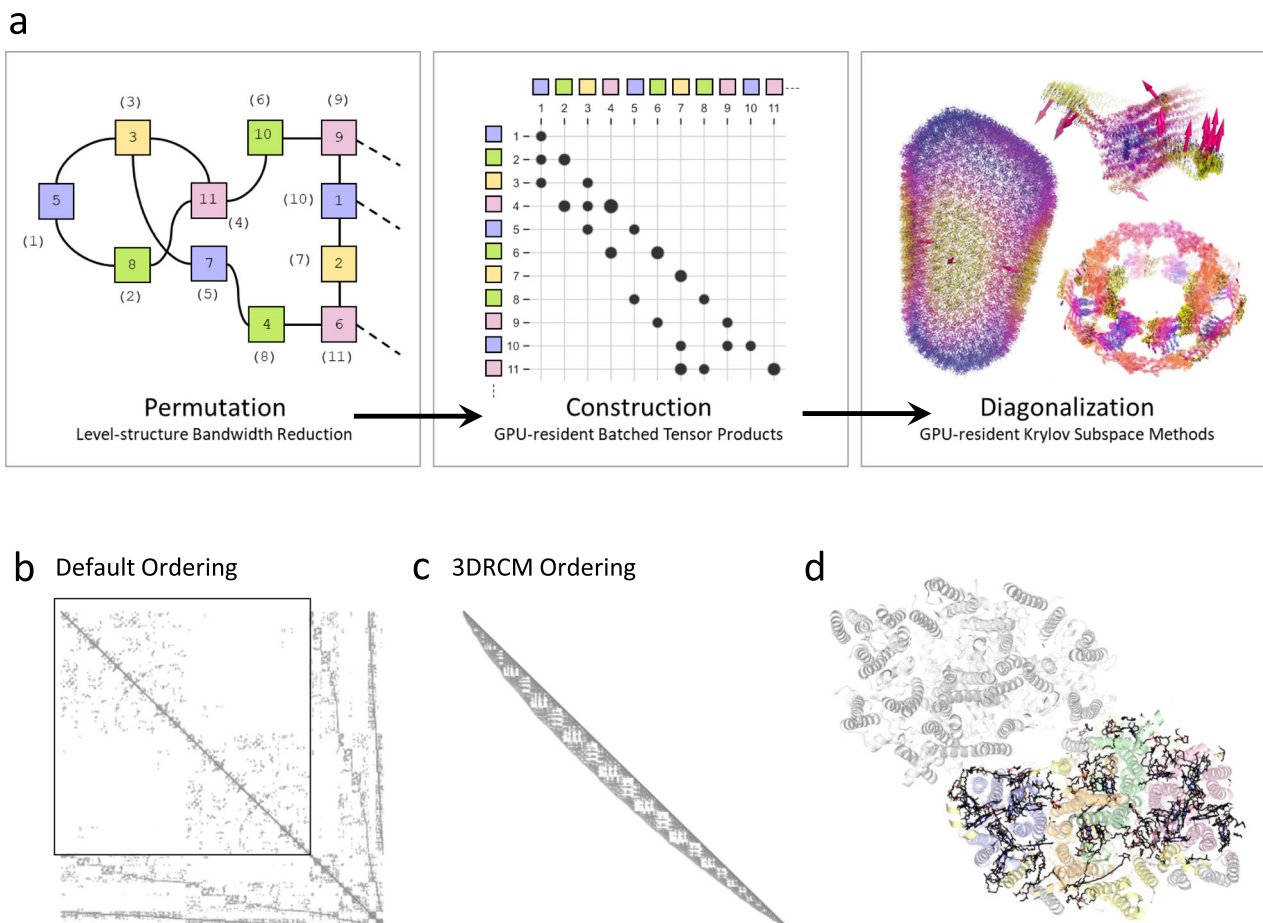
### Overview of the INCHING algorithm

The NMA is a study of a high-dimensional potential energy surface under harmonic approximation. In this work, the INCHING algorithm is applied to the elastic network model (ENM), a practical form of NMA introduced by Tirion<sup>15</sup>. As described in Methods, NMA is based on solving a standard eigenproblem,  $\mathbf{H}\mathbf{Q} = \mathbf{Q}\mathbf{\Omega}$ , concerning a Hessian matrix  $\mathbf{H} \in \mathbb{R}^{3N \times 3N}$ , where  $\mathbf{\Omega} \in \mathbb{R}^{3N \times 3N}$  and  $\mathbf{Q} \in \mathbb{R}^{3N \times 3N}$  are the eigenvalue matrix and the eigenvector matrix, respectively. For large collective conformational changes, only the eigenpairs with the smallest non-zero eigenvalues (the lowest-frequency modes) are of interest. However, challenges in the construction and diagonalization of the Hessian grow with the system size, and are notoriously resistant to parallelization. Consequently, current approaches to NMA are not readily adaptable for GPU computing, which imposes even stricter memory limits and algorithmic requirements to fully leverage on-chip parallelism. As illustrated in Fig. 1a, the INCHING algorithm was implemented in three synergistic stages—permutation, construction, and diagonalization—to resolve the challenges of calculating NMA on GPU. In the following sections, we will describe each stage in detail.

### Permutations to achieve bandwidth reduction

GPUs are hardware engineered to efficiently access contiguous memory locations, thus inherently favoring dense computations. However, in practice, the experience of memory access is dictated by the sparsity pattern of the data presented. Since the inception of NMA, it has been recognized that the Hessian is globally-sparse with prevailing zeros, primarily due to the long-range cut-offs applied to molecular mechanics potentials<sup>11</sup>. These globally-sparse Hessians are also locally-sparse, requiring very large bandwidths with non-uniform adjacencies (Fig. 1b–f, c.) due to sequentially distal segments of polymer(s) interacting in tertiary and/or quaternary structure(s)<sup>39,40</sup> (Supplementary Fig. 1), and the situation is further complicated by chemicals (e.g., cofactors, ions, water, lipids) with no inherent order being integrated into the polymer structure (Supplementary Fig. 2). All these factors contribute to degrade on-chip parallelism when dense row-sweeps were performed in batches.

The primary objective of the INCHING algorithm is to strategically permute columns and rows of the Hessian, prior to all its computations, such that a globally-sparse-yet-locally-dense computation can be achieved in subsequent stages. By permuting the atom ordering, we can always generate a graph isomorph of the original elastic network, while allowing retrieval of its unpermuted eigenpairs in linear complexity if a bijection mapping is provided. However, finding a permutation that exactly minimizes the bandwidth of a matrix is NP-



**Fig. 1 | Overview of the INCHING algorithm.** **a** Our Isomorphic NMA Calculations Harnessing 1 Necessary Gpu algorithm (INCHING) consists of three stages, namely, permutation, construction, and diagonalization. In the permutation stage, the 3-Dimensional Reverse Cuthill-McKee algorithm (3DRCM) takes atom coordinates as the input and produce a bandwidth-reduced indexing as output. The initial atom ordering is denoted within square-shaped nodes, while the 3DRCM-permuted atom ordering is represented using parentheses. In the construction stage, the lower triangle of a Hessian with globally-sparse-yet-locally-dense pattern is constructed incrementally on a Graphic Processing Unit (GPU) with tensor products and broadcasts. In the diagonalization stage, the eigenproblem is solved on GPU to produce depictable mode shapes. **b** Sparsity pattern of the Hessian with default

atom ordering for the PsbM-deletion mutant of photosystem II. The top left square highlighted corresponds to the proteins in the macromolecular complex, which also comprises cofactors, lipids and water molecules. **c** Sparsity pattern of the Hessian lower triangle with 3DRCM permuted atom ordering for the same complex (**d**) overall spatial distribution of the chemicals intercalated in multiple protein chains in default sequential order (green, blue, pink, orange, yellow). Only one of the dimeric halves of the structure is colored, the rest of the macromolecule is shown as a gray surface; note that the complex is asymmetric due to slight difference in lipids and cofactors intercalated. The Hessian matrix in default ordering has a much longer bandwidth everywhere than that with the 3DRCM ordering.

complete<sup>41</sup>. To approximate this reduction, we have designed a level-structure algorithm called 3DRCM, which extends on the well-established Reverse Cuthill-McKee (RCM)<sup>42,43</sup> algorithm, to handle 3-D coordinate data. Traditional RCM and variants<sup>44–46</sup> operate under the assumption that the input –a matrix– is realized and stored once-and-for-all. This framework is suitable for testing new offline analyses on the preserved matrix. However, a major problem overlooked is that it may not even be practical to realize the matrix efficiently, when its parallel computation spans a large bandwidth. In 3DRCM, a 3-D-tree data structure<sup>47</sup> is taken as input instead of a matrix to eliminate this circular reference. This modification enables efficient dynamic neighborhood lookups in the permutation process and avoids premature realization of the Hessian matrix. In Supplementary Fig. 3a, we showed that for macromolecules with over 100 thousand atoms, the mean bandwidth for the Hessian in default PDB sequence ordering can vary between 10% and 90% of  $N$ . However, by permuting the Hessian with 3DRCM before its construction, the mean bandwidth of the Hessian is consistently reduced and remains below 10% of  $N$ , hence confirming the local density of the matrix. In Method and Supplementary Fig. 3b, we also showed that the time complexity of 3DRCM is practically

dominated by a linear term proportional to the packing density of the macromolecule, thereby ensuring the cost-effectiveness of the algorithm.

### GPU-resident construction of the Hessian with batched tensor products

The proposed 3DRCM permutation enables spatial neighbors from each batch of atoms to be retrieved as locally-dense indexing slices. In our Method, we further showed that, by extending on the tensorial representation suggested by Koehl<sup>30</sup>, the Hessian can be constructed in batches through vectorized operations, such as tensor products and broadcasts, to leverage efficient on-chip parallelism in GPU. This replaces the need for explicit index loading and external storage of the distance matrix to enable fast computations. Besides, in our implementation, only the lower triangle of the Hessian was incrementally stored, and later accessed, in Compressed Sparse Row (CSR) format, this halves the memory requirement. In Supplementary Fig. 4, we showed that even for a system comprising 2.4 million atoms (PDBID: 3J3Q), the computation of the 3DRCM permutation ordering and the subsequent construction of the Hessian on GPU took only 5.2 min and

6.9 min in wall-clock time, respectively. Our method shows significant speedup over ProDy2.4<sup>48</sup>, an open-source implementation intended for coarse-grained protein representation in ANM tasks. For the largest case that ProDy2.4 handled (PDBID: 6NCL), containing 305 thousand atoms, ProDy2.4 took 7.8 h to construct the Hessian, whereas our INCHING approach delivered the matrix in just 93 seconds, showing a 302-fold acceleration.

### GPU-resident diagonalization of the Hessian matrix

In NMA calculations, a significant portion of computational resources is often dedicated to diagonalizing the Hessian, where the throughput is dependent on the performance of the eigensolver. One popular method in solving large symmetric eigenproblem is the Implicitly Restarted Lanczos Method (IRLM)<sup>49</sup>, which is available in the ARPACK package<sup>50</sup>, a backend incorporated into ProDy2.4. In ARPACK, the matrix-vector multiplications are accelerated by threaded BLAS sub-routines. However, despite the utilization of a high-end AMD processor with 64 threads, ProDy2.4 routines (ProDy-ARPACK-FullSparse) fail to converge within 48 h when systems exceed 305 thousand atoms (PDBID: 6HIV, 311 thousand atoms), indicating that we have reached the limit of hardware improvement. Note that ProDy2.4 with the default LAPACK<sup>51</sup> backend (ProDy-LAPACK-FullDense) working on a full dense Hessian is faster than all other programs when there are less than 500 atoms, but we were not able to proceed once exceeding 29 thousand atoms due to rapid rise in memory consumption. We also implemented a version of INCHING (INCHING-ARPACK-FullSparse) that uses our fast Hessian construction routine on GPU followed by a one-time device-to-host transfer to carry out diagonalization with ARPACK, but we were not able to proceed beyond 1.2 million atoms in 48 h, reflecting the bottleneck in calculation is fundamentally the diagonalization process.

GPU-accelerated approaches for solving large sparse symmetric eigenproblems continues to be a vibrant area of research<sup>52–57</sup>. To explore this next-generation technology, we implemented, optimized and compared several diagonalization methods on a single NVIDIA® GPU device, including the Implicitly Restarted Lanczos Method (IRLM)<sup>49,50</sup>, the Thick Restart Lanczos Method (TRLM)<sup>58</sup>, the Jacobi Davidson Method (JDM)<sup>31,32</sup> and some of their Chebyshev-filtered versions e.g. the Chebyshev-filtered Thick Restart Lanczos Method (CTRLM)<sup>34</sup> and Chebyshev-Davidson Method (CDM)<sup>35</sup>. In general, these methods all share the objective of computing eigenpairs within an interval (e.g., those with the smallest eigenvalues), but they differ in their approaches to update the Krylov subspace that approximates the eigenpairs. (See Methods for detail). The IRLM and TRLM are variants of the Hermitian Lanczos Method (HLM), but they differ in how they utilize the solutions of a much smaller tridiagonal eigenproblem to reinitialize the Krylov subspace at restarts. On the other hand, the JDM directly corrects the Krylov subspace by solving for an approximate fit to eliminate residuals and subsequently also works over solutions of a smaller projected eigenproblem. The CTRLM and CDM uses filters (low- or band-pass) constructed from Chebyshev polynomials to magnify wanted interval in the spectrum. The choice among these methods is not obvious, though their speed all depends on the cost of matrix-vector multiplication. In our INCHING protocols (INCHING-TRLM-HalfSparse, INCHING-IRLM-HalfSparse, INCHING-JDM-HalfSparse, INCHING-CTRLM-HalfSparse, INCHING-CDM-HalfSparse), these multiplications are accelerated by the SpMV CSR kernel in cuSPARSE<sup>59</sup>. The GPU computation was instructed through CuPy<sup>60</sup>, a python API to NVIDIA®'s CUDA<sup>61</sup>, cuBLAS<sup>62</sup>, cuSPARSE<sup>63</sup> and custom kernel programming. To accommodate the lower memory capacity of GPU, the Hessian is accessed as a sum decomposition of its lower triangle, halving the memory requirement. Alternatively, the calculations can also be done in the full sparse matrix with a further 40% speedup, when GPU memory is not exhausted. (See Supplementary Fig. 5) In Supplementary Fig. 6, we further showed that an explicit

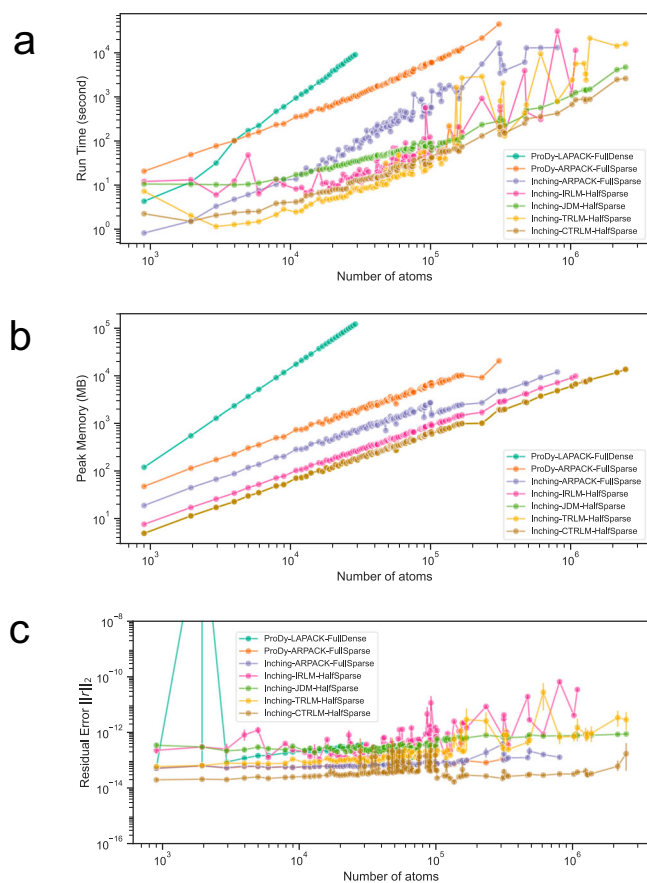
external deflation<sup>36</sup> of the first 6 rigid modes with zero eigenvalues (rotations and translations) can help to improve runtime of the remaining 58 non-zero eigenpairs in sub-megascale regime while sharing very similar memory footprint. Finally, by incorporating low- and band-pass Chebyshev filters<sup>34,57</sup>, we also lifted the memory limit regarding the number of modes to be resolved. (See INCHING-CTRLM-HalfSparse in Fig. 2 and Supplementary Figs. 7 and 8)

### Benchmarks

Benchmarks on correctness, memory and speed for calculating the first 64 modes are illustrated in Fig. 2. All reported runtimes are wall-clock time, and the tests were conducted on a computer with a 64-threads AMD EPYC™ 7513 processor and a single NVIDIA® A100 Tensor Core GPU, unless specifically noted. The accuracy is measured by the 2-norm of the residual error. Across 116 benchmark cases encompassing macromolecules ranging from around a thousand atoms to around 2.4 million atoms, including PDBID:3J3Q, the largest experimentally resolved atomic structure in PDB at <10Å range, our INCHING protocols were able to afford accuracy at  $10^{-12}$  level and achieved a peak memory consumption consistently lower than the ProDy2.4 routines. Importantly, with memory requirement halved by accessing only the lower triangle of the Hessian with a sum decomposition, the throughputs of our INCHING protocols, including construction and diagonalization of the Hessian, are still 146–251 times faster than ProDy-ARPACK-FullSparse and are 265–1290 times faster than ProDy-LAPACK-FullDense, depending on our choice of diagonalization algorithm. Remarkably, for a 305-thousand-atoms system (PDBID: 6NCL), all of our INCHING protocols took at most 5 min to converge with the fastest convergence being achieved by INCHING-TRLM-HalfSparse at 2.9 min. The same task took ProDy-ARPACK-FullSparse 12.4 h to converge. INCHING-TRLM-HalfSparse, INCHING-JDM-HalfSparse and their Chebyshev-filtered versions (INCHING-CTRLM-HalfSparse, INCHING-CDM-HalfSparse) also converged for the HIV-1 capsid structure (PDBID:3J3Q, 2.4 million atoms), with the fastest convergence achieved by INCHING-CTRLM-HalfSparse within 44 min. In Supplementary Fig. 7, we show that, at the same accuracy level as JDM ( $10^{-12}$ ), moderate speed-up in sub-megascale regime (mean at 1.21) can be achieved by CDM with an 80-degree polynomial, though the speed-up diminished to a slow-down in the megascale regime (mean at 0.91). This contrasts with applying an optimized low-pass filter to TRLM, where we showed that the Chebyshev-filtered TRLM (CTRLM) can steadily deliver a speed-up over TRLM in the megascale regime (mean at 7.44). In Supplementary Fig. 8, we further showed that linear or better scaling in runtime has been achieved in terms of the radii  $R_C = 6.8, 14\text{Å}$  to be considered, except for several cases at a lower radius  $R_C = 6\text{Å}$ , likely due to poorer conditioning of the matrix. In Fig. 3, by incorporating a band-pass Chebyshev filter developed in recent works<sup>34,57</sup> into our TRLM implementation (INCHING-CTRLM-HalfSparse), we also achieved linear time scaling in the number modes to be solved, while keeping memory usage constant. This ultimately allows 1000 modes of all the benchmarks to be solved, in batches of 64 eigenvectors on a standard A100 NVIDIA® GPU, without compromising memory or run time or atomic details. The method was also tested on a much less expensive RTX4090 NVIDIA® GPU in batches of 28 eigenvectors in Supplementary Fig. 10.

### Anisotropic vibrations of some megascale systems

To illustrate the usage of our software, we have applied INCHING to some of the largest atomic objects available, both natural and artificial. In Fig. 4, we illustrate the first non-zero mode of the mature HIV-1 capsid structure (PDBID:3J3Q), the largest experimentally resolved atomic object to date at <10Å resolution range, containing 2.4 million atoms and 1.6 billion non-zero entries in its Hessian. The cone-shaped capsid is an assembly of 186 hexamers and 12 pentamers, and it was suggested that these pentamers located at its hemispherical ends induce stable

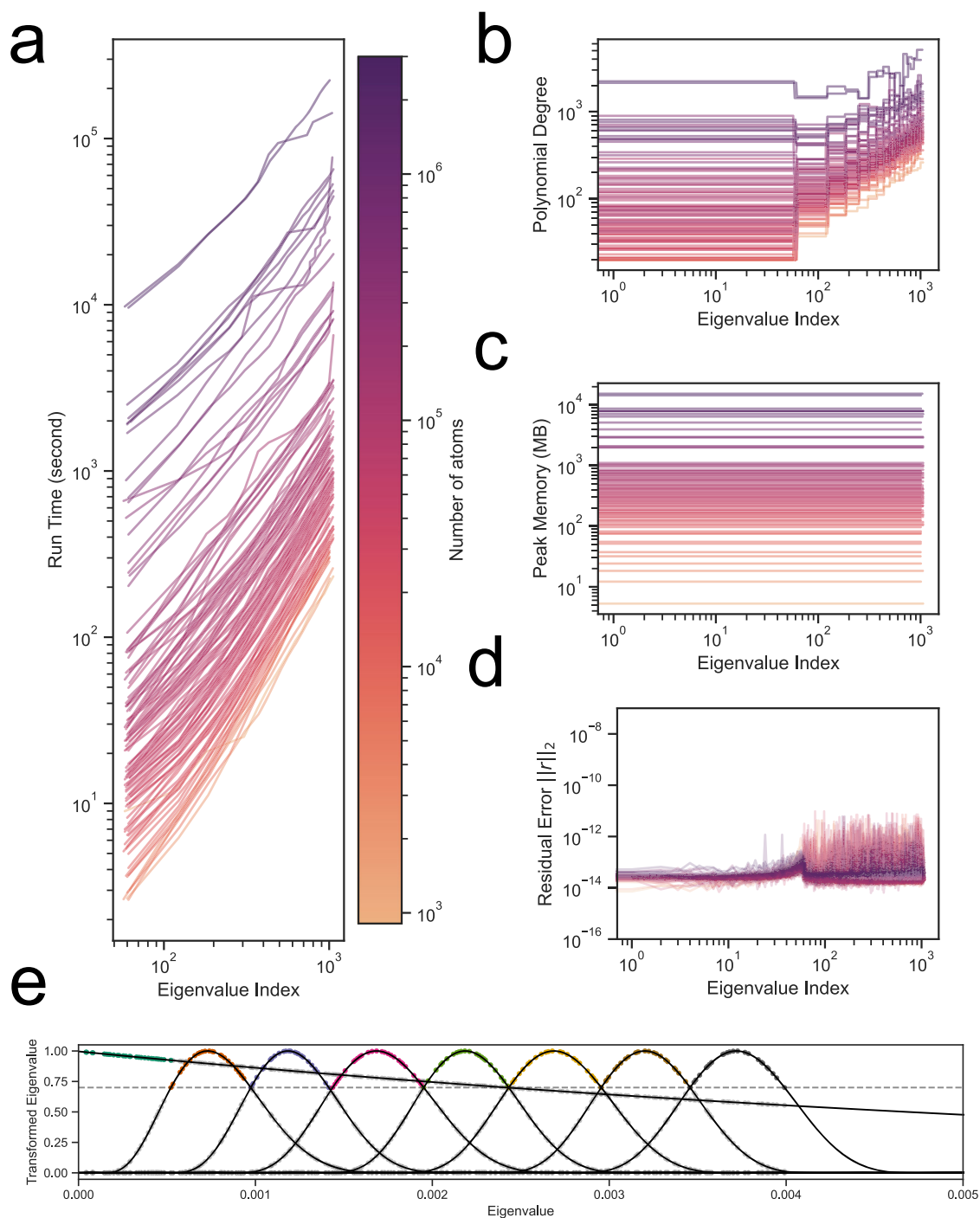


**Fig. 2 | Benchmark on throughput, memory consumption and correctness for each macromolecular structure in the benchmark dataset.** For all methods, the first 64 eigenpairs were calculated. Methods tested includes 2 ProDy methods (i.e., “ProDy-LAPACK-FullDense” and “ProDy-ARPACK-FullSparse”) and 5 INCHING methods (i.e., “INCHING-TRLM-HalfSparse”, “INCHING-CTRLM-HalfSparse”, “INCHING-JDM-HalfSparse”, “INCHING-IRLM-HalfSparse”, “INCHING-ARPACK-FullSparse”). Text in the name describe the eigensolver (e.g., ARPACK, LAPACK, JDM, TRLM, CTRLM, IRLM), the access of the Hessian matrix (e.g., “Full” means that the whole Hessian matrix is accessed and stored; “Half” means that only the lower triangle of the Hessian matrix is accessed and stored.), and the storage format of the Hessian matrix (e.g., “Dense” means a 2-D array is stored in double precision; “Sparse” means a Compressed Sparse Row (CSR) format is stored in double precision). “TRLM”, “CTRLM”, “JDM”, “IRLM” are our implementation of the Thick Restart Lanczos Method (TRLM), Chebyshev-filtered Thick Restart Lanczos Method (CTRLM), Jacobi-Davidson Method (JDM) and Implicitly Restarted Lanczos Method (IRLM). **a** Benchmark on overall run time including Hessian realization and subsequent diagonalization. The run time is wall-clock time to complete all the calculation. **b** Benchmark on peak memory consumption. Note that for “ProDy-LAPACK-FullDense”, we were not able to proceed once there are more than 29 thousand atoms due to memory overflow. Note that for “ProDy-ARPACK-FullSparse”, we were not able to proceed once there are more than 305 thousand atoms as it takes more than 48 h to converge. Also note that “INCHING-TRLM-HalfSparse”, “INCHING-CTRLM-HalfSparse” and “INCHING-JDM-HalfSparse” share very similar peak memory requirement. **c** Benchmark on residual error  $\|r\|_2$ , defined as the 2-norm of residual vector  $r$ . The error bar presented is the 95% confidence interval ( $n = 64$ ) calculated from the residual error of all the eigenvalues of a macromolecular structure in the benchmark dataset. Note that for “ProDy-LAPACK-FullDense”, the second macromolecular structure (PDBID: 1ASL) were not able to converge within  $10^{-10}$  likely due to severe fill-ins. All programs were stopped if run time exceeds 48 h. All INCHING programs were stopped if number of restarts exceeded 15000 rounds. Source data are provided as a Source Data file.

closure of the capsid by allowing sharp bite angles at the surface<sup>64</sup>. With our INCHING-JDM-HalfSparse protocol, we were able to resolve its first 64 normal modes at residual error  $10^{-12}$  within 1.3 h on an NVIDIA® A100 Tensor Core GPU, where the first non-zero mode corresponds to the oscillation of its hexameric surface roughly anchored at the pentamers supporting the theory of quasiequivalence<sup>65</sup>. In Fig. 5, the normal modes of a dilated human nuclear pore complex (NPC)<sup>66</sup> resolved at 50 Å were shown (EMDB: EMD14321, PDBID: 7R5J). The structure contains 4.8 million atoms and 3.5 billion non-zero entries in its Hessian. We were able to obtain the first 30 non-zero modes of this NPC structure within 12.5 h at accuracy of  $10^{-12}$ . In this calculation, an explicit external deflation<sup>36</sup> was applied to remove the first six rigid modes. We observe that while the first 14 non-zero modes mostly concern motions in the cytoplasmic and nuclear ring, the constriction of the inner ring can be observed at the fifteenth non-zero mode, reflecting the key functionally relevant conformational change. This calculation is at the edge of 80 GB of memory limit for the hardware used, though more powerful GPU systems on the market could handle even larger macromolecule superstructures. In Fig. 6, we applied the same methodology on the largest artificial DNA origami nano-structure, a DNA airplane made of 33kbp at its relaxed state<sup>67</sup> containing around 1.8 million atoms, including hydrogens, with a Hessian containing 1.8 billion non-zero entries. The structure has an apparent bilateral symmetry, where the joints leading to the wings are not chemically symmetric<sup>67</sup>. Interestingly, while its first non-zero mode presents a symmetric flop in its wings, the second non-zero mode demonstrates a complicated non-symmetric twisting motion involving its wings and stabilizers. In Supplementary Fig. 9, we also solved the first 64 modes of a 5-million pseudo-atoms coarse-grained representation of a Faustovirus capsid<sup>68</sup> (PDBID: 5J7V 26 million atoms resolved at 15.5 Å), illustrating the potential to incorporate coarse-graining strategies in handling systems that cannot fit into memory.

## Discussion

In recent years, there has been a significant shift in the focus of NMA methodologies, with an emphasis on accommodating larger atomic systems. The aim of an NMA is to gain insight into how some macroscopic motions involving communicating domains or subunits in macromolecular ensembles, can arise from microscopic interactions at the atomic level. In this respect, the mode shapes of NMA, which orchestrate collective motions concerning distal parts of the macromolecule, can often provide hints to understanding the biological structure at hand. Existing methods to analyze NMA can readily handle medium-sized structures with a few thousand atoms, but beyond this application of NMA is limited by computing resources and is not easily scalable by parallelization. This fundamentally restricts the molecular plasticity analysis in experimental techniques such as Cryo-ET and Cryo-EM, where macromolecular structures studied are often composed of millions of atoms and rather flexible. In these applications, fast and accurate megascale NMA would greatly facilitate modeling conformational dynamics in the refinement of the composite density maps<sup>19–21</sup> and the flexible fitting of atomic models<sup>22–24</sup> when EM volume or atomic-resolution template structures are available. Two major challenges in NMA scaling to megascale macromolecular complexes are the construction of the Hessian matrix and solution of the large-scale eigenproblem entailed. In this work, we have developed and implemented several advanced numerical recipes optimized for GPU computing, including a bandwidth-reduction algorithm and several alternative eigensolvers, to handle these challenges. This allows us to resolve normal modes, under the practical form of elastic network model, without coarse-graining the provided atomic structure for systems with several million atoms. In many cases, coarse-graining of biomolecular residues is expert-driven, and the appropriate level of granularity can be hard to determine when the chemical moiety is elongated (e.g., lipids) or is planar in shape (e.g., aromatics), especially

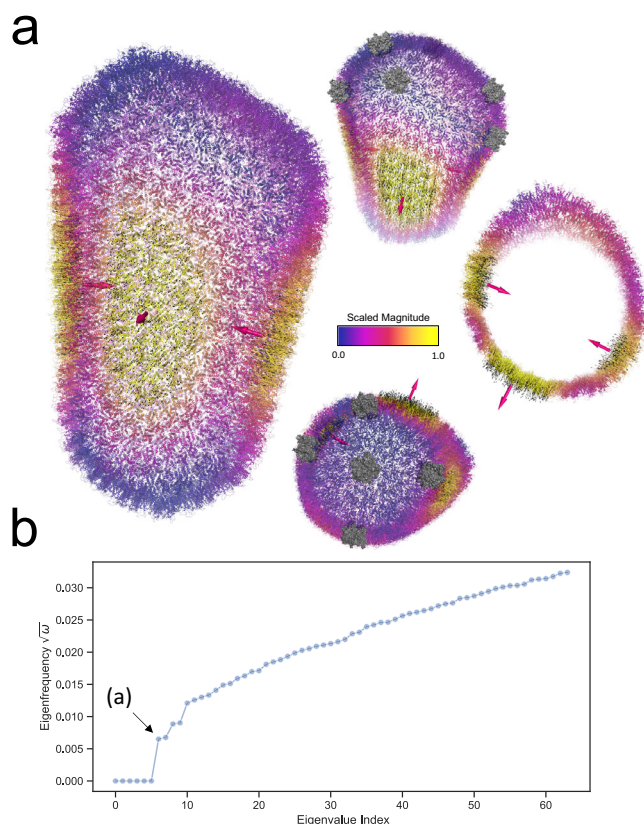


**Fig. 3 | Scaling in number of eigenmodes in Chebyshev-filtered thick restart Lanczos method.** Up to 1000 or slightly more eigenpairs were calculated using our INCHING-CTRLM-HalfSparse implementation of the Chebyshev-filtered Thick Restart Lanczos Method with radius of interaction  $R_C = 8\text{\AA}$ . The coloring of the lines, as shown in the color bar, refers to the number of atoms in the system. **a** Run time to complete until the indexed eigenvalue **(b)** optimized polynomial degrees involved in calculating the indexed eigenvalue. **c** Peak memory consumed when

calculating the spectrum slice containing the indexed eigenvalue. **d** Residual error in each of the eigenvalues. **e** Mapping of the lowest 224 eigenvalues in batches of 28 eigenvectors when low- and band-pass Chebyshev filters were applied to the HIV capsid system (PDBID: 3J3Q) without coarse-graining. Colors of the points indicate eigenvalues resolved in different batches. Source data are provided as a Source Data file.

in heterogeneous complexes comprised of protein, nucleic acid, lipid membranes, and sugars. The coarse-grained particle motions must be broadcasted into an all-atom construct if an atomic-scale understanding of the system is desired. In this respect, the implementation of an all-atom NMA, faithful to the formulation of Tirion<sup>15</sup>, is a direct response to this shortcoming as all the atoms are now included in the

system without neglecting each of their degrees of freedom (DOF). Several recent approaches to reduce DOF of atomic systems are interesting and indispensable to further scale up though. One of the first and the most popular approaches is the Anisotropic Network Model (ANM), which consider a subset of  $n_b$  atoms from the all-atom system, usually only the phosphorus P in nucleic acids and/or the



**Fig. 4 | Vibrations of the mature HIV-1 capsid structure (PDBID: 3J3Q).** **a** The first non-zero mode of the capsid. The black arrows are the displacement field of 1000 atoms randomly chosen from those in the top 90% quantile of vibration magnitude in the eigenvector. The arrow in magenta indicates an average direction for local clusters of the displacement field. Color scale in the cartoon from blue to red indicates increasing magnitude of vibration. Note that a logistic kernel is applied to the eigenvector to control the magnitude in visualization. (See Methods). The top-right inset is a clipped view of the capsid. The bottom-right inset indicates the locations of the pentamers (Gray opaque surface). **b** The first 64 eigenfrequency of the macromolecule, the eigenfrequency is the square root of eigenvalue  $\omega$ . The black arrow indicates the first non-zero eigenmode displayed in (a). Source data are provided as a Source Data file.

backbone carbon C $\alpha$  of the proteins<sup>15,25</sup>, thus effectively reducing  $N$  for 8-10 times. Agreement of ANM with experimental data<sup>18,25</sup> is satisfactory, though the displacement information of the discounted atoms is lost. This issue was alleviated by the rotation-translation block (RTB) method<sup>26</sup>, where a projection operator was developed to coarse-grain the atomic system into a system of  $n_b$  rigid blocks, each with its own translational and rotational DOF, hence effectively reducing the  $3N \times 3N$  Hessian matrix to a  $6n_b \times 6n_b$  RTB matrix. Further refinement along this line is the Block Normal Mode (BNM)<sup>27</sup> method, which surrogates the initial realization of the peak-memory-consuming  $3N \times 3N$  Hessian matrix; this was done by exploiting the block structure of the projection operator and by constructing only part of the sparse full atomic Hessian on-the-fly. Very recently, it was also shown that an extrapolation of RTB modes can effectively predict nonlinear motions at large amplitudes<sup>69</sup>. Depending on the granularity of the system, more than tenfold reduction in the size of the Hessian matrix can be achieved. However, the cost of diagonalization can still be very prohibitive.

In this respect, parallel computing on GPU, as we implemented in this work, is an effective way to amortize the cost at a fundamental level, without compromising accuracy or atomistic detail. We also expect that advances in GPU hardware<sup>70</sup> in memory and processing

rate, integrated with techniques in solving large-scale eigenproblems, in particular recent trends in spectrum slicing<sup>34,52-57,71</sup>, would eventually allow even faster NMA calculations on systems exceeding 20 million atoms, such as structure of a Faustovirus<sup>68</sup> available in low resolution (PDBID: 5J7V, 15.5 Å). Nonetheless, given the applicability of NMA in biological systems, we believe our framework will be useful in the exploration of megascale structural dynamics. The growing complexity of the megacomplexes resolved by Cryo-EM and/or cryo-ET now calls for methods that can rapidly capture conformational and functional plasticity, potentially as an intrinsic part of the refinement pipeline. Such methods will play a crucial role in advancing our understanding of these macromolecular machines.

## Methods

### Normal mode analysis of an elastic network model

The Normal Mode Analysis (NMA) is a classic approach to derive motions from static structures at local minima of a potential energy surface. The potential energy  $V$  is dependent on the conformation  $\mathbf{X} \in \mathbb{R}^{3N}$  for a structure with  $N$  atoms at time  $t$ . Without loss of generality, we begin with a particular conformation  $\mathbf{X}^{(0)} \in \mathbb{R}^{3N}$ . For small displacements, we may then tolerate a second-order Taylor expansion,

$$V(\mathbf{X}|\mathbf{X}^{(0)}) = V(\mathbf{X}^{(0)}|\mathbf{X}^{(0)}) + \nabla V(\mathbf{X}^{(0)}|\mathbf{X}^{(0)}) (\mathbf{X} - \mathbf{X}^{(0)}) + \frac{1}{2} (\mathbf{X} - \mathbf{X}^{(0)})^T \nabla^2 V(\mathbf{X}^{(0)}|\mathbf{X}^{(0)}) (\mathbf{X} - \mathbf{X}^{(0)}) + \dots \quad (1)$$

By choosing the energy level  $V(\mathbf{X}^{(0)}|\mathbf{X}^{(0)}) = 0$  and assuming local minimum  $\nabla V(\mathbf{X}^{(0)}|\mathbf{X}^{(0)}) = 0$ , we are left with the quadratic form,

$$V_{NMA} = \frac{1}{2} \Delta \mathbf{X}^T \mathbf{H} \Delta \mathbf{X} \quad (2)$$

where  $\mathbf{H} = : \nabla^2 V(\mathbf{X}^{(0)}|\mathbf{X}^{(0)})^T \in \mathbb{R}^{3N \times 3N}$  is the Hessian matrix and  $\Delta \mathbf{X} = : (\mathbf{X} - \mathbf{X}^{(0)}) \in \mathbb{R}^{3N}$  is the deviation from the minimum. We will defer the layout of  $\mathbf{H}$  to the next section when the form of potential energy  $V$  is defined and proceed to discuss the outcome of an NMA. Substituting  $V_{NMA}$  into the equations of motion gives a partial differential system

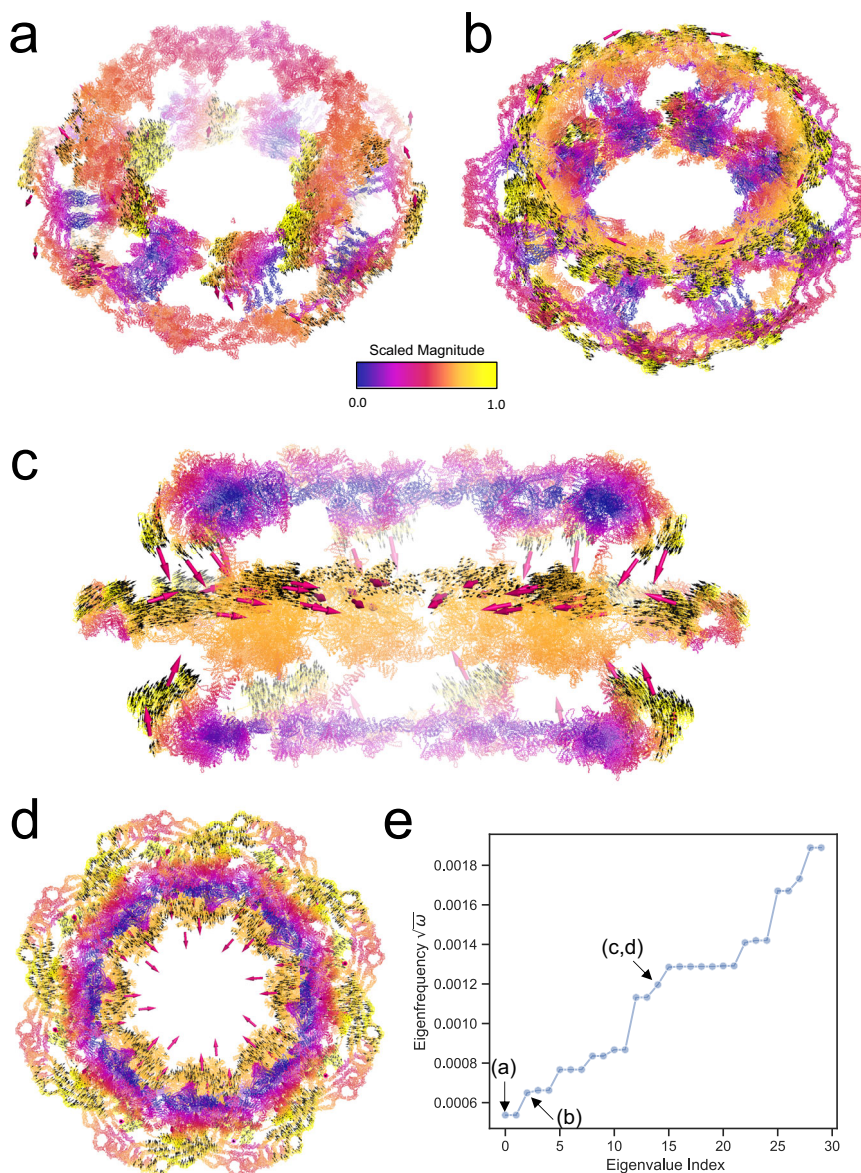
$$\mathbf{M} \frac{d^2 \Delta \mathbf{X}}{dt^2} = - \mathbf{H} \Delta \mathbf{X} \quad (3)$$

And, the general solution  $\Delta \mathbf{X} = \mathbf{Q} e^{-i\omega t}$  gives a generalized eigenproblem

$$\mathbf{H} \mathbf{Q} = \mathbf{M} \mathbf{Q} \boldsymbol{\omega} \quad (4)$$

Where the eigenvector matrix  $\mathbf{Q} \in \mathbb{R}^{3N \times 3N}$  can be viewed as an eigentensor  $\mathbf{Q} \in \mathbb{R}^{3N \times N \times 3}$ . The physical meaning of this eigentensor  $\mathbf{Q}$  is that each of its atom slice (the second index) gives a displacement vector, hence  $3N$  linearized mode shapes, i.e., the collective motions, of the static structure  $\mathbf{X}^{(0)}$  can be obtained. In NMA, the first 6 modes will always have zero eigenvalues corresponding to rigid translational and rotational displacements. The eigentensor  $\mathbf{Q}$  is also the best displacement under orthonormality constraint  $\min_{\mathbf{Q}} V_{NMA} s.t. \mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ . Eigenfrequency  $\sqrt{\omega_i}$  is the square root of the eigenvalue  $\omega_i$ .

A variety of  $V$  exists. In the standard NMA, potential energies from a detailed all-atom forcefield can be used. However, this approach suffers from the tedious, and sometimes virtually unachievable, requirement of energy minimization<sup>10</sup>. This requirement was surrogated in the work of Tirion<sup>15</sup>, where the potential energy of an atomic system was considered as an elastic network model (ENM) connecting



**Fig. 5 | Vibrations of the dilated human Nuclear Pore Complex (NPC).** Overview of motions in several representative non-zero modes. **a** The first non-zero normal mode (indexed as Mode 0) **(b)** the third non-zero normal mode (indexed as Mode 2) **(c)** clipped view of the fifteenth non-zero normal mode (indexed as Mode 14) **(d)** bird-eye view of the fifteenth non-zero normal mode (indexed as Mode 14) The arrows in magenta indicate an average directions for the local clusters of the displacement field, while the tiny black arrows show more detailed displacement fields of 10,000 atoms

randomly chosen from those in the top 80% quantile of vibration magnitude in the eigenvector. Color scale in the cartoon from blue to red indicates increasing magnitude of vibration. Note that a logistic kernel is applied to the eigenvector to control the magnitude in visualization. (See “Methods”). **e** The first 30 non-zero eigenfrequencies of the macromolecule, the eigenfrequency is the square root of eigenvalue  $\omega$ . The black arrows indicate the first non-zero eigenmode displayed in **(a)**, **(b)**, **(c)** and **(d)**. Source data are provided as a Source Data file.

atoms  $i$  and  $j$

$$V(\mathbf{X}) = : \sum_{ij} \frac{1}{2} k_{ij} (r_{ij} - r_{ij}^{(0)})^2 \quad (5)$$

Where  $r_{ij}$  is the variable Euclidean distance between atom pairs;  $r_{ij}^{(0)}$  is the equilibrium Euclidean distance from  $X^{(0)}$ . Importantly, a Heaviside function  $k_{ij}$ , parameterized on  $r_{ij}^{(0)}$ , was used to eliminate the long-range interactions beyond the threshold  $R_C$ . This can also be considered as applying an adjacency matrix  $k_{ij} \in K$  to the network.

$$k_{ij}(r_{ij}^{(0)}) = \begin{cases} 1, & r_{ij}^{(0)} \leq R_C \\ 0, & r_{ij}^{(0)} > R_C \end{cases} \quad (6)$$

As in the formulation of Tirion, a unified mass  $\mathbf{M} = \mathbf{I}$  were taken, hence reducing the generalized eigenproblem to a standard eigenproblem.

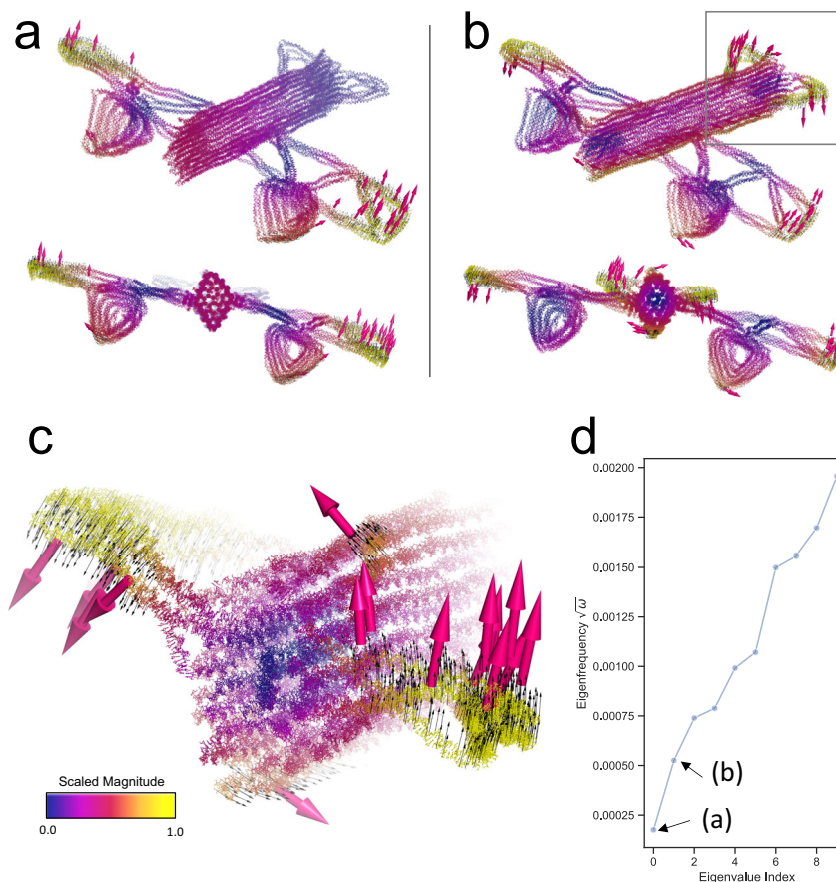
$$\mathbf{H}\mathbf{Q} = \mathbf{Q}\mathbf{\Omega} \quad (7)$$

#### Permuting the Hessian to produce graph isomorphs

In our later exposition, we will permute the Hessian matrix such that the resultant matrix is locally dense. Observe that a similarity transform of  $\mathbf{H}$  with permutation matrix  $\mathbf{P} \in \mathbb{R}^{3N \times 3N}$  preserves the eigenvalues  $\mathbf{\Omega}$  as

$$(\mathbf{P}\mathbf{H}\mathbf{P}^T)(\mathbf{P}\mathbf{Q}) = (\mathbf{P}\mathbf{Q})\mathbf{\Omega} \quad (8)$$





**Fig. 6 | Vibrations of a DNA origami airplane.** Overview of motions in the first two non-zero modes of the airplane. **a** The first non-zero eigenmode. **b** The second non-zero eigenmode; the square border indicates the stabilizer, which is also shown in (c). Color scale in the cartoon from blue to red indicates increasing magnitude of vibration. Note that a logistic kernel is applied to the eigenvector to control the magnitude in visualization. (See “Methods”). The bottom inset show the airplane along its apparent bilateral symmetry axis. **c** A zoom into the motion of the stabilizer in the second non-zero mode. The black

arrows are the displacement field of 1000 atoms randomly chosen from those in the top 90% quantile of vibration magnitude in the eigenvector. The arrow in magenta indicates an average direction for local clusters of the displacement field, which also correspond to (a). **d** The first 10 non-zero eigenfrequencies of the macromolecule, the eigenfrequency is the square root of eigenvalue  $\omega$ . Note that explicit external deflation was applied to remove the rigid modes. Source data are provided as a Source Data file.

and the eigenvector  $\mathbf{P}^T \mathbf{PQ}$  of the original problem can be recovered by simply applying inverse of the unitary permutation matrix, which is its transpose. Importantly, if we restrict the permutation to atom ordering, i.e., permuting every 3 consecutive indices as an immutable group in the  $3N$  indices, then the resultant Hessian is a representation of a graph isomorph of the original elastic network, where the bijection is provided by correspondence between the original atom ordering and the permuted atom ordering. This property means that the similarity transform can be computed by simply initiating with a permuted coordinate and similarly  $\mathbf{Q}$  can be recovered by permuting the atom slice of  $\mathbf{PQ}$  using the same bijection with linear cost.

### Constructing the Hessian using batched tensor products and broadcasts

We begin with the following un-vectorized notation, where the superscript (0) is dropped for readability.

$$\mathbf{X}^{(0)} = : [\mathbf{X}_1^{(0)}, \dots, \mathbf{X}_n^{(0)}]; \mathbf{Vec}(\mathbf{X}^{(0)}) = : [x_1 y_1 z_1, \dots, x_n y_n z_n] \quad (9)$$

Applying differentiations, the off-diagonal force constant block is specified by

$$\begin{aligned} \mathbf{H}_{ij} &= -\frac{k_{ij}}{(r_{ij}^{(0)})^2} (\mathbf{X}_i^{(0)} - \mathbf{X}_j^{(0)}) (\mathbf{X}_i^{(0)} - \mathbf{X}_j^{(0)})^T \\ &= -\frac{k_{ij}}{(r_{ij}^{(0)})^2} \begin{bmatrix} (x_i - x_j)^2 & (x_i - x_j)(y_i - y_j) & (x_i - x_j)(z_i - z_j) \\ (y_i - y_j)(x_i - x_j) & (y_i - y_j)^2 & (y_i - y_j)(z_i - z_j) \\ (z_i - z_j)(x_i - x_j) & (z_i - z_j)(y_i - y_j) & (z_i - z_j)^2 \end{bmatrix} \\ \mathbf{H}_{ii} &= -\sum_j \mathbf{H}_{ij} \end{aligned} \quad (10)$$

For any matrix  $\mathbf{A}$  with  $n$  rows, we can define for a row in it,

$$\beta_i = : |i - \min_j (|j| a_{ij} \neq 0)| \quad (11)$$

Then the bandwidth of the matrix is  $\max_i(\beta_i)$  and the profile of the matrix is  $\sum_i \beta_i$ , the mean bandwidth is  $\frac{\sum_i \beta_i}{n}$ . Clearly, due to the cutoff  $R_C$ , the layout of the Hessian tensor  $\mathbf{H} \in \mathbb{R}^{N \times 3 \times N \times 3}$ , viewing from the atom slices, will have the same bandwidth as the adjacency matrix  $\mathbf{K}$  of the macromolecule. The tensor is globally sparse in most cases. This observation also applies to standard NMA as long range cut-offs were used for van der Waals, electrostatic, and hydrogen-bonding interactions<sup>12</sup>. In general, it is not advisable to construct a dense Hessian matrix that contains a large number of zero entries due to the quadratic storage consumption and the resultant increase in fill-ins. To construct a sparse Hessian matrix, we find the tensorial layout by Koehl<sup>30</sup> a good starting notation. In which, he defined the following  $N \times 3$  matrix, which can also be vectorized as  $3N$  elements.

$$\mathbf{U}_{ij}(\mathbf{X}^{(0)}) = : \left( 0, \dots, 0, \frac{\mathbf{X}_i^{(0)} - \mathbf{X}_j^{(0)}}{r_{ij}^{(0)}}, 0, \dots, 0, \frac{\mathbf{X}_j^{(0)} - \mathbf{X}_i^{(0)}}{r_{ij}^{(0)}}, 0, \dots, 0 \right) \quad (12)$$

The distance  $r_{ij}^{(0)}$  can be precomputed and filtered by  $k_{ij}(r_{ij}^{(0)})$  for pairs within  $R_C$ . Then, the off-diagonal block of the entire Hessian tensor is given as a sum of Kronecker products.

$$\mathbf{H}_{ij} = \sum_i \sum_j k_{ij} \mathbf{U}_{ij} \otimes \mathbf{U}_{ij}^T \quad (13)$$

More importantly, when Hessian-vector multiplication  $\mathbf{H}\mathbf{w}$  is accessed in eigensolvers, the following can be done separately for each  $(i,j)$  pairs.

$$\mathbf{H}\mathbf{w} = \sum_{ij} k_{ij} (\mathbf{U}_{ij} \otimes \mathbf{U}_{ij}^T) \mathbf{w} = \sum_{ij} k_{ij} (\mathbf{U}_{ij} \mathbf{w}) \mathbf{U}_{ij} \quad (14)$$

Given explicit indexing for locations of the non-zeros  $k_{ij}$ , the above Hessian-vector product can be effectively parallelized on a hyperthreaded computer.

In our implementation, the tensorial representation above was refined to fully exploit on-chip parallelism on GPU and to avoid explicit index loading which tends to thread divergence on GPU. To facilitate discussion, we define a difference matrix  $\mathbf{G}_{ij} \in \mathbb{R}^{N \times 3}$

$$\mathbf{G}_j(\mathbf{X}^{(0)}) = : \left( \mathbf{0}, \dots, \mathbf{0}, \mathbf{X}_{b_3}^{(0)} - \mathbf{X}_j^{(0)}, \dots, \mathbf{X}_i^{(0)} - \mathbf{X}_j^{(0)}, \dots, \mathbf{X}_{b_4}^{(0)} - \mathbf{X}_j^{(0)}, \mathbf{0}, \dots, \mathbf{0} \right) \quad (15)$$

Supposed, we know  $|b_3 - b_4|$  is the bandwidth of an off-diagonal batch of  $H$ , then, for two contiguous atom indices  $(b_1, b_2]$  and  $(b_3, b_4]$ , the batch  $\mathbf{G} \in \mathbb{R}^{3|b_1 - b_2| \times 3|b_3 - b_4|}$  can be given by the four index corners  $(b_1, b_2, b_3, b_4)$ .

$$\mathbf{G}_{(b_1, b_2, b_3, b_4)} = \sum_{i \in (b_1, b_2]} \sum_{j \in (b_3, b_4]} \mathbf{G}_i \otimes \mathbf{G}_j^T \quad (16)$$

In practice,  $(b_3, b_4]$  could include intervals that lack neighboring atoms for the range  $(b_1, b_2]$ , but the indexing slice  $(b_3, b_4]$  can be further refined into multiple contiguous indexing slices to eliminate unnecessary calculations. Nevertheless, filtering for non-zero elements in the batch can be done by observing Eq. (10) in that the squared interatomic distance is given by the trace of each block  $(r_{ij}^{(0)})^2 = \text{tr}(\mathbf{G}_{ij})$ , then we can consider the following block-wise operations, readily

parallelizable on GPU by in-place tensorial broadcasts.

$$\begin{aligned} Y_{ij} &= \text{tr}(\mathbf{G}_{ij}) \\ K_{ij} &= \begin{cases} 1, & Y_{ij} \leq R_C^2 \\ 0, & Y_{ij} > R_C^2 \end{cases} \\ Y_{ij} &= \begin{cases} Y_{ij}^{-1}, & \forall i \neq j \\ 0, & \forall i = j \end{cases} \end{aligned} \quad (17)$$

$$\mathbf{H}_{(b_1, b_2, b_3, b_4)} = \left( \left( \mathbf{K}_{(b_1, b_2, b_3, b_4)} \odot \mathbf{Y}_{(b_1, b_2, b_3, b_4)} \right) \otimes \mathbf{1}_{3 \times 3} \right) \odot \mathbf{G}$$

The  $\odot$  is the Hadamard product operator;  $\mathbf{1}_{3 \times 3}$  is the all-ones matrix presenting the broadcast. This is followed by row-sum in the diagonal  $\mathbf{H}_{ii} = -\sum_{j \in (b_3, b_4]} \mathbf{H}_{ij}$  accordingly. Only the lower triangle is incrementally stored in the Column-Sparse-Row (CSR) format for subsequent calculation, hence consuming  $O(9p)$  memory for the data content;  $p$  is the number of non-zero pairwise interaction for  $i \geq j$ . The dense batch  $\mathbf{G}_{(b_1, b_2, b_3, b_4)}$  presented above will require  $O(9|b_1 - b_2| |b_3 - b_4|)$  erasable memory, but we will show later that the local bandwidth  $|b_3 - b_4|$  can be reduced.

### Bandwidth and layout of the Hessian

Obviously, the cost of the parallel computation will depend on the local bandwidth  $|b_3 - b_4|$  of the batch  $(b_1, b_2]$ . The Hessian obtained with the default ordering of atoms can result in very large bandwidths with non-uniform patterns meaning interactions among sequentially distal parts within the tertiary structure or individual chains within a quaternary structure. A hypothetical minimal example is a macromolecule with 3 peptide chains  $\alpha$ ,  $\beta$  and  $\gamma$ , where only  $\alpha$ - $\gamma$  and  $\beta$ - $\gamma$  interactions were found, but not  $\alpha$ - $\beta$ ; in this case, the order  $\alpha$ - $\beta$ - $\gamma$  will have a much larger bandwidth than  $\alpha$ - $\gamma$ - $\beta$ . Many reasons can attribute to this observation. An example found in the benchmark set is a chimeric Sesbania mosaic virus coat protein<sup>72</sup> (PDBID: 4YSZ) composed of 12 sets of pentamers arranged on the vertices of an icosahedron, where each protein chain in the pentamer interacts with 7 other chains within 8 Å. (See Supplementary Fig. 1) A systematic order to build up this icosahedron, as done by the authors, is to place a pentamer on each of the four corners of the three orthogonal golden rectangles, done one rectangle after another starting at the shorter sides of the golden rectangles. (See Supplementary Fig. 1a) In this case, pentamers on the shorter sides interact within 8 Å among each other, but pentamers on the long sides of the same rectangle do not. This creates exactly the situation where the aforementioned  $\alpha$ - $\beta$ - $\gamma$  order arises. Indeed, the vertices on a icosahedron can never be clustered satisfactorily and there are specialized algorithm to calculate their vibrational dynamics<sup>73</sup>. Besides, due to the inconsistent presence of water molecules, the exact symmetry is destroyed in the crystal structure, which is commonly encountered. Nevertheless, the degeneracies in modes due to the apparent symmetry of this macromolecule is captured in our program. (See Supplementary Fig. 1d) There are also cases where no meaningful sequential ordering is possible when intact chemical structures (e.g. cofactors, ions, water, lipids) are intercalated between sequentially distal parts of the polymer. A practical example in the benchmark set is a PsbM-deletion mutant of photosystem II<sup>74</sup> (PDBID: 5H2F), where elongated lipids and cofactors are integral part of the macromolecule complex surrounded by several protein chains. In this case, it is not obvious as for how to rearrange the atoms to reduce the bandwidth, but we can show that the Hessian can always be permuted to reduce bandwidth by the algorithm described in the next section. (See Supplementary Fig. 2)

### Bandwidth reduction of the Hessian matrix

Finding a permutation to minimize the bandwidth is NP-complete<sup>41</sup>, but a reduction of bandwidth and its overall profile can be approximated by algorithms such as the Reverse Cuthill McKee (RCM)

algorithm. The RCM is a greedy approximation with a breadth-first level structure. The input of a standard RCM algorithm is an adjacency graph, and its outputs is a permuted node ordering. In RCM, starting from a peripheral node with the lowest degree of connection, adjacent unvisited nodes are collected as a level structure and re-ordered by their degree of connection. Provided that an adjacency matrix  $\mathbf{K}$  is precomputed, the time complexity of the RCM is bounded by  $O(4|E| + 2cm|E| + N)$ , where  $m = \max_j \sum_i K_{ij}$  is the maximum degree among all nodes and  $2|E|$  is the number of edges in  $\mathbf{K}$ . A detail proof is in reference<sup>43</sup>. The term  $4|E|$  is referring to  $2|E|$  operations to determine the degree of each node plus another  $2|E|$  operations to sweep through the adjacency matrix to locate the neighbors. The term  $2cm|E|$  refers to the insertion sorting of the degree in retrieved neighbors. The last term  $N$  is due to reversal of order. The input of RCM is an adjacency matrix  $\mathbf{K}$ . For 3-D coordinates, this may be obtained by a cell-linked list data structure in  $O(27Nn_c)$  where  $n_c$  is the average number of particles per cell in volume  $R_c^3$  or simply an all-to-all calculation as done in ProDy2.4 for small-sized systems. Therefore, a standard RCM algorithm operating on 3-D coordinate data will require a total time complexity of  $O(27Nn_c + 4|E| + 2cm|E| + N)$  and ideally a  $O(mN)$  memory to operate due to the adjacency matrix. To surrogate the storage of the adjacency matrix, we supplemented the RCM with a k-D tree data structure<sup>47</sup>, which takes the coordinate data  $\mathbf{X}^{(0)}$  directly as input. The algorithm is labeled as 3DRCM to avoid confusion with the standard RCM. A pseudocode of our 3DRCM implementation is provided in the Supplementary Information as Algorithm 1. A 3-D tree is a balanced space-partitioning binary tree, which can be constructed and stored by finding and storing hyperplanes that split the number of points in halves. As such, the construction of a 3-D tree takes  $O(N \log N)$  time and  $O(3N)$  space. On a 3-D tree, each range search for adjacent neighbors within threshold distance takes  $O(3N^{2/3})$ <sup>75</sup>. At this stage, the total time complexity of 3DRCM appears to be  $O(2 \cdot 3N^{2/3} + 2cm|E| + N)$ , where the first term  $2 \cdot 3N^{2/3}$  is due to the collection of degree in each node and the search for neighbor, but we will show that in both 3DRCM and RCM the term  $O(2cm|E|)$  is dominating. A tighter bound on  $O(2cm|E|)$  in the context of NMA can be obtained as follows. For a stable macromolecule without atomic clashes, the shortest interatomic distance is due to covalent bonds at around 1Å. The Kepler's conjecture proven by Hales<sup>76</sup> states that the maximum volume ratio occupied by equidistant sphere packing is  $\frac{\pi}{3\sqrt{2}} < 0.7405$ . Hence, the maximum number of atoms allowable without clashes within radius  $R_c$  is  $\rho = R_c^3 \frac{\pi}{3\sqrt{2}} \geq m$ . This packing density  $\rho$  presents an upper bound to the number  $m$  of non-zero entries in a row on the adjacency matrix. This bounds the total number of edges as

$$\rho N \geq mN > 2|E| \quad (18)$$

For an 8 Å radius, there can only be less than 380 atoms in absence of clash. Therefore, the time complexity of 3DRCM and RCM in the context of NMA are bounded by  $O(6N^{2/3} + cp^2N + N)$  and  $O(27N\rho + 2\rho N + cp^2N + N)$  respectively. Taking  $c=1$ , the 3DRCM is dominated by the pseudolinear term  $O(\rho^2N)$  unless  $N$  is beyond 3.7 million atoms. From experience, for a system containing 2.4 million atoms, computing the permutation ordering with 3DRCM and the calculation of the Hessian took only 5.2 min and 6.9 min in wall-clock time respectively. See Supplementary Fig. 4. Hence, for most practical purpose, the trade-offs in time complexity in 3DRCM compared to RCM is negligible.

### Low-complexity Krylov-subspace eigensolvers

In this work we have implemented several iterative methods<sup>31,32,49,58</sup> to solve large eigenvalue problems on the GPU. Specifically, in NMA, we are mostly interested in the smallest non-zero eigenpairs. The Hessian

matrix of concern is a sparse positive semidefinite real symmetric matrix without weak diagonal dominance. To facilitate communication, we adopt more generic notations here. The matrix of concern is notated as  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ; the exact and the approximate eigenpairs are  $(\mathbf{A} \in \mathbb{R}^{m \times m}, \mathbf{U} \in \mathbb{R}^{n \times m})$  and  $(\tilde{\mathbf{A}} \in \mathbb{R}^{m \times m}, \tilde{\mathbf{U}} \in \mathbb{R}^{n \times m})$  respectively;  $n$  and  $m$  denotes the dimension of the basis set and the number of basis respectively. The standard eigenproblem is thus  $\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{A}$  and we assumed orthonormality among the exact eigenvectors. The approximate  $\tilde{\lambda} \in \tilde{\mathbf{A}}$  can be obtained from the Rayleigh quotient  $\tilde{\lambda} = \tilde{\mathbf{u}}^T \mathbf{A} \tilde{\mathbf{u}}$  and the residual  $\mathbf{r} = \mathbf{A}\tilde{\mathbf{u}} - \tilde{\lambda}\tilde{\mathbf{u}}$  can always be evaluated by its 2-norm, the residual error  $\|\mathbf{r}\|_2$ .

In the following paragraphs, we will progressively introduce two branches of iterative methods implemented, namely the Hermitian Lanczos Methods followed by the Jacobi-Davidson Method. Only rationales and key equations were presented. A common central idea is to find  $\tilde{\mathbf{u}}$  by refining an initial guess  $\mathbf{v}_0 \in \mathbb{R}^n$  to build up an orthonormal basis  $\mathbf{V} \in \mathbb{R}^{n \times m}$  in the Krylov subspace

$$\mathcal{K}_k(\mathbf{A}, \mathbf{v}_0) = : \text{span}\{\mathbf{v}_0, \mathbf{A}^1\mathbf{v}_0, \mathbf{A}^2\mathbf{v}_0, \dots, \mathbf{A}^{k-1}\mathbf{v}_0\} \quad (19)$$

The residual is minimized when the Galerkin condition

$$(\mathbf{v}_e)_j^T (\mathbf{A}\tilde{\mathbf{u}} - \tilde{\lambda}\tilde{\mathbf{u}}) = 0 \forall j \in (0, \dots, m-1) \quad (20)$$

is satisfied, where  $\tilde{\mathbf{u}} = \mathbf{V}\mathbf{y}$ . The  $\mathbf{y} \in \mathbb{R}^{m \times 1}$  is an unknown component to combine  $\mathbf{V}$ , but if  $\mathbf{V} \in \mathcal{K}_k$  has its orthonormality maintained satisfactorily, for example by incorporating the Modified Gram Schmidt algorithm (MGS), then  $\mathbf{y}$  is the eigenvector of a symmetric tridiagonal eigenproblem<sup>77</sup> of a much smaller size  $m \times m$

$$(\mathbf{V}^T \mathbf{A} \mathbf{V}) \mathbf{y} = \tilde{\lambda} \mathbf{y}, \quad (21)$$

and  $\tilde{\mathbf{u}} \in \mathbb{R}^n$  can be recovered by definition. This orthogonal projection technique, an example of Rayleigh-Ritz process, is common to all three methods implemented.

**Technical remarks.** Given the recurring application of certain techniques in the implemented methods, an assortment of technical remarks is presented here before progressing further. The MGS can be applied more than once to prevent loss of orthogonality due to float point roundoffs and the procedure is called a full reorthogonalization (FRO). A pseudocode of orthogonalizing a vector against a basis with MGS, and similarly with an Iterative Classical Gram-Schmidt (ICGS) algorithm is provided in the Supplementary Information as Algorithm 2 and Algorithm 3. The low-complexity  $O(n^3)$  cost of the matrix-vector multiplications in producing  $\mathcal{K}_k$  can be amortized by parallelisms in GPU. When only the lower triangle  $\mathbf{L}$  of  $\mathbf{A}$  is available, the multiplication with an arbitrary vector  $x$  can be done as a sum decomposition  $\mathbf{L}x + \mathbf{L}^T x - \mathbf{D}x$ , where  $\mathbf{D} = : \text{diag}(\mathbf{A})$ . In our implementation, a custom kernel that utilizes the SpMV algorithm in CuSPARSE is written to accommodate the availability of the lower triangle, halving the memory requirement to access the Hessian. In all cases, the convergence rate of the  $i$ -th eigenpair is dictated by the ratio of adjacent exact eigenvalues  $|\frac{\lambda_i}{\lambda_{i+1}}| \leq 1$ , the smaller the better. Note that  $\mathbf{A}$  was shifted upward as  $\mathbf{A} + \mathbf{I}$  to avoid poor condition number. Several strategies to improve the convergence rate for the smallest eigenpairs includes (1) the Shift-and-Inverse technique, where the Ritz pairs of  $(\mathbf{A} - \sigma\mathbf{I})^{-1}$  closest to  $\sigma$  was sought instead and the inverse is incorporated into the matrix-vector product by solving for  $\mathbf{v}^{(k+1)}$  in  $(\mathbf{A} - \sigma\mathbf{I})\mathbf{v}^{(k+1)} = \mathbf{v}^{(k)}$ , (2) the implicitly shifted QR algorithm by Francis<sup>78</sup>, where multiple shifts were applied to the subspace iteration problem typically concerning a small eigenproblem (See part b of Algorithm 5) and (3) filter diagonalization<sup>33,34,57</sup> that magnifies convergence rate for regions of a

spectrum (See Algorithm 4,8 and 9 and later sections). Finally, the six exact eigenvectors  $\mathbf{U}_6$  corresponding to the rigid modes can be removed from  $\mathcal{H}_k$  by an explicit external deflation<sup>36</sup>; the dense deflated matrix  $\mathbf{A} - \sigma_H \mathbf{U}_6 \mathbf{I}_6 \mathbf{U}_6^T$ , with eigenvalues corresponding to  $\mathbf{U}_6$  raised to  $\sigma_H$ , is not stored but incorporated into the matrix-vector multiplications; the sparse  $\mathbf{U}_6$  is stored in CSR format. This is implemented in all three methods and can be optionally called.

**Hermitian Lanczos method (HLM).** The  $\mathbf{V} \in \mathcal{H}_k$  described earlier can be obtained iteratively as  $\mathbf{v}_{k+1}$  by projecting away components of previously obtained  $\mathbf{v}_i, i \in 0 \dots k$  from the current matrix-vector product  $\mathbf{A}\mathbf{v}_k$ . In the passing, the smaller eigenproblem  $\mathbf{T} := \mathbf{V}^T \mathbf{A} \mathbf{V}$  is also produced. Importantly, for a symmetric matrix  $\mathbf{A}$  and  $\mathbf{V} \in \mathcal{H}_k$ , the  $\mathbf{T}$  is symmetric tridiagonal. Therefore, the recursion to obtain  $\mathbf{v}_{k+1}$  simplifies to three terms

$$\mathbf{v}_{k+1} = \mathbf{A}\mathbf{v}_k - \beta_k \mathbf{v}_{k-1} - \alpha_k \mathbf{v}_k \tag{22}$$

Where  $\alpha_k = : \mathbf{v}_k^T \mathbf{A} \mathbf{v}_k$  and  $\beta_k = : \mathbf{v}_{k-1}^T \mathbf{A} \mathbf{v}_k$  are the main diagonal and the first subdiagonal of  $\mathbf{T}$ . This is the essence of the HLM proposed by Lanczos<sup>79</sup>. Collecting the recursion and a rearrangement reveals the Lanczos factorization stopping at the  $k$ -th step.

$$\mathbf{A}\mathbf{V}_k = \mathbf{V}_k \mathbf{T}_k + \beta_k \mathbf{v}_{k+1} \mathbf{e}_k^T \tag{23}$$

A pseudocode of a  $p$ -step Lanczos Factorization is provided in the Supplementary Information as Algorithm 4. Note that the approximate eigenvectors obtained from the HLM will converge to those with extremal eigenvalues, i.e., both the largest and the smallest eigenpairs, but only the smallest eigenpairs are wanted in our problem setting. Further, while the basis  $\mathbf{V}_k$  builds up in each iterative step, eventually we will experience overflow in storage, likely before convergence. To address these concerns, two practical variants of the HLM with different restarting techniques were proposed.

- **Implicitly Restarted Lanczos Method (IRLM).** The IRLM algorithm was proposed by Lehoucq and Sorensen<sup>49</sup>. In IRLM, the Lanczos factorization is supplemented with an implicitly shifted QR algorithm by Francis<sup>78</sup>. The purpose is to shift the smaller eigensystem  $\mathbf{T}^{(j-1)}$  at each restart round  $j$  by the current approximation of the  $m - k$  unwanted eigenvalues such that the convergence rate of the  $k$  wanted smallest eigenvalues is improved. The only information we need to determine these shifts is the sorted approximated eigenvalues and the number of wanted eigenvalues we desired. As a result, after obtaining each QR factorization  $\mathbf{Q}^{(j)} \mathbf{R}^{(j)} = \mathbf{T}^{(j-1)} - \tilde{\lambda}_j \mathbf{I}$ , the Lanczos factorization in Eq. (23) is modified as

$$\mathbf{A}\hat{\mathbf{V}}_k = \hat{\mathbf{V}}_k \hat{\mathbf{T}}_k + \mathbf{v}_{k+1} \hat{\mathbf{b}} \tag{24}$$

where  $\hat{\mathbf{b}} = \mathbf{Q}^{(j)T} \mathbf{e}_k$  and  $\hat{\mathbf{T}} = \mathbf{R}^{(j)} \mathbf{Q}^{(j)} + \tilde{\lambda}_k \mathbf{I}$  and the first  $k$  orthonormal basis is updated as  $\hat{\mathbf{V}}_k = \mathbf{V}_k \mathbf{Q}^{(j)}$ . This is equivalent to performing  $m - k$  simple polynomial filterings<sup>80</sup>. In our implementation, to promote parallelism, the  $m - k$  QR factorizations was performed before a cumulative update of the matrices; this is followed by a full reorthogonalization of the updated basis  $\mathbf{V}_k$ . Note that in IRLM the residual error  $\|\mathbf{r}\|_2$  is bounded by the  $\beta_{k+1}$ , a result due to Paige<sup>81</sup> such that  $\|\mathbf{r}\|_2$  needs not be computed. A pseudocode of our IRLM implementation is provided in the Supplementary Information as Algorithm 5.

- **Thick-Restart Lanczos Method (TRLM).** The TRLM algorithm was proposed by Wu and Simon<sup>58</sup>. In TRLM, the smaller eigenproblem  $\mathbf{T}$  is again solved with its eigenpairs sorted, but rather than performing implicit shift, the  $\mathbf{T}$  is projected onto its wanted eigenvectors  $\mathbf{Y}_k$  as  $\hat{\mathbf{T}} = \mathbf{Y}_k^T \mathbf{T} \mathbf{Y}_k$  and the basis updated to  $\hat{\mathbf{V}} = \mathbf{V} \mathbf{Y}_k$

accordingly. As a result, the Lanczos factorization (in Eq. (23)) is modified as

$$\mathbf{A}\hat{\mathbf{V}}_k = \hat{\mathbf{V}}_k \hat{\mathbf{T}}_k + \beta_k \mathbf{v}_{k+1} \mathbf{e}_k^T \mathbf{Y} \tag{25}$$

When we want to find the next  $\mathbf{v}_{k+1}$  by orthogonalizing  $\mathbf{A}\mathbf{v}_k$  against the previous  $\hat{\mathbf{V}}_k$ , these transforms allow us to compute it with  $\beta_{k-1} (\mathbf{Y}^T \mathbf{e}_{k-1})$  known from the restart

$$\hat{\mathbf{v}}_{k+1} = \mathbf{A}\hat{\mathbf{v}}_k - \hat{\alpha} \hat{\mathbf{v}}_k - \hat{\mathbf{V}}_{k-1} \beta_{k-1} (\mathbf{Y}^T \mathbf{e}_{k-1}) \tag{26}$$

Full reorthogonalization was done in the basis  $\hat{\mathbf{V}}_k$ . A pseudocode of our TRLM implementation is provided in the Supplementary Information as Algorithm 6.

**Jacobi-Davidson method (JDM).** The JDM algorithm was proposed by Fokkema, Sleijpen and Van der Vorst<sup>32</sup>. Similar to the HLM, the JDM also considers approximations in the Krylov subspace, but rather than refining the approximations with a Lanczos factorization, it attempts to solve the following equation

$$\mathbf{A}(\tilde{\mathbf{u}} + \mathbf{z}) = (\tilde{\lambda} + \eta)(\tilde{\mathbf{u}} + \mathbf{z}) \tag{27}$$

where an approximate solution pair  $(\eta, \mathbf{z})$  to correct the approximate eigenpair  $(\tilde{\lambda}, \tilde{\mathbf{u}})$  can be obtained by incorporating the Galerkin condition  $(\tilde{\mathbf{u}} + \mathbf{z})^T \mathbf{z} = 0$  into projectors resulting in the following correction equation

$$(\mathbf{I} - \tilde{\mathbf{u}} \tilde{\mathbf{u}}^T) (\mathbf{A} - \tilde{\lambda} \mathbf{I}) (\mathbf{I} - \tilde{\mathbf{u}} \tilde{\mathbf{u}}^T) \mathbf{z} = -\mathbf{r} \tag{28}$$

The correction equation does not need to be solved exactly such that  $\mathbf{A}$  in this equation can be replaced by a preconditioner of  $\mathbf{A}$  or by simply finding a solution of  $\mathbf{z}$  to certain extent of precision. For the NMA eigenproblems, we cannot find a satisfactory preconditioner that is not dense such that GPU memory is not overflowed. Therefore, in our implementation, the latter strategy was adopted by restricting the number of steps or precision in the solution in the generalized minimal residual iteration (GMRES), which is also a Shift-and-Inverse example. An implementation of the GMRES taking a sparse matrix input was modified from the CuPy package. A pseudocode of our JDM implementation is provided in the Supplementary Information as Algorithm 7.

**Filter diagonalization**

The Lanczos and Davidson methods can be adjusted by filters to isolate certain eigenvalues. The motivation of this Filter Diagonalization (FD) approach is to transform the spectrum such that eigenvalues from wanted intervals become dominant<sup>33,82</sup>. (See Fig. 6e for an illustration.) For example, to obtain the lower spectrum, the Chebyshev-Davidson Method<sup>35</sup> (CDM) removed the Jacobi correction in JDM (which approximates a moving rational filter) and applied a moving low-pass, fixed-degree, Chebyshev polynomial of the first kind to the input matrix. (Algorithm 8) FD can also be applied to solve for the interior rather than the extremal eigenpairs<sup>34,57</sup>. This can be useful when we are extracting a large amount of eigenpairs from the lower end of the spectrum, as then the interior slices of the spectrum can be obtained without storing and orthogonalizing against the filtered subspace from the very lowest end. To achieve this, a band-pass filter, which is an  $M$ -degree Chebyshev expansion of a Dirac-delta-like function damped by a Jackson kernel<sup>33,83</sup>, was developed in the EVSL library<sup>34,57</sup> and incorporated into our INCHING-CTRLM protocol (Algorithm 6, shared with INCHING-TRLM protocol). The band-pass filter  $p(t)$  were obtained following the three-term recurrence of  $T_j$ , the  $j$ -th degree Chebyshev

polynomial of the first kind.

$$T_{j+2}(t) = 2tT_{j+1}(t) - T_j(t); T_0(t) = 1; T_1(t) = t;$$

$$p(t) = \sum_{j=0}^{j=M} g_j^{(M)} \mu_j T_j(t) / \sum_{j=0}^{j=M} g_j^{(M)} \mu_j T_j(\cos(\theta)) \quad (29)$$

The coefficients  $\mu_j$  in the expansion and the damping kernel  $g_j^{(M)}$  were precalculated with the following equations, where  $\delta_{j0}$  is the Kronecker delta and  $\tilde{\pi} = \pi / (M + 2)$  and  $\theta$  is the adjusted center of the arccosine of the corresponding transformed wanted interval  $[\alpha_s, \beta_s]$  within  $[-1, 1]$ . (See reference<sup>34,57</sup> for further details.)

$$g_j^{(M)} = \frac{\sin((j+1)\tilde{\pi})}{(M+2)\sin(\tilde{\pi})} + \left(1 - \frac{j+1}{M+2}\right) \cos(j\tilde{\pi})$$

$$\mu_j = -\frac{1}{2}\delta_{j0} + \cos(j\theta) \quad (30)$$

The filter was applied during matrix-vector multiplication (Algorithm 9), where the spectrum of  $A$  were scaled and shifted to the range of cosine  $[-1, 1]$  using spectral bounds from the Lanczos process (i.e. Algorithm 4) following reference<sup>35</sup>.

### Implementation notes

In many parts of the algorithms presented above, matrix-vector multiplication, whether dense or sparse, is often the calculation bottleneck. In our implementations, this shortcoming is amortized by capitalizing on the high degree of parallelism in GPU and existing techniques built around NVIDIA® GPUs, including CUDA<sup>61</sup> (v. 11.6.2), cuBLAS<sup>62</sup>, and cuSPARSE<sup>63</sup>. To facilitate installation, code readability and version control, the CuPy package (v. 11.5) was used as an application programming interface to access these technologies. In all the algorithms implemented, only the lower triangle of the Hessian matrix was required as input in CSR format. Our INCHING-JDM and INCHING-IRLM implementations were written with reference to the thesis of Geus<sup>84</sup> and the ARPACK package<sup>50</sup> respectively. Our INCHING-TRLM implementation is modified from the CuPy library's default, where memory footprint was optimized by exploiting sum decomposition of the lower triangle; further speedup was achieved by reducing the number of transposes done to the basis set. Where a smaller projected eigenproblem has to be solved, the calculation was done on CPU by calling 'numpy.linalg.eigh' from the NumPy (v. 1.23.5) package<sup>85</sup>. For methods with Chebyshev filter, the implementations in the EVSL library<sup>57</sup> and the Chebyshev-Davidson method<sup>35</sup> were referenced. When a band-pass Chebyshev filter is invoked, the converged eigenvectors are sorted again before off-loading for storage. For benchmarks with ProDy, the ProDy (v.2.4) package<sup>48</sup> was installed with the NumPy (v. 1.23.5) package<sup>85</sup> and the SciPy (v.1.8.0) package<sup>86</sup>. Dense eigenproblems in ProDy are solved using a LAPACK backend called through NumPy. Sparse eigenproblems in ProDy are solved using a ARPACK backend called through SciPy 'scipy.sparse.linalg.eigsh'. Benchmarks on correctness, memory and speed were conducted on a single piece of A100 NVIDIA® GPU with tensor core activated and 80GB GPU memory capacity and the AMD EPYC™ 7513 processor with 64 threads. For systems with more than 2.1 billion non-zero entries, 64-bit integers were used for indexing in the CSR format, otherwise 32-bit integers were used by default. Methods were also tested on a RTX4090 NVIDIA® GPU with 24GB memory and 64-bit indexing used throughout.

### Hyperparameters in calculations

In all cases, the following hyperparameters were used unless otherwise stated. We followed the work of Koehl<sup>30</sup> taking the neighborhood cutoff threshold  $R_C = 8\text{Å}$  for atomic systems. After performing 3DRCM, the indexing slice for each batch of atom was analyzed for presence of

gaps, i.e., regions where no neighbor of the batch is present. The indexing slice was then split into multiple contiguous indexing slices by removing gaps that extend for more than 100 atoms. For INCHING-IRLM, the tolerance of error bound is set to  $10^{-10}$ . For both INCHING-TRLM and INCHING-JDM, the tolerance of residual error is  $10^{-12}$ . The correction equation in INCHING-JDM was solved using the GMRES algorithm, where only 20 steps of minimization at max were allowed. The default maximum allowed number of restarts for INCHING-IRLM, INCHING-TRLM and INCHING-JDM is 15000 steps, but for  $R_C = 6\text{Å}$  benchmarks, at max 30000 steps were allowed to accommodate difficult convergence in some cases. The maximum allowed basis in the Krylov subspace is three times the desired number of output eigenpairs, i.e. which is  $64 \times 3 = 192$  vectors in the benchmark. The 1000 modes presented in Fig. 6 (200 modes in Supplementary Fig. 10) were calculated using INCHING-CTRLM in batches of 64 modes with 128 basis (28 modes with 56 basis) for all structures, though larger basis were affordable; external explicit deflation of the free modes were applied when low-pass filter is used. To scan through the lower end of the spectrum using INCHING-CTRLM, we begin by applying the low-pass filter to obtain the largest of the smallest  $m$  eigenvalues  $\alpha_s$  and set  $\alpha_s - 10^{-10}$  as the left bound of the next interior wanted interval  $[\alpha_s, \beta_s]$ ; the process continues with the band-pass filter until a desired number of eigenpairs are obtained. Note that a binary search was performed to locate  $\beta_s$  such that  $[\alpha_s, \beta_s]$  contains  $< m$  eigenvalues for the band-pass case and  $< 5m$  eigenvalues for the low-pass case. The polynomial degrees were also maximized to maintain the smallest transformed eigenvalues in  $[\alpha_s, \beta_s]$  as 0.7 for both filters such that convergence rate remained constant. (See Fig. 6e.). All programs were stopped if convergence was not achieved within 48 h, wall-clock time. All INCHING programs were stopped if the maximum number of restarts were exceeded. For all benchmarks with  $R_C = 14\text{Å}$ , 64-bit integers indexing were used for consistency.

### Benchmark coordinates

While ideally all the structures on the Protein Data Bank can be considered, we randomly selected structures at an increasing interval. For systems with less than 100 thousand atoms, structures were downloaded in the PDB format; the interval of increase is around 1 thousand atoms. For systems with more than 100 thousand atoms, structures were downloaded as mmCIF format; the interval of increase is around 10 thousand atoms. Note that while the iterative methods all work very well, some structures can contain components not connected within 8 Å threshold, resulting in more than 6 rigid modes. This happens in a small portion of PDB structures when crystallographic water(s) or protein chain(s) with no neighbor within 8 Å were scrupulously put into the crystallographic map. These structures were removed from consideration to avoid confusion. Overall, 85 structures with less than 100 thousand atoms and 31 structures with more than 100 thousand atoms were tested. This benchmark set contains 116 structures in total with number of atoms ranging from around 1 thousand atoms to 2.4 million atoms.

### Megascale atomic structures

We applied our method to some of the largest atomic objects available. The mature HIV-1 capsid structure (PDBID:3J3Q) containing around 2.4 million atoms was obtained from the Protein Data Bank<sup>37</sup> without any modification. The human Nuclear Pore Complex structure (EMDB: EMD14321, PDBID: 7R5J) was downloaded as a bioassembly from the Protein Data Bank<sup>37</sup> and chains LA, MA, NA, OA with clashing linkers were removed; the final structure contains 4.8 million atoms. The atomic structure of the 26 million atoms Faustovirus capsid (PDBID:5J7V) was downloaded from the Protein Data Bank; C $\alpha$  were extracted from the structure resulting in a 5 million pseudo-atom coarse-grained representation; NMA with  $R_C = 14\text{Å}$  was performed with INCHING-CTRLM-HalfSparse with external explicit deflation of the 6

free modes; the first 64 non-zero modes were obtained. The atomic structure of the DNA airplane<sup>67</sup> was obtained from Nanobase.org<sup>87</sup> and the atomic coordinates were reconstructed using TacoxDNA<sup>88</sup>; the final structure includes hydrogens and contains 1.8 million atoms. The NMA of HIV-1 capsid structure were calculated using INCHING-JDM-HalfSparse and the same hyperparameters of numerical methods outlined above. For the DNA airplane and the NPC complex, the NMA was calculated using INCHING-JDM-HalfSparse with external explicit deflation of the 6 free modes, the number of modes to resolve, the maximum allowed basis in the Krylov subspace, the maximum number of allowed steps in GMRES and the maximum number of restarts allowed were revised to the first 30 non-zero normal modes, 120 vectors, 150 GMRES steps and 550000 restarts respectively, otherwise all other hyperparameters are the same. The megascale structures and two movies corresponding to a mature HIV-1 capsid structure (PDBID:3J3Q, Supplementary Movie 1) and a ribosome bound to elongation factor G (PDBID:4V9H, Supplementary Movie 2) were displayed using PyMOL (v2.4.1).

### A logistic kernel to visualize atomic motion

In all the iterative methods, the resultant eigenvector matrix is orthonormalized, which means as size of the system increases a decreasing magnitude of atomic displacement vector will be experienced. In our visualization module, to avoid this scaling problem, a logistic kernel is applied to fine-tune the atomic magnitude  $x_i$  of the eigentensor. First, to avoid eccentricity in the atomic magnitudes, they are clipped between (0.025, 0.975) quantile of the magnitudes. Second, a logistic kernel is applied on the clipped magnitude  $x_i$  for the  $i$ -th atom

$$y_i = G(x_i|\vartheta) = \frac{1}{(1 + \vartheta e^{-x_i})} < 1 \quad (31)$$

where  $\vartheta = 0.05$ . To remove the effect of scaling on the normalized eigenvector as the number of atom increase, the output magnitudes are further centered at  $\tilde{y}_i = \frac{y_i - y^-}{y^+ - y^-} \leq 1$ , where  $(y^-, y^+)$  are the minimum and maximum of  $y_i \forall i$ .

### Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

### Data availability

The data generated in this study, including the atomic coordinates and benchmarks, have been deposited in the Zenodo database under accession code 8087817<sup>89</sup>. Source data are provided with this paper.

### Code availability

The INCHING source code is available on GitHub<sup>90</sup> <https://github.com/jhmlam/Inching>. Jupyter notebooks and Python scripts for the experiments and analyses presented in the paper are available. Frozen versions of the software and associated code for analysis are also available in the Zenodo database under accession code 10645601<sup>91</sup>. All source codes are provided under an Apache License 2.0 license.

### References

- Delarue, M. & Dumas, P. On the use of low-frequency normal modes to enforce collective movements in refining macromolecular structural models. *Proc. Natl. Acad. Sci.* **101**, 6957–6962 (2004).
- Tama, F., Miyashita, O. & Brooks, C. L. III Normal mode based flexible fitting of high-resolution structure into low-resolution experimental data from cryo-EM. *J. Struct. Biol.* **147**, 315–326 (2004).
- Gur, M., Madura, J. D. & Bahar, I. Global transitions of proteins explored by a multiscale hybrid methodology: application to adenylate kinase. *Biophys. J.* **105**, 1643–1652 (2013).
- Franklin, J., Koehl, P., Doniach, S. & Delarue, M. MinActionPath: maximum likelihood trajectory for large-scale structural transitions in a coarse-grained locally harmonic energy landscape. *Nucleic Acids Res.* **35**, W477–W482 (2007).
- Bakan, A. & Bahar, I. The intrinsic dynamics of enzymes plays a dominant role in determining the structural changes induced upon inhibitor binding. *Proc. Natl. Acad. Sci.* **106**, 14349–14354 (2009).
- Shrivastava, I. H. & Bahar, I. Common mechanism of pore opening shared by five different potassium channels. *Biophys. J.* **90**, 3929–3940 (2006).
- Lezon, T. R. & Bahar, I. Constraints imposed by the membrane selectively guide the alternating access dynamics of the glutamate transporter GltPh. *Biophys. J.* **102**, 1331–1340 (2012).
- Bahar, I., Lezon, T. R., Bakan, A. & Shrivastava, I. H. Normal mode analysis of biomolecular structures: Functional mechanisms of membrane proteins. *Chem. Rev.* **110**, 1463–1497 (2010).
- Levitt, M. Conformation analysis of proteins. (University of Cambridge, 1972).
- Levitt, M., Sander, C. & Stern, P. S. The normal modes of a protein: Native bovine pancreatic trypsin inhibitor. *Int. J. Quantum Chem.* **24**, 181–199 (1983).
- Brooks, B. & Karplus, M. Normal modes for specific motions of macromolecules: application to the hinge-bending mode of lysozyme. *Proc. Natl. Acad. Sci.* **82**, 4995–4999 (1985).
- Brooks, B. & Karplus, M. Harmonic dynamics of proteins: normal modes and fluctuations in bovine pancreatic trypsin inhibitor. *Proc. Natl. Acad. Sci. USA.* **80**, 6571–6575 (1983).
- Ma, J. Usefulness and limitations of normal mode analysis in modeling dynamics of biomolecular complexes. *Structure* **13**, 373–380 (2005).
- Hayward, S. & Groot, B. L. De. Normal modes and essential dynamics. in *Molecular Modeling of Proteins* 89–106 (Humana Press). [https://doi.org/10.1007/978-1-59745-177-2\\_5](https://doi.org/10.1007/978-1-59745-177-2_5). (2008)
- Tirion, M. M. Large amplitude elastic motions in proteins from a single-parameter, atomic analysis. *Phys. Rev. Lett.* **77**, 1905–1908 (1996).
- Hinsen, K., Thomas, A. & Field, M. J. Analysis of domain motions in large proteins. *Proteins: Struct. Funct. Bioinforma.* **34**, 369–382 (1999).
- Hinsen, K. Analysis of domain motions by approximate normal mode calculations. *Proteins: Struct. Funct. Bioinforma.* **33**, 417–429 (1998).
- Tama, F. & Sanejouand, Y.-H. Conformational change of proteins arising from normal mode calculations. *Protein Eng., Des. Select.* **14**, 1–6 (2001).
- Brink, J. et al. Experimental verification of conformational variation of human fatty acid synthase as predicted by normal mode analysis. *Structure* **12**, 185–191 (2004).
- Schilbach, S. et al. Structures of transcription pre-initiation complex with TFIID and mediator. *Nature* **551**, 204–209 (2017).
- Jin, Q. et al. Iterative elastic 3D-to-2D alignment method using normal modes for studying structural dynamics of large macromolecular complexes. *Structure* **22**, 496–506 (2014).
- Krieger, J. M., Sorzano, C. O. S., Carazo, J. M. & Bahar, I. Protein dynamics developments for the large scale and cryoEM: case study of ProDy 2.0. *Acta Cryst. D.* **78**, 399–409 (2022).
- Zhang, Y. et al. State-dependent sequential allostery exhibited by chaperonin TRiC/CCT revealed by network analysis of Cryo-EM maps. *Prog. Biophys. Mol. Biol.* **160**, 104–120 (2021).
- Vuillemot, R., Miyashita, O., Tama, F., Rouiller, I. & Jonic, S. NMMD: Efficient Cryo-EM flexible fitting based on simultaneous normal mode and molecular dynamics atomic displacements. *J. Mol. Biol.* **434**, 167483 (2022).
- Atilgan, A. R. et al. Anisotropy of fluctuation dynamics of proteins with an elastic network model. *Biophys. J.* **80**, 505–515 (2001).

26. Tama, F., Gadea, F. X., Marques, O. & Sanejouand, Y.-H. Building-block approach for determining low-frequency normal modes of macromolecules. *Proteins: Struct. Funct. Bioinforma.* **41**, 1–7 (2000).
27. Li, G. & Cui, Q. A coarse-grained normal mode approach for macromolecules: An efficient implementation and application to Ca<sup>2+</sup>-ATPase. *Biophys. J.* **83**, 2457–2474 (2002).
28. Perahia, D. & Mouawad, L. Computation of low-frequency normal modes in macromolecules: Improvements to the method of diagonalization in a mixed basis and application to hemoglobin. *Comput. Chem.* **19**, 241–246 (1995).
29. Marques, O. & Sanejouand, Y.-H. Hinge-bending motion in citrate synthase arising from normal mode calculations. *Proteins: Struct. Funct. Bioinforma.* **23**, 557–560 (1995).
30. Koehl, P. Large eigenvalue problems in coarse-grained dynamic analyses of supramolecular systems. *J. Chem. Theory Comput.* **14**, 3903–3919 (2018).
31. Sleijpen, G. L. G., & Van der Vorst, H. A. A Jacobi-Davidson iteration method for linear eigenvalue problems. *SIAM Rev.* **42**, 267–293 (2000).
32. Sleijpen, G. L. G., Booten, A. G. L., Fokkema, D. R. & van der Vorst, H. A. Jacobi-Davidson type methods for generalized eigenproblems and polynomial eigenproblems. *Bit Numer Math.* **36**, 595–633 (1996).
33. Wall, M. R. & Neuhauser, D. Extraction, through filter-diagonalization, of general quantum eigenvalues or classical normal mode frequencies from a small number of residues or a short-time segment of a signal. I. Theory and application to a quantum-dynamics model. *J. Chem. Phys.* **102**, 8011–8022 (1995).
34. Li, R., Xi, Y., Vecharynski, E., Yang, C. & Saad, Y. A Thick-Restart lanczos algorithm with polynomial filtering for hermitian eigenvalue problems. *SIAM J. Sci. Comput.* **38**, A2512–A2534 (2016).
35. Zhou, Y. & Saad, Y. A Chebyshev–Davidson algorithm for large symmetric eigenproblems. *SIAM J. Matrix Anal. Appl.* **29**, 954–971 (2007).
36. Parlett, B. N. 5. Deflation. in *The Symmetric Eigenvalue Problem* 87–92 (Society for Industrial and Applied Mathematics). <https://doi.org/10.1137/1.9781611971163.ch5> (1998).
37. Berman, H. M. et al. The protein data bank. *Nucleic Acids Res.* **28**, 235–242 (2000).
38. Zhao, G. et al. Mature HIV-1 capsid structure by cryo-electron microscopy and all-atom molecular dynamics. *Nature* **497**, 643–646 (2013).
39. Townshend, R. J. L. et al. Geometric deep learning of RNA structure. *Science* **373**, 1047–1051 (2021).
40. Baek, M. et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science* **373**, 871–876 (2021).
41. Papadimitriou, Ch. H. The NP-Completeness of the bandwidth minimization problem. *Computing* **16**, 263–270 (1976).
42. Cuthill, E. & McKee, J. Reducing the bandwidth of sparse symmetric matrices. in *Proceedings of the 1969 24th national conference* 157–172 (Association for Computing Machinery, New York, NY, USA). <https://doi.org/10.1145/800195.805928> (1969).
43. George, A. & Liu, J. W. Computer Solution of Large Sparse Positive Definite Systems (Alan George and Joseph W. Liu. *SIAM Rev.* **26**, 289–291 (1984).
44. Azad, A., Jacquelin, M., Buluç, A. & Ng, E. G. The Reverse Cuthill-McKee Algorithm in Distributed-Memory. in *2017 IEEE International Parallel and Distributed Processing Symposium (IPDPS)* 22–31. <https://doi.org/10.1109/IPDPS.2017.85> (2017).
45. Mlakar, D., Winter, M., Parger, M. & Steinberger, M. Speculative Parallel Reverse Cuthill-McKee Reordering on Multi- and Many-core Architectures. in *2021 IEEE International Parallel and Distributed Processing Symposium (IPDPS)* 703–713. <https://doi.org/10.1109/IPDPS49936.2021.00080> (2021).
46. Karantasis, K. I., Lenharth, A., Nguyen, D., Garzarán, M. J. & Pingali, K. Parallelization of Reordering Algorithms for Bandwidth and Wavefront Reduction. in *SC '14: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* 921–932. <https://doi.org/10.1109/SC.2014.80>. (2014)
47. Bentley, J. L. Multidimensional binary search trees used for associative searching. *Commun. ACM* **18**, 509–517 (1975).
48. Bakan, A., Meireles, L. M. & Bahar, I. ProDy: Protein dynamics inferred from theory and experiments. *Bioinformatics* **27**, 1575–1577 (2011).
49. Lehoucq, R. B. & Sorensen, D. C. Deflation techniques for an implicitly restarted arnoldi iteration. *SIAM J. Matrix Anal. Appl.* **17**, 789–821 (1996).
50. Lehoucq, R., Sorensen, D. & Yang, C. *ARPACK Users' Guide: Solution of large-scale eigenvalue problems with implicitly restarted arnoldi methods*. 43–66 (SIAM, 1998). <https://doi.org/10.1137/1.9780898719628>.
51. Anderson, E. et al. *LAPACK Users' Guide*. (Society for Industrial and Applied Mathematics). <https://doi.org/10.1137/1.9780898719604>. (1999)
52. Yu, V. W. et al. GPU-acceleration of the ELPA2 distributed eigensolver for dense symmetric and hermitian eigenproblems. *Comput. Phys. Commun.* **262**, 107808 (2021).
53. Yu, V. W. et al. ELSI — An open infrastructure for electronic structure solvers. *Comput. Phys. Commun.* **256**, 107459 (2020).
54. Wu, X., Davidović, D., Achilles, S. & Di Napoli, E. ChASE: a distributed hybrid CPU-GPU eigensolver for large-scale hermitian eigenvalue problems. in *Proceedings of the Platform for Advanced Scientific Computing Conference* 1–12 (Association for Computing Machinery, New York, NY, USA). <https://doi.org/10.1145/3539781.3539792>. (2022)
55. Sgherzi, F., Parravicini, A. & Santambrogio, M. D. A mixed precision, multi-GPU design for large-scale Top-K sparse eigenproblems. in *2022 IEEE International Symposium on Circuits and Systems (ISCAS)* 1259–1263. <https://doi.org/10.1109/ISCAS48785.2022.9937893>. (2022)
56. Anzt, H. et al. Optimizing Krylov Subspace Solvers on Graphics Processing Units. in *2014 IEEE International Parallel & Distributed Processing Symposium Workshops* 941–949. <https://doi.org/10.1109/IPDPSW.2014.107>. (2014)
57. Li, R., Xi, Y., Erlandson, L. & Saad, Y. The eigenvalues slicing library (EVL): Algorithms, implementation, and software. *SIAM J. Sci. Comput.* **41**, C393–C415 (2019).
58. Wu, K. & Simon, H. Thick-restart lanczos method for large symmetric eigenvalue problems. *SIAM J. Matrix Anal. Appl.* **22**, 602–616 (2000).
59. Bell, N. & Garland, M. Efficient Sparse Matrix-Vector Multiplication on CUDA | Research. *NVIDIA Technical Report NVR-2008-004* [https://research.nvidia.com/publication/2008-12\\_efficient-sparse-matrix-vector-multiplication-cuda](https://research.nvidia.com/publication/2008-12_efficient-sparse-matrix-vector-multiplication-cuda) (2008).
60. Okuta, R., Unno, Y., Nishino, D., Hido, S. & Loomis, C. CuPy: A NumPy-Compatible Library for NVIDIA GPU Calculations. *Proceedings of Workshop on Machine Learning Systems (LearningSys) in The Thirty-first Annual Conference on Neural Information Processing Systems (NIPS)* (2017).
61. Nickolls, J., Buck, I., Garland, M. & Skadron, K. Scalable Parallel Programming with CUDA: Is CUDA the parallel programming model that application developers have been waiting for? *Queue* **6**, 40–53 (2008).
62. cuBLAS. <https://docs.nvidia.com/cuda/cublas/>.
63. cuSPARSE. <https://docs.nvidia.com/cuda/cusparse/>.
64. Pornillos, O., Ganser-Pornillos, B. K. & Yeager, M. Atomic-level modelling of the HIV capsid. *Nature* **469**, 424–427 (2011).
65. Caspar, D. L. D. & Klug, A. Physical principles in the construction of regular viruses. *Cold Spring Harb. Symp. Quant. Biol.* **27**, 1–24 (1962).

66. Mosalaganti, S. et al. AI-based structure prediction empowers integrative structural analysis of human nuclear pores. *Science* **376**, eabm9506 (2022).
67. Huang, C.-M., Kucinic, A., Johnson, J. A., Su, H.-J. & Castro, C. E. Integrated computer-aided engineering and design for DNA assemblies. *Nat. Mater.* **20**, 1264–1271 (2021).
68. Klose, T. et al. Structure of faustovirus, a large dsDNA virus. *Proc. Natl Acad. Sci.* **113**, 6206–6211 (2016).
69. Hoffmann, A. & Grudin, S. NOLB: Nonlinear rigid block normal-mode analysis method. *J. Chem. Theory Comput.* **13**, 2123–2134 (2017).
70. Pandey, M. et al. The transformational role of GPU computing and deep learning in drug discovery. *Nat. Mach. Intell.* **4**, 211–221 (2022).
71. Ikegami, T., Sakurai, T. & Nagashima, U. A filter diagonalization for generalized eigenvalue problems based on the Sakurai–Sugiura projection method. *J. Computational Appl. Math.* **233**, 1927–1936 (2010).
72. Gulati, A. et al. Structural studies on chimeric Sesbania mosaic virus coat protein: Revisiting SeMV assembly. *Virology* **489**, 34–43 (2016).
73. Yang, Z., Bahar, I. & Widom, M. Vibrational dynamics of icosahedrally symmetric biomolecular assemblies compared with predictions based on continuum elasticity. *Biophys. J.* **96**, 4438–4448 (2009).
74. Uto, S. et al. Mutual relationships between structural and functional changes in a PsbM-deletion mutant of photosystem II. *Faraday Discuss.* **198**, 107–120 (2017).
75. Lee, D. T. & Wong, C. K. Worst-case analysis for region and partial region searches in multidimensional binary search trees and balanced quad trees. *Acta Inform.* **9**, 23–29 (1977).
76. Hales, T. et al. A formal proof of the Kepler conjecture. *Forum Math., Pi* **5**, e2 (2017).
77. Cuppen, J. J. M. A divide and conquer method for the symmetric tridiagonal eigenproblem. *Numer. Math.* **36**, 177–195 (1980).
78. Francis, J. G. F. The QR transformation A unitary analogue to the LR transformation—Part 1. *Comput. J.* **4**, 265–271 (1961).
79. Lanczos, C. An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. *J. Res. Natl Bur. Stand.* **45**, 255–282 (1950).
80. Sorensen, D. C. Implicit application of polynomial filters in a k-Step Arnoldi method. *SIAM J. Matrix Anal. Appl.* **13**, 357–385 (1992).
81. Paige, C. C., Parlett, B. N. & van der Vorst, H. A. Approximate solutions and eigenvalue bounds from Krylov subspaces. *Numer. Linear Algebra Appl.* **2**, 115–133 (1995).
82. Saad, Y. 7. Filtering and restarting techniques. in *Numerical Methods for Large Eigenvalue Problems* 163–191 (Society for Industrial and Applied Mathematics). <https://doi.org/10.1137/1.9781611970739.ch7>. (2011)
83. Weiße, A., Wellein, G., Alvermann, A. & Fehske, H. The kernel polynomial method. *Rev. Mod. Phys.* **78**, 275–306 (2006).
84. Geus, R. The Jacobi-Davidson algorithm for solving large sparse symmetric eigenvalue problems with application to the design of accelerator cavities. (ETH Zurich). <https://doi.org/10.3929/ethz-a-004469464>. (2002)
85. Harris et al. Array programming with NumPy. *Nature* **585**, 357–362 (2020).
86. Virtanen, P. et al. SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nat. Methods* **17**, 261–272 (2020).
87. Poppleton, E., Mallya, A., Dey, S., Joseph, J. & Šulc, P. Nanobase.org: a repository for DNA and RNA nanostructures. *Nucleic Acids Res.* **50**, D246–D252 (2022).
88. Suma, A. et al. TacoxDNA: A user-friendly web server for simulations of complex DNA structures, from single strands to origami. *J. Comput. Chem.* **40**, 4 (2019).
89. Lam, J. H. Scalable computation of anisotropic vibrations for large macromolecular assemblies, InchingDataRepo. <https://doi.org/10.5281/zenodo.8087817>, (2023).
90. Lam, J. H. Scalable computation of anisotropic vibrations for large macromolecular assemblies, Inching, <https://doi.org/10.5281/zenodo.10645600>, (2024).
91. Lam, J. H. Scalable computation of anisotropic vibrations for large macromolecular assemblies, Inching: Release zenodo, <https://doi.org/10.5281/zenodo.10645601> (2024).

## Acknowledgements

This study was supported by internal funding from USC Dornsife college to V.K., and an NSF grant, OAC-2118061 to A.N. The authors acknowledge the Center for Advanced Research Computing (CARC) at the University of Southern California for providing computing resources that have contributed to the research results reported within this publication. URL: <https://carc.usc.edu>.

## Author contributions

J.H.L., A.N. and V.K. conceived the research. A.N. and V.K. supervised the research. J.H.L. developed the algorithms and led the project under the guidance of A.N. and V.K. All authors discussed the results and contributed to the writing of the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-024-47685-8>.

**Correspondence** and requests for materials should be addressed to Aiichiro Nakano or Vsevolod Katritch.

**Peer review information** *Nature Communications* thanks Giuseppe M. J. Barca, and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

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

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024