# Persistent topology for cryo-EM data analysis

Kelin Xia[1] and Guo-Wei Wei[1,2,3] *
[1]Department of Mathematics
Michigan State University, MI 48824, USA
[2]Department of Electrical and Computer Engineering
Michigan State University, MI 48824, USA
[3]Department of Biochemistry and Molecular Biology
Michigan State University, MI 48824, USA

April 11, 2018

**Abstract**

In this work, we introduce persistent homology for the analysis of cryo-electron microscopy (cryo-EM) density maps. We identify the topological fingerprint or topological signature of noise, which is widespread in cryo-EM data. For low signal to noise ratio (SNR) volumetric data, intrinsic topological features of biomolecular structures are indistinguishable from noise. To remove noise, we employ geometric flows which are found to preserve the intrinsic topological fingerprints of cryo-EM structures and diminish the topological signature of noise. In particular, persistent homology enables us to visualize the gradual separation of the topological fingerprints of cryo-EM structures from those of noise during the denoising process, which gives rise to a practical procedure for prescribing a noise threshold to extract cryo-EM structure information from noise contaminated data after certain iterations of the geometric flow equation. To further demonstrate the utility of persistent homology for cryo-EM data analysis, we consider a microtubule intermediate structure (EMD-1129). Three helix models, an alpha-tubulin monomer model, an alpha- and beta-tubulin model, and an alpha- and beta-tubulin dimer model, are constructed to fit the cryo-EM data. The least square fitting leads to similarly high correlation coefficients, which indicates that structure determination via optimization is an ill-posed inverse problem. However, these models have dramatically different topological fingerprints. Especially, linkages or connectivities that discriminate one model from another one, play little role in the traditional density fitting or optimization, but are very sensitive and crucial to topological fingerprints. The intrinsic topological features of the microtubule data are identified after topological denoising. By a comparison of the topological fingerprints of the original data and those of three models, we found that the third model is topologically favored. The present work offers persistent homology based new strategies for topological denoising and for resolving ill-posed inverse problems.

Key words: Cryo-EM, Topological signature, Geometric flow, Topological denoising, Topology-aided structure determination.

---

*Address correspondences to Guo-Wei Wei. E-mail:wei@math.msu.edu

**Contents**

# 1 Introduction

The quantitative understanding of structure, function, dynamics and transport of biomolecules is a fundamental theme in contemporary life sciences. Geometric analysis and associated biophysical modeling have been the main workhorse in revealing the structure-function relationship of biomolecules and contribute enormously to the present understanding of biomolecular systems. However, biology encompasses over more than twenty orders of magnitude in time scales from electron transfer and ionization on the scale of femtoseconds to organism life spanning over tens of years, and over fifteen orders of magnitude in spatial scales from electrons and nuclei to organisms. The intriguing complexity and extraordinarily large number of degrees of freedom of biological systems give rise to formidable challenges to their quantitative description and theoretical prediction. Most biological processes, such as signal transduction, gene regulation, DNA specification, transcription and post transcriptional modification, are essentially intractable for atomistic geometric analysis and biophysical simulations, let alone *ab-initio* quantum mechanical descriptions. Therefore, the complexity of biology and the need for its understanding offer an extraordinary opportunity for innovative theories, methodologies, algorithms and tools.

The study of subcellular structures, organelles and large multiprotein complexes has become one of the major trends in structural biology. Currently, one of the most powerful tools for the aforementioned systems is cryo-electron microscopy (cryo-EM), although other techniques, such as macromolecular X-ray crystallography, nuclear magnetic resonance (NMR), electron paramagnetic resonance (EPR), multiangle light scattering, confocal laser-scanning microscopy, small angle scattering, ultra fast laser spectroscopy, etc., are useful for structure determination in general.[1–5] In cryo-EM experiments, samples are bombarded by electron beams at cryogenic temperatures to improve the signal to noise ratio (SNR). The working principle is based on the projection (thin film) specimen scans collected from many different directions around one or two axes, and the Radon transform for the creation of three-dimensional (3D) images. One of major advantages of cryo-EM is that it allows the imaging of specimens in their native environment. Another major advantage is its capability of providing 3D mapping of entire cellular proteomes together with their detailed interactions at nanometer or subnanometer resolution.[1–4] The resolution of cryo-EM maps has been improved dramatically in the past two decades, thanks to the technical advances in experimental hardware, noise reduction and image segmentation techniques. By further taking the advantage of symmetric averaging, many cryo-EM based virus structures have already achieved a resolution that can be interpreted in terms of atomic models. There have been a variety of elegant methods[6–11] and software packages in cryo-EM structural determination.[12–18]

Most biological specimens are extremely radiation sensitive and can only sustain a limited electron dose of illumination. As a result, cryo-EM images are inevitably of low SNR and limited resolution.[5] In fact, the SNRs of cryo-tomograms for subcellular structures, organelles and large multi-protein complexes are typically in the neighborhood of 0.01.[5] To make the situation worse, the image contrast, which depends on the difference between electron scattering cross sections of cellular components, is also very low in most biological systems. Consequently, cryo-EM maps often do not contain adequate information to offer unambiguous atomic-scale structural reconstruction of biological specimens. Additional information obtained from other techniques, such as X-ray crystallography, NMR and computer simulation, is indispensable to achieve subnanometer resolutions. However, for cryo-EM data that do not have much additional information obtained from other techniques, the determination of what proteins are involved can be a challenge, not to mention subnanometer structural resolution.

To improve the SNR and image contrast of cryo-EM data, a wide variety of denoising algorithms has been employed.[19–26] Standard techniques, such as bilateral filter[23–25] and iterative median filtering[26] have been utilized for noise reduction. Additionally, wavelets and related techniques have also been developed for cryo-EM noise removing.[19] Moreover, anisotropic diffusion[20,21] or Beltrami flow[22] approach has been proposed for cryo-EM signal recovering. However, cryo-EM data denoising is far from adequate and remains a challenge due to the extremely low SNRs and other technical complications.[5,27–29] For example, one of difficulties is how to distinguish signal from noise in cryo-EM data. As a result, one does not know when to stop or how to apply a threshold in an iterative noise removing process. There is a pressing need for innovative mathematical approaches to further tackle this problem.

Recently, persistent homology has been advocated as a new approach for dealing with big data sets.[30–33] In general, persistent homology characterizes the geometric features with persistent topological invariants by

defining a scale parameter relevant to topological events. The essential difference between the persistent homology and traditional topological approaches is that traditional topological approaches describe the topology of a given object in truly metric free or coordinate free representations, while persistent homology analyzes the persistence of the topological features of a given object via a filtration process, which creates a family of similar copies of the object at different spatial resolutions. Technically, a series of nested simplicial complexes is constructed from a filtration process, which captures topological structures continuously over a range of spatial scales. The involved topological features are measured by their persistent intervals. Persistent homology is able to embed geometric information to topological invariants so that "birth" and "death" of isolated components, circles, rings, loops, pockets, voids or cavities at all geometric scales can be monitored by topological measurements. The basic concept of persistent homology was introduced by Frosini and Landi[34] and by Robins[35] in 1999 independently. Edelsbrunner et al.[36] introduced the first efficient computational algorithm, and Zomorodian and Carlsson[37] generalized the concept. A variety of elegant computational algorithms has been proposed to track topological variations during the filtration process.[38–42] Often, the persistent diagram can be visualized through barcodes,[32] in which various horizontal line segments or bars are the homology generators lasted over filtration scales. It has been applied to a variety of domains, including image analysis,[43–46] image retrieval,[47] chaotic dynamics verification,[48,49] sensor network,[50] complex network,[51,52] data analysis,[31,53–56] computer vision,[45] shape recognition[57] and computational biology.[58–60]

The concept of persistent homology has also been used for noise reduction. It is generally believed that short lifetime events (or bars) are of less importance and thus regarded as "noise" while long lifetime ones are considered as "topological signals",[61] although this idea was challenged in a recent work.[62] In topological data analysis, pre-processing algorithms are needed to efficiently remove these noise. Depending on the scale of a feature, a simple approach is to pick up a portion of landmark points as a representative of topological data.[63] The points can be chosen randomly, spatially evenly, or from extreme values. More generally, certain functions can be defined as a guidance for node selection to attenuate the noise effect, which is known as thresholding. Clustering algorithms with special kernel functions can also be employed to recover topological signal.[61] All of these methods can be viewed as a process of data sampling without losing the significant topological features. They rely heavily on the previous knowledge of the geometric or statistic information. In contrast, topological simplification,[36,64,65] which is to remove the simplices and/or the topological attributes that do not survive certain threshold, focuses directly on the persistence of topological invariant. In contrast, Gaussian noise is known to generate a band of bars distributes over a wide range in the barcode representation.[66] Thank to the pairing algorithm, persistence of a homology group is measured through an interval represented by a simplex pair. If the associated topological invariant is regarded less important, simplices related to this simplex pair are reordered. This approach, combined with Morse theory, proves to be a useful tool for denoising,[64,65] as it can alters the data locally in the region defined as noise. Additionally, statistical analysis has been carried out to provide confidence sets for persistence diagram. However, persistent homology has not been utilized for cryo-EM data noise reduction, to our knowledge.

A large amount of experimental data for macroproteins and protein-protein complexes has been obtained from cryo-EM. To analyze these structural data, it is a routine procedure to fit them with the available high-resolution crystal structures of their component proteins. This approach has been shown to be efficient for analyzing many structures and has been integrated into many useful software packages such as Chimera.[67] However, this docking process is limited by data quality. For some low resolution data, which usually also suffer from low SNRs, there is enormous ambiguity in structure fitting or optimization, i.e., a mathematically ill-posed inverse problem. Sometimes, high correlation coefficients can be attained simultaneously in many alternative structures, while none of them proves to be biologically meaningful. Basically, the fitting or optimization emphasizes more on capturing "bulk" regions, which is reasonable as greater similarities in density distributions imply higher possibility. However, little attention is paid to certain small "linkage" regions, which play important roles in biological system especially in macroproteins and protein-protein complexes. Different linkage parts generate different connectivity, and thus directly influence biomolecular flexibility, rigidity, and even its functions. Since persistent homology is equally sensitive to both bulk regions and small linkage regions, it is able to make a critical judgment on the model selection in structure determination, However, nothing has been reported on persistent homology based solution to ill-posed inverse problems, to our knowledge.

Although persistent homology has been applied to a variety of fields, the successful use of persistent homology is mostly limited to characterization, identification and analysis (CIA). Indeed, persistent homology has

seldom employed for quantitative prediction. Recently, we have introduced molecular topological fingerprints (MTFs) based on persistent homology analysis of molecular topological invariants of biomolecules.[62] We have utilized MTFs to reveal the topology-function relationship of macromolecules. It was found that protein flexibility and folding stability are strongly correlated to protein topological connectivity, characterized by the persistence of topological invariants (i.e., accumulated bar lengths).[62] Most recently, we have employed persistent homology to analyze the structure and function of nano material, such as nanotubes and fullerenes. The MTFs are utilized to quantitatively predict total curvature energies of fullerene isomers.[68]

The overall objective of this work is to explore the utility of persistent homology for cryo-EM analysis. First, we propose a topology based algorithm for cryo-EM noise reduction and clean-up. We study the topological fingerprint or topological signature of noise and its relation to the topological fingerprint of cryo-EM structures. We note that the histograms of topological invariants of the Gaussian random noise have Gaussian distributions in the filtration parameter space. Contrary to the common belief that short barcode bars correspond to noise, it is found that there is an inverse relation between the SNR and the band widths of topological invariants, i.e., the lower SNR, the larger barcode band width is. Therefore, at a low SNR, noise can produce long persisting topological invariants or bars in the barcode presentation. Moreover, for cryo-EM data of low SNRs, intrinsic topological features of the biomolecular structure are hidden in the persistent barcodes of noise and indistinguishable from noise contributions. To recover the topological features of biomolecular structures, geometric flow equations are employed in the present work. It is interesting to note that topological features of biomolecular structures persist, while the topological fingerprint of noise moves to the right during the geometric flow iterations. As such, "signal" and noise separate from each other during the geometric flow based denoising process and make it possible to prescribe a precise noise threshold for the noise removal after certain iterations. We demonstrate the efficiency of our persistent homology controlled noise removal algorithm for both synthetic data and cryo-EM density maps.

Additionally, we introduce persistent homology as a new strategy for resolving the ill-posed inverse problem in cryo-EM structure determination. Although the structure determination of microtubule data EMD 1129 is used as an example, similar problems are widespread in other intermediate resolution and low resolution cryo-EM data. As EMD 1129 is contaminated by noise, a preprocess of denoising is carried out by using our persistent homology controlled geometric flow algorithm. A helix backbone is obtained for the microtubule intermediate structure. Based on the assumption that the voxels with high electron density values are the centers of tubulin proteins, we construct three different microtubule models, namely a monomer model, a two-monomer model, and a dimer model. We have found that all three models give rise to essentially the same high correlation coefficients, i.e., 0.9601, 0.9607 and 0.9604, with the cryo-EM data. This ambiguity in structure fitting is very common with intermediate and low resolution data. Fortunately, after our topology based noise removal, the topology fingerprint of microtubule data is very unique, which is true for all cryo-EM data or data generated by using other molecular imaging modalities. It is interesting to note that although three models offer the same correlation coefficients with the cryo-EM data, their topological fingerprints are dramatically different. It is found that the topological fingerprint of the microtubule intermediate structure (EMD 1129) can be captured only when two conditions are simultaneously satisfied: first, there must exist two different types of monomers, and additionally, two type of monomers from dimers. Therefore, based on topological fingerprint analysis, we can determine that only the third model is a correct model for microtubule data EMD 1129.

The rest of this paper is organized as follows. The essential methods and algorithms for geometric and topological modelings of biomolecular data are presented in Section 2. Approaches for geometric modeling, which are necessary for topological analysis, are briefly discussed. Methods for persistent homology analysis are described in detail. We illustrate the use of topological methods with both synthetic volumetric data and cryo-EM density maps. Their persistence of topological invariants is represented by barcodes. The geometric interpretation of the topological features is given. Section 3 is devoted to the persistent homology based noise removal. The characterization of Gaussian noise is carried out over a variety of SNRs to understand noise topological signature. Based on this understanding, we design a persistent homology monitored and controlled algorithm for noise removal, which is implemented via the geometric flow. Persistent homology guided denoising is applied to the analysis of a supramolecular filamentous complex. In Section 4, we demonstrate topology-aided structure determination of microtubule cryo-EM data. Several aspects are considered including helix backbone evaluation, coarse-grained modeling and topology-aided structure design and evaluation. We
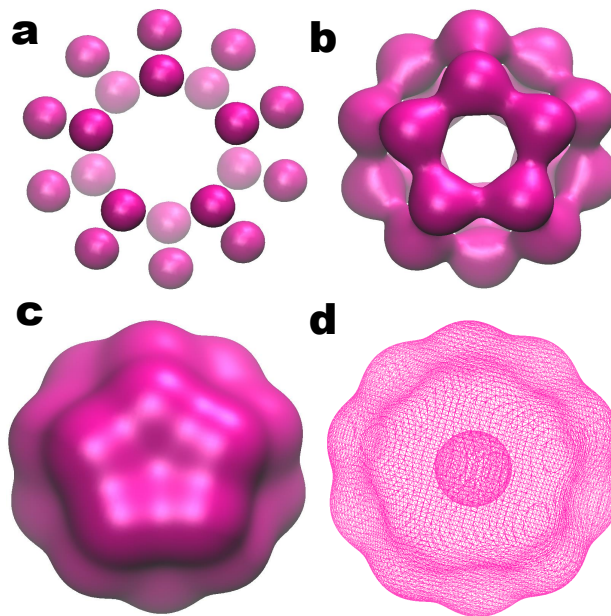
Figure 1: A series of surfaces extracted from fullerene $C_{20}$ density function (1) with $w_j = 1$, $\kappa = 2$ and $\sigma_j = 0.5$Å. The isovalues for subfigures **a**, **b** and **c** are 0.7, 0.6 and 0.4, respectively. Subfigure **d** is wire-frame surface representation of subfigure **c**. The Betti numbers can be directly obtained through identifying the numbers of connected components, circles or loops, and voids. In **a**, $\beta_0$ is 20, $\beta_1$ and $\beta_2$ are 0; In **b**, $\beta_0 = 1$, $\beta_1 = 11$ and $\beta_2 = 0$; In **c** and **d**, $\beta_0 = 1$, $\beta_1 = 0$ and $\beta_2 = 1$. Attention should be paid to $\beta_1$ for the structure in **b**. The $\beta_1$ is not exactly equal to the total number ($N_1$) of the circles or loops, and instead it is equal to $N_1 - 1$. This is due to the fact that the corresponding homology group has a basis with only $N_1 - 1$ elements. Therefore, one of elements can be expressed as the " linear combination" of the rest.

show that topology is able to resolve ill-posed inverse problem. This paper ends with a conclusion.

## 2  Geometric and topological modelings of biomolecular data

Persistent homology has been utilized to analyze biomolecular data, which are collected by different experimental means, such as macromolecular X-ray crystallography, NMR, EPR, etc. Due to their different origins, these data may be available in different formats, which requires appropriate topological tools for their analysis. Additionally, their quality, i.e., resolution and SNR varies for case to case, and thus, a preprocessing may be required. Moreover, although biomolecular structures are not a prerequisite for persistent homology analysis, the understanding of biomolecular structure, function and dynamics is crucial for the interpretation of topological results. As a consequence, appropriate geometric modeling[69] is carried out in a close association with topological analysis. Furthermore, information from geometric and topological modelings is in turn, very valuable for data preprocessing and denoising. Finally, topological information is shown to be crucial for geometric modeling, structural determination and ill-posed inverse problems.

### 2.1  Geometric modeling of biomolecules

Geometric modeling of biomolecules gives rise to their structural information, which is of paramount importance for biological understanding and structure-function relationship.[69, 70] Geometric modeling typically begins with experimental data. There are two major repositories, namely, Protein Data Bank (PDB) and Electron Microscopy Data Bank (EMDB), for storing biomolecular experimental data. The PDB lists detailed information of atomic coordinates, occupancy and Debye-Waller factor (or B-factor) of all the atoms in proteins, DNAs, RNAs and their complexes. In the persistent homology terminology, PDB provides point cloud data for biomolecules. In contrast, the EMDB typically offers volumetric data, or density maps of large biomolecular systems, such as multi-proteins, subcellular structures and organelles from cryo-EM at the molecular level resolution.

6

Molecular structures and their visualization can be generated by using a variety of molecular models, such as the atom and bond model of molecules,[71] the van der Waals surface, the solvent-excluded surface (SES) (also known as molecular surface (MS)) and the solvent-accessible surface, have been proposed.[72,73] These models have been widely applied to the analysis of biomolecular structure, function, and interaction, such as ligand-receptor binding, protein specification, drug design, macromolecular assembly, protein-nucleic acid and protein-protein interactions, and enzymatic mechanism.[74] Nevertheless, they admit geometric singularities, i.e., tips, cusps and self-intersecting surfaces which are troublesome in simulations[75-78] and *ad hoc* in physical foundation because electron density decays gradually at the molecular boundary.[79,80]

Recently, we have introduced the differential geometry theory of surfaces to address the above-mentioned problems in biomolecular geometric modeling by curvature control PDEs,[81] mean curvature flows[82,83] and potential driven geometric flows.[84] The minimization principle was utilized for the biomolecular surface construction. We have further generalized these ideas to incorporate multiscale and multiphysical descriptions of biomolecules.[79,80,85,86] Our approaches have been adopted and/or generalized by many others.[87-90]

Most recently, we have proposed flexibility and rigidity index (FRI) for flexibility analysis and B-factor prediction of proteins and other biomolecules.[91,92] In FRI, protein topological connectivity is measured by rigidity index and flexibility index. In particular, the rigidity index represents protein density profile. Consider a protein with $N$ atoms. Their locations are represented by $\{\mathbf{r}_j | \mathbf{r}_j \in \mathbb{R}^3, j = 1, 2, \cdots, N\}$. We denote $\|\mathbf{r}_i - \mathbf{r}_j\|$ the Euclidean space distance between $i$th atom and the $j$th atom. We define a position ($\mathbf{r}$) dependent rigidity or density function[91,92]

$$\mu(\mathbf{r}) = \sum_{j=1}^{N} w_j(\mathbf{r}_j)\Phi(\|\mathbf{r} - \mathbf{r}_j\|; \sigma_j), \tag{1}$$

where $w_j(\mathbf{r})$ is an atom type dependent weight function and $\sigma_j$ is an atomic type dependent scale parameter which is proportional to the atomic radius. Here $\Phi(\|\mathbf{r} - \mathbf{r}_j\|; \sigma_j)$ is a correlation kernel, which is, in general, a real-valued monotonically decreasing function satisfying

$$\Phi(\|\mathbf{r} - \mathbf{r}_j\|; \sigma_j) = 1 \quad \text{as} \quad \|\mathbf{r} - \mathbf{r}_j\| \to 0 \tag{2}$$
$$\Phi(\|\mathbf{r} - \mathbf{r}_j\|; \sigma_j) = 0 \quad \text{as} \quad \|\mathbf{r} - \mathbf{r}_j\| \to \infty. \tag{3}$$

Although Delta sequences of the positive type discussed in an earlier work[93] are all good choices, generalized exponential functions

$$\Phi(\|\mathbf{r} - \mathbf{r}_j\|; \sigma_j) = e^{-(\|\mathbf{r} - \mathbf{r}_j\|/\sigma_j)^\kappa}, \quad \kappa > 0 \tag{4}$$

and generalized Lorentz functions

$$\Phi(\|\mathbf{r} - \mathbf{r}_j\|; \sigma_i) = \frac{1}{1 + (\|\mathbf{r} - \mathbf{r}_j\|/\sigma_j)^\upsilon}, \quad \upsilon > 0. \tag{5}$$

have been commonly used in our recent work.[91,92] We refer these general classes of volumetric surface definition as rigidity surfaces and density profiles. We use a surface extraction procedure, such as marching cubes, to extract a Lagrangian surface from their volumetric data, or the Eulerian representation of surface density profiles. Obviously, when $\kappa = 2$ in Eq. (4), Eq. (1) gives rise to a representation of Gaussian surfaces, which have many formulations.[70,91,94-98] In general, Gaussian surfaces are quite smooth and free of geometric singularity. The generation of Gaussian surfaces can be very fast and readily available in the Cartesian representation.[96] Other geometric modeling approaches include curvature analysis and symmetry analysis. Mean curvature and Gauss curvature can be estimated in both Lagrangian representations[69] and Eulerian representation.[70] Maximal and minimal principle curvatures can be utilized for drug binding site prediction.[70] Symmetric analysis is frequently employed in biophysical modeling.[96] Utilizing symmetry leads to the reduction of number of genes, which is very common in viral complexes. Many protein complexes, such as microtubules, are very symmetric as well.

Figure 1 illustrates surfaces extracted from density function Eq. (1) with $w_j = 1$, $\kappa = 2$ and $\sigma_j = 0.5$Å. In this work, we use density function (1) as a mathematical model for cryo-EM density maps. A series of
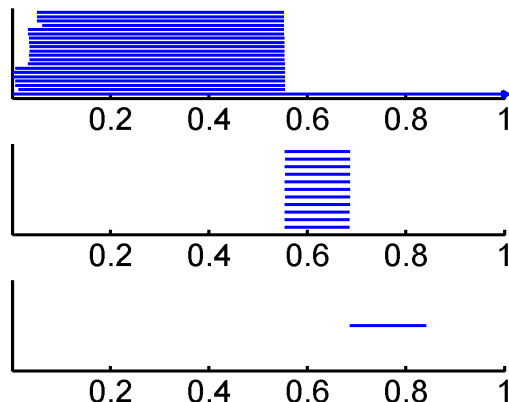
Figure 2: The intrinsic topological fingerprints of fullerene $C_{20}$. The top, middle, and bottom panels are for the barcodes of $\beta_0$, $\beta_1$ and $\beta_2$, respectively. The filtration process is based on the decrement of the density isovalues. The filtration barcodes demonstrate pentagonal structures and central void. It should be notice that $\beta_0$ bars do not emerges simultaneously, which is due to the Cartesian representation of data.

surfaces is plotted in Fig. 1 to demonstrate some typical structures in the filtration procedure for fullerene $C_{20}$ density function. The isovalues for Figs. 1 **a**, **b** and **c** are 0.7, 0.6 and 0.4, respectively. Figure 1 **d** is a wire-frame surface representation of Fig. 1 **c**.

## 2.2 Topological modeling of biomolecules

Persistent homology theory and algorithm can be found in the literature[32, 36–40, 42] as well as our papers.[62, 68] In this work, we focus on electron density maps of macro-protein or protein-protein complexes available as volumetric data deposited in the EMDB. Unlike point cloud data which are commonly studied with simplicial complex using Javaplex,[99] volumetric data are usually analyzed by the discrete Morse theory. For all the density map based volumetric data used in this work, we use the same filtration process that is built based on the decrement of the electron density value. More specifically, cryo-EM density maps or density functions generated from Eq. (1) is used as the filtration parameter. The filtration process goes from the highest density isovalue to the lowest one. After the filtration, the data are analyzed with Perseus.[100] We first consider a benchmark test to illustrate the persistent homology analysis of density function (1).

**Topological persistence of $C_{20}$** The Betti numbers can be directly obtained through identifying the number of connected components, circles or loops, and voids or holes. For example, in Fig. 1 **a**, $\beta_0$ is 20, $\beta_1$ and $\beta_2$ are 0. In Fig. 1 **b**, $\beta_0 = 1$, $\beta_1 = 11$ and $\beta_2 = 0$. In Fig. 1 **c** and **d**, $\beta_0 = 1$, $\beta_1 = 0$ and $\beta_2 = 1$. Attention should be paid to $\beta_1$ for the structure in Fig. 1 **b**. The $\beta_1$ is not exactly equal to the total number ($N_1$) of the circles or loops, instead it equals to $N_1 - 1$. This is due to the reason that the corresponding homology group has a basis with only $N_1 - 1$ independent elements. Roughly speaking, one of the circles can be expressed as the "linear combination" of the rest.

Barcodes provide a systematic representation of topological persistence.[32] Figure 2 shows the intrinsic topological patterns of fullerene $C_{20}$. Just the same as we counted above, there are 11 $\beta_1$ bars and only one $\beta_1$ bar. It should be noticed that $\beta_0$ bars do not emerge simultaneously, which is due to the discretization effect, namely, the density function is discretized with only a finite resolution.

**Topological persistence of EMD 1776** After demonstrating the persistent homology analysis for the density function (1) of a known structure ($C_{20}$), we further consider realistic cryo-EM data, EMD 1776, which is for eye lens chaperone $\alpha$-crystallin assemblies. Figure 3 depicts the surfaces extracted from different isovalues for EMD 1776. The isovalues for Figs. 3 **a**, **b**, **c** and **d** are 0.150, 0.100, 0.081 and 0.050 respectively. Similarly, Betti numbers can be directly obtained through counting the numbers of connected components, circles, and voids. In **a**, $\beta_0$ is 12, $\beta_1$ and $\beta_2$ are 0. In Fig. 3 **b**, $\beta_0 = 4$, $\beta_1 = 4$ and $\beta_2 = 0$. In Fig. 3 **c**, $\beta_0 = 1$, $\beta_1 = 13$
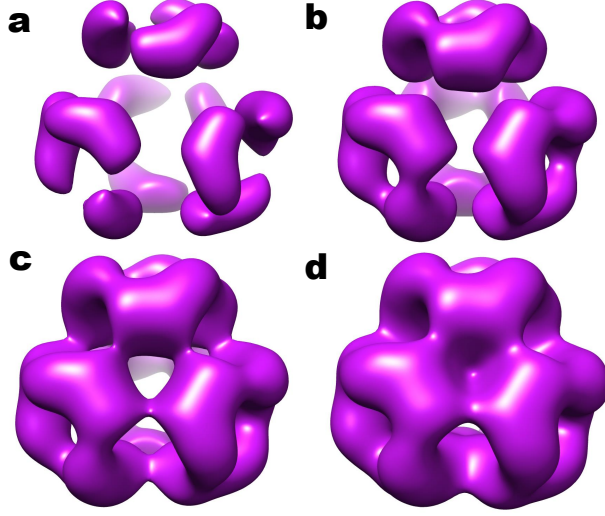
Figure 3: Surfaces extracted with different isovalues for EMD 1776. The isovalues for subfigures **a**, **b**, **c** and **d** are 0.150, 0.100, 0.081 and 0.050 respectively. The Betti numbers can be directly obtained through identifying the number of connected components, circles or loops, and voids or holes. In **a**, $\beta_0$ is 12, $\beta_1$ and $\beta_2$ are 0; In **b**, $\beta_0 = 4$, $\beta_1 = 4$ and $\beta_2 = 0$; In **c**, $\beta_0 = 1$, $\beta_1 = 13$ and $\beta_2 = 0$; In **d**, $\beta_0 = 9$, $\beta_1 = 13$ and $\beta_2 = 0$. In **c** and **d**, the $\beta_1$ is one count fewer than the total number ($N_1$) of the circles or loops, because the corresponding homology group has a basis with only $N_1 - 1$ elements.

and $\beta_2 = 0$. Also in Fig. 3 **d**, $\beta_0 = 9$, $\beta_1 = 13$ and $\beta_2 = 0$. As discussed above, in Figs. 3 **c** and **d**, the $\beta_1$ value is $N_1 - 1$, rather than the total number ($N_1$) of the circles, due to $N_1 - 1$ independent elements in the corresponding homology group. The barcode representation is demonstrated in Fig. 4, which is consistent with our analysis.

It should be noticed that we only consider the regions with density values larger than 0.03, namely, the filtration goes from the largest value (0.28) to a threshold value (0.03). For density values smaller than 0.03, data suffer from lower SNRs as discussed in Section 3.2. Denoising techniques are indispensable for extracting more information from low isovalues. In the next section, we apply persistent homology to noise reduction and topological feature identification.

## 3 Persistent homology based noise reduction

In this section, we present persistent homology based cryo-EM data noise reduction, which is a crucial process in cryo-EM analysis. Protein data in EMDB are mostly obtained by cryo-EM. As discussed earlier, the cryo-EM data suffer from low resolution and low SNR. Therefore a denoising process is always a necessity before carrying out geometric modeling and/or topological analysis. We focus on noise reduction based on persistent homology analysis. Specifically, we employ persistent homology to discriminate signal from noise and utilize this information for denoising thresholding. In our benchmark test, we assume a known object is contaminated with Gaussian noise. Geometric flows proposed in our earlier work[84, 101] are used for noise reduction of realistic cryo-EM data.

### 3.1 Topological fingerprints of Gaussian noise

We first analyze the topological fingerprint or topological signature of noise. We use Gaussian noise as an example for the present study. Other noise can be analyzed in a similar manner. The Gaussian white noise is generated by randomly selecting values from a normal distribution,

$$n(t) = \frac{A_n}{\sqrt{2\pi}\sigma_n} e^{-\frac{(t-\mu_n)^2}{2\sigma_n^2}}, \tag{6}$$

where $A_n$, $\mu_n$ and $\sigma_n$ are the amplitude, mean value and standard deviation of the noise, respectively. We denote $\mu_s$ as the mean value of signal. Then the degree of noise contamination can be described by SNR,
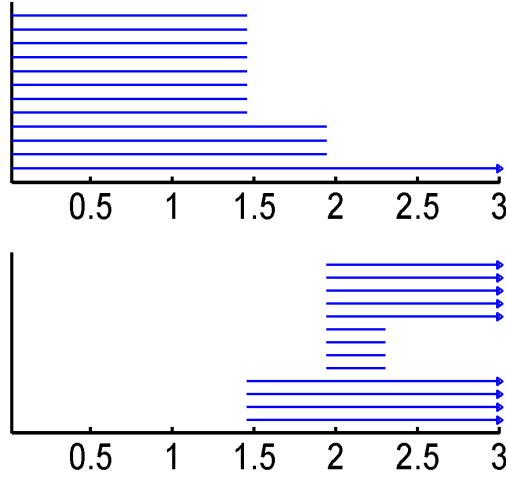
Figure 4: The intrinsic topological patterns of the EMD 1776 structure. The top and bottom panels are for the barcodes of $\beta_0$ and $\beta_1$, respectively. The filtration process is based on the decrement of density isovalues. The filtration goes from the largest isovalue (0.28) to isovalue threshold (0.03).

$SNR = \mu_s/\sigma_n$. Please note that the present definition of SNR is by no means unique. Based on the physical properties of the signal, SNR can be defined in terms of average power, amplitude, variance of the signal, and so on. In our discussion, the signal information represented in volumetric data can be easily analyzed. We generate the noise data with specified SNR by adding suitable amplitude of Gaussian white noise. Stated differently, the noise contaminated data are generated by adding different levels of Gaussian white noise to the original data, i.e., the density function or density map. We then investigate their corresponding persistent barcode patterns.

The density function described in Eq. (1) with the generalized exponential kernel shown in Eq. (4) is used to simulate the density of fullerene $C_{20}$. We choose $\kappa = 1.0$ and $\sigma = 0.5$Å in our study. Figure 5 demonstrates the barcode representation for contaminated fullerene $C_{20}$ data with different SNRs. The SNRs for Figs. 5 **a**, **b**, **c** and **d** are 0.1, 1.0, 10.0 and 100.0, respectively. It can be seen from Figs. 5 **a** and **b**, that when SNR is low, i.e., SNR= 0.1 or 1.0, fullerene atoms are invisible. When the SNR is increased in Figs. 5 **c** and **d**, the molecular intrinsic patterns begin to emerge.

The topological signature of the above four cases can be analyzed as shown in Figs. 5 **e**, **f**, **g** and **h**. First, we note that all of three topological invariants, namely $\beta_0$, $\beta_1$ and $\beta_2$, are very sensitive to and essentially dominated by noise. Additionally, noise gives rise to continuous bands (or stripes) of bars in all the three topological invariants. Moreover, noise provides similarly numbers of bars in $\beta_0$, $\beta_1$ and $\beta_2$ panels. Moreover, contradicting to the common belief that noise only contributes to short lived bars, the noise bars can be very long. The average band length of noise bars proportions to noise magnitude. In fact, in our recent work,[62] we regard all bars, including short lived bars, of a protein as molecular topological fingerprints and as being of equal importance. Furthermore, at low SNRs (i. e., Figs. 5 **e** and **f**), the bars of three invariants maintain Gaussian-like distributions with respect to the isovalue filtration (i. e., the $x$-axis) as shown in Fig. 6. However, at relatively high SNRs, bars do not have the Gaussian-like distributions due to a relatively high $C_{20}$ density, see Fig. 6.

To further analyze the topological signature of noise, we consider a protein structure made of only beta sheets obtained from 2GR8. Again its total density distribution is approximated by using the density function given in Eq. (1) realized by using the generalized exponential kernel of Eq. (4) with $\kappa = 1.0$ and $\sigma = 2.0$Å. Figure 7 shows our results. The SNRs for Figs. 7 **a**, **b**, **c** and **d** are 0.1, 0.5, 1.0 and 10.0, respectively. The corresponding barcodes are given in Figs. 7 **e**, **f**, **g** and **h**. The topological signature of noise in Fig. 7 is quite similar to that in Fig. 6, Essentially, Gaussian noise induces continuous bands (or stripes) of bars, whose
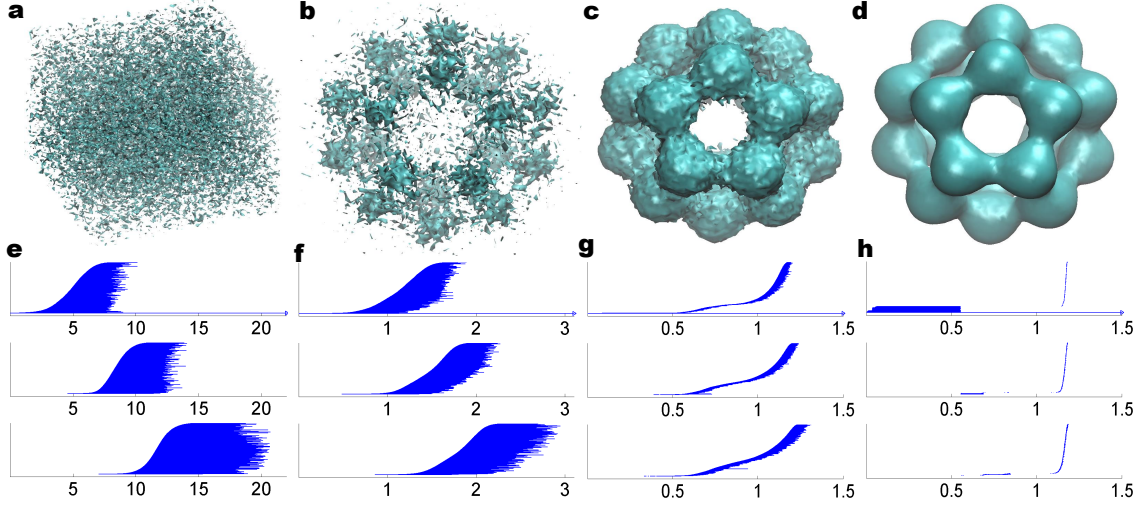
Figure 5: The barcode representation for contaminated fullerene $C_{20}$ data with different SNRs. The SNRs for **a**, **b**, **c** and **d** are 0.1, 1.0, 10.0 and 100.0, respectively, and isovalues used for their visualization are 4.00, 1.00, 0.60, and 0.60, respectively. The corresponding barcodes are given in **e**, **f**, **g** and **h**, respectively. The top, middle, and bottom panels are for the barcodes of $\beta_0$, $\beta_1$ and $\beta_2$, respectively. It can be seen from the barcodes that when SNR is low, i.e., SNR = 0.1 or 1.0, fullerene molecule is invisible. At a low noise level, molecular pattern emerges. More importantly, the persistence of the psudo-topological structure is directly related to the SNR. In the barcode representation, noise tends to induce a continuous band (or stripe) of bars, of which the width or relative persistent length is determined by the magnitude of the noise.
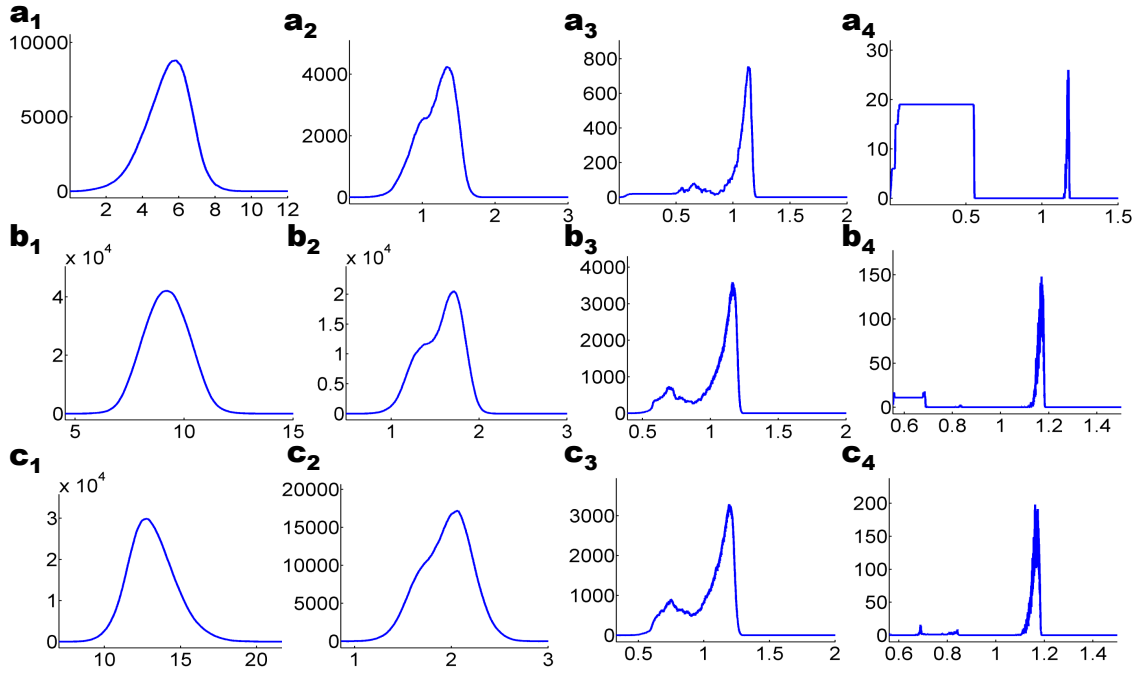


Figure 6: The histograms of topological invariants over the filtration process for contaminated fullerene $C_{20}$ data with different SNRs. The $\mathbf{a}_i$, $\mathbf{b}_i$ and $\mathbf{c}_i$ rows are the counts of $\beta_0$, $\beta_1$ and $\beta_2$, respectively, where subscripts $i = 1, 2, 3$ and $4$ correspond to SNRs 0.1, 1.0, 10.0 and 100.0, respectively. It can be seen that as the SNR increases, barcodes for noise and signal gradually separate from each other. Since the Gaussian noise is used, the noise parts typically assume Gaussian distributions in shape.
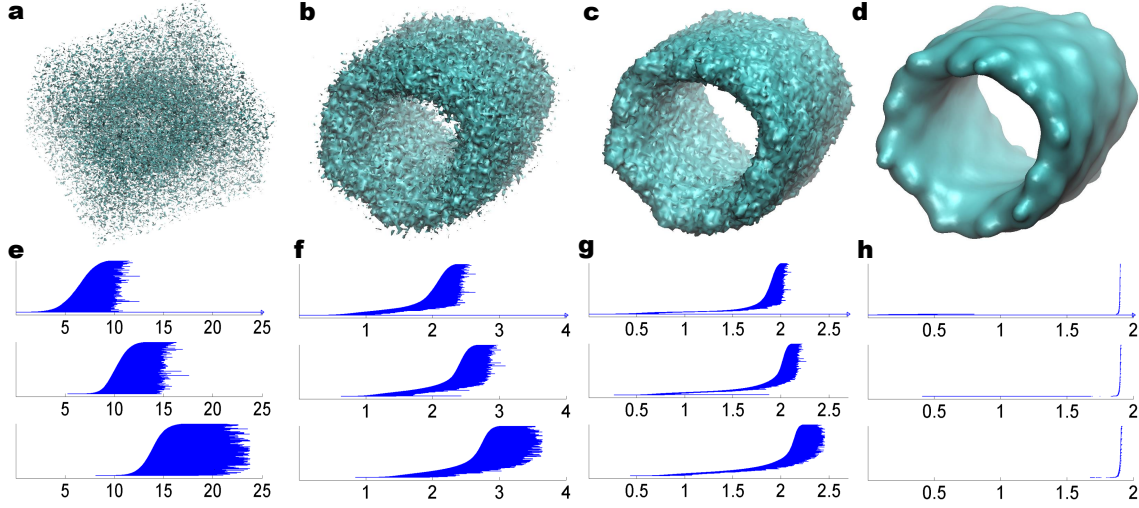
11

Figure 7: The barcodes representation for protein 2GR8 beta segment with different SNRs. The SNRs for **a**, **b**, **c** and **d** are 0.1, 0.5, 1.0 and 10.0, respectively, and isovalues used for their visualization are 5.00, 1.00, 1.00, and 1.00, respectively. The corresponding barcodes are presented in **e**, **f**, **g** and **h**, respectively. The top, middle, and bottom panels are for the barcodes of $\beta_0$, $\beta_1$ and $\beta_2$, respectively. It can be seen from the barcodes that when the SNR is large, i.e., SNR=0.1 or 0.5, the original topological properties are blurred. When the noise effect dwindles, the intrinsic patterns begin to emerge. More importantly, the persistence of the psudo-topological structure is directly related to the SNR. In the barcode representation, noise tends to induce a continuous band (or stripe) of bars, of which the width or relative persistent length is determined by the magnitude of the noise.

width or relative persistent length is determined by noise intensity.

As demonstrated in above examples, barcodes present a topological description of Gaussian noise. The related band width of the noise can be used to assess noise magnitude. This qualitative description can be further used as a guidance for topological noise reduction and topological fingerprint identification.

## 3.2 Topological denoising

**Geometric flows** Geometric PDEs offer an efficient approach for noise reduction. High order geometric PDEs were first introduced for edge-preserving image restoration in 1999 and have a general form[101]

$$\frac{\partial u(\mathbf{r}, t)}{\partial t} = -\sum_q \nabla \cdot \mathbf{j}_q + e(u(\mathbf{r}, t), |\nabla u(\mathbf{r}, t)|, t), \quad q = 0, 1, 2, \cdots \tag{7}$$

where the nonlinear hyperflux term $\mathbf{j}_q$ is given by

$$\mathbf{j}_q = -d_q(u(\mathbf{r}, t), |\nabla u(\mathbf{r}, t)|, t) \nabla \nabla^{2q} u(\mathbf{r}, t), \quad q = 0, 1, 2, \cdots \tag{8}$$

where $\mathbf{r} \in \mathbb{R}^n$, $\nabla = \frac{\partial}{\partial \mathbf{r}}$, $u(\mathbf{r}, t)$ is the processed image function, $d_q(u(\mathbf{r}, t), |\nabla u(\mathbf{r}, t)|, t)$ are edge sensitive diffusion coefficients and $e(u(\mathbf{r}, t), |\nabla u(\mathbf{r}, t)|, t)$ is a nonlinear operator. The original noise data $X(\mathbf{r})$ is used as the initial input $u(\mathbf{r}, 0) = X(\mathbf{r})$. The hyper-diffusion coefficients $d_q(u, |\nabla u|, t)$ in Eq. (8) can also be chosen as the Gaussian form

$$d_q(u(\mathbf{r}, t), |\nabla u(\mathbf{r}, t)|, t) = d_{q0} \exp\left[-\frac{|\nabla u|^2}{2\sigma_q^2}\right], \tag{9}$$

where $d_{q0}$ is chosen as a constant with value depended on the noise level, and $\sigma_0$ and $\sigma_1$ are local statistical variance of $u$ and $\nabla u$

$$\sigma_q^2(\mathbf{r}) = \overline{|\nabla^q u - \overline{\nabla^q u}|^2} \quad (q = 0, 1). \tag{10}$$
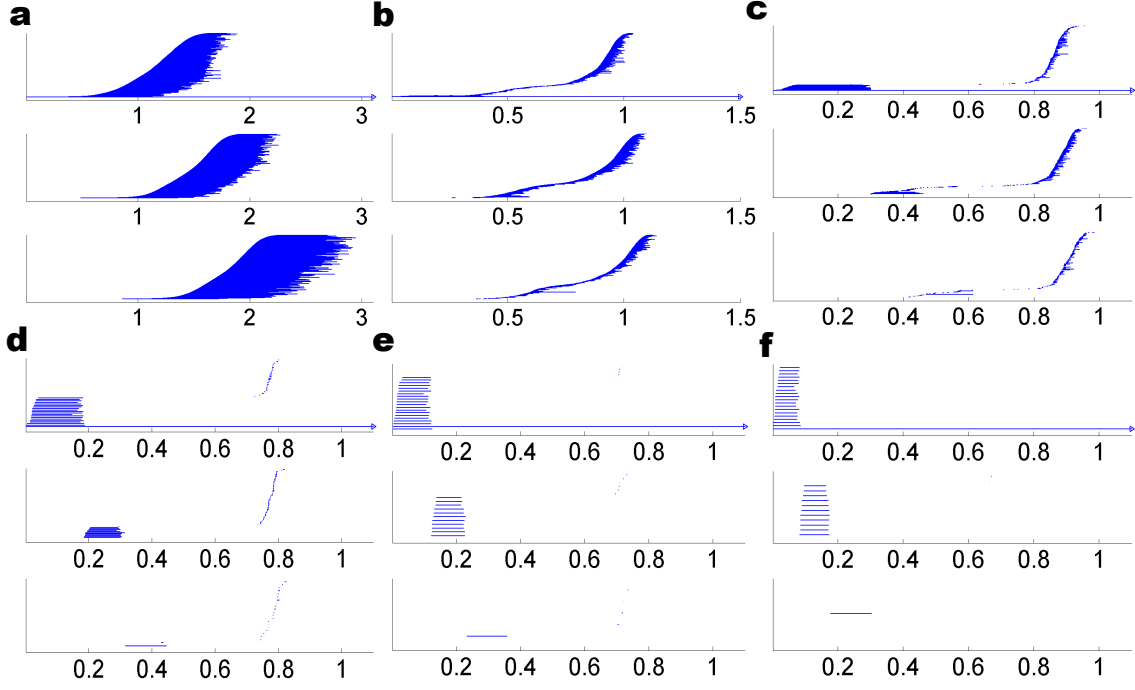
Figure 8: The barcodes representation for denoising contaminated fullerene $C_{20}$ with SNR 1.0. The barcodes for fullerene $C_{20}$ with SNR 1.0 is demonstrated in **a**. The denoising steps for **b**, **c**, **d**, **e** and **f** are 20, 40, 60, 80 and 100, respectively. The noise induced topological invariants have been gradually weakened and finally eradicated. Compared with the original noise-polluted barcodes in **a**, the noise effect has been enormously scaled down after only 20 steps of denoising as indicated in **b**. In **c**, there is a clear separation between the intrinsic topological features of fullerene $C_{20}$ and noise induced topological invariants. From **d** to **f**, the noise effect is further reduced, and we are able to identify a persistent barcode pattern, which is an indication of the intrinsic topological invariants of fullerene $C_{20}$. It should also be noticed that the denoising process fades noise intensity, therefore noise induced topological patterns gradually move to the right of filtration parameter and finally disappear.

ants have been gradually weakened and finally eradicated. Compared with the original noise-polluted barcodes in **a**, the noise effect has been enormously scaled down after only 20 steps of denoising as indicated in **b**. In **c**, there is a clear separation between the intrinsic topological features of fullerene $C_{20}$ and noise induced topological invariants. From **d** to **f**, the noise effect is further reduced, and we are able to identify a persistent barcode pattern, which is an indication of the intrinsic topological invariants of fullerene $C_{20}$. It should also be noticed that the denoising process fades noise intensity, therefore noise induced topological patterns gradually move to the right of filtration parameter and finally disappear.

Here the notation $\overline{Y(\mathbf{r})}$ represents the local average of $Y(\mathbf{r})$ centered at position $\mathbf{r}$.

High order geometric PDEs have many practical applications.[101–103] They have been specifically modified for molecular surface formation and evolution[84] as,

$$\frac{\partial S}{\partial t} = (-1)^q \sqrt{g(|\nabla \nabla^{2q} S|)} \nabla \cdot \left( \frac{\nabla(\nabla^{2q} S)}{\sqrt{g(|\nabla \nabla^{2q} S|)}} \right) + P(S, |\nabla S|), \tag{11}$$

where $S$ is the hypersurface function, $g(|\nabla \nabla^{2q} S|) = 1 + |\nabla \nabla^{2q} S|^2$ is the generalized Gram determinant and $P$ is a generalized potential term. When $q = 0$ and $P = 0$, a Laplace-Beltrami equation is obtained,[83]

$$\frac{\partial S}{\partial t} = |\nabla S| \nabla \cdot \left( \frac{\nabla S}{|\nabla S|} \right). \tag{12}$$

We employ this Laplace-Beltrami equation for the noise reduction in this paper.
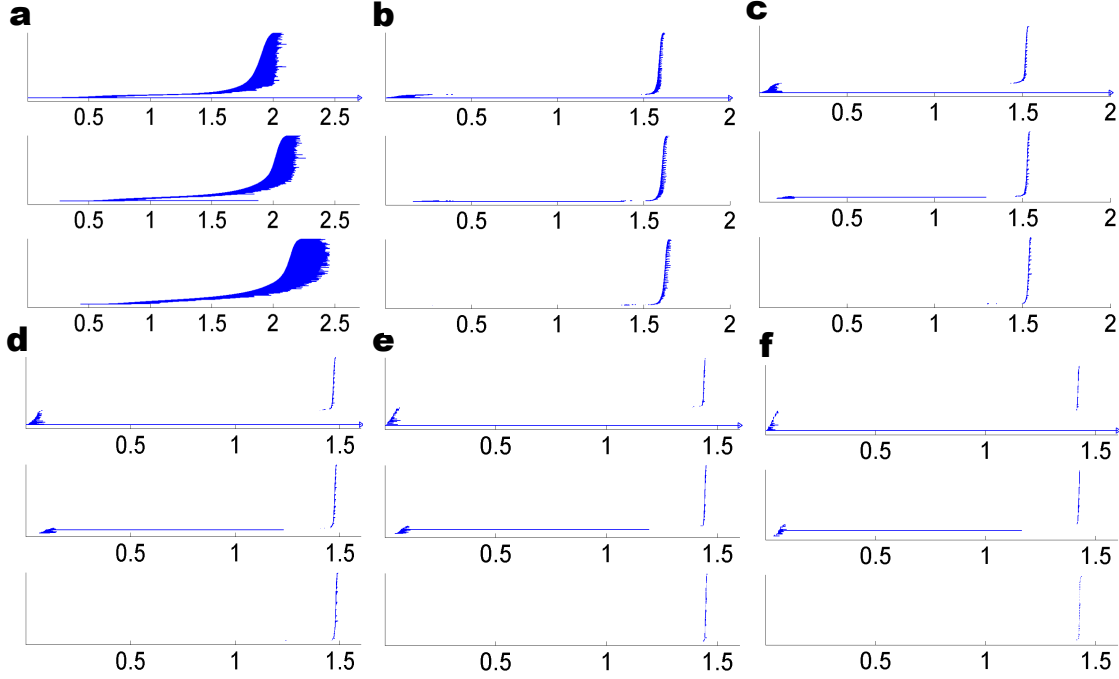
Figure 9: The barcodes representation for the noise reduction of contaminated protein 2GR8 beta segment with SNR 1.0. The barcodes for 2GR8 with SNR 1.0 is demonstrated in **a**. The denoising steps for **b**, **c**, **d**, **e** and **f** are 10, 20, 30, 40 and 50, respectively. In this case, the noise induced topological invariants have been gradually weakened but not eradicated. Compared with the original noise-polluted barcodes in **a**, the noise effect has been enormously scaled down after 10 steps of denoising as indicated in **b**. From **c** to **f**, there is a clear separation between the intrinsic topological features of protein segment and noise induced topological invariants. The noise effect is continuously reduced, and we are able to identify a persistent barcode pattern, which is an indication of the intrinsic topological invariants of the protein segment. However, unlike the fullerene $C_{20}$, further denoising of 2GR8 data will remove both the noise and intrinsic structure related topological information.

**Topological fingerprint identification**  Computationally, the finite different method is used to discretize Eq. (12). Suitable time interval $\Delta t$, and grid spacing $h$ are needed. For cryo-EM data, its voxel spacing is related to the data resolution and varies greatly. For example, The voxel spacings of EMD1776, EMD1229 and EMD5729 are 1.69 Å, 4.00 Å, and 4.16 Å, respectively. In our simulated fullerene $C_{20}$ , $C_{60}$ and protein 2GR8 examples, the voxel spacings are 0.06 Å, 0.08 Å, and 0.40 Å, respectively. To avoid confusion and control the noise reduction process systematically, we simply ignore the voxel spacing and use the unified parameters $\Delta t = 1.0E - 5$ and $h = 1.0E - 2$. The intensity of denoising is then described by the number of iteration steps.

The noise reduction effectiveness is commonly validated by a visual comparison with the original results. Quantitative assessment usually proves to be difficult, as noise and signal or image information can be tightly entangled. In this section, we propose a topological method for monitoring the evolution of relative behaviors of noise and signal during the denoising process. More specifically, persistent barcodes from a series of denoising data are compared. The noise signature and topological fingerprint in these barcodes are carefully studied. The present emphasis is on the topological fingerprint identification.

During the denoising process, topological features of both signal and noise are constantly evolving. To quantitatively analyze their behaviors, three cases with Gaussian white noise and a case with Cryo-EM data are considered. The first case is a noise contaminated fullerene $C_{20}$ with SNR 1.0. The persistent barcode results are demonstrated in Fig. 8. We vary the number of denoising steps in our study. For Figs. 8**b**, **c**, **d**, **e** and **f**, the numbers of iterations are 20, 40, 60, 80 and 100, respectively. The noise induced topological persistence has been gradually weakened and finally cleaned up. Compared with the original noise-polluted barcodes in Fig. 8**a**, the noise effect has been enormously scaled down after only 10 steps of denoising as
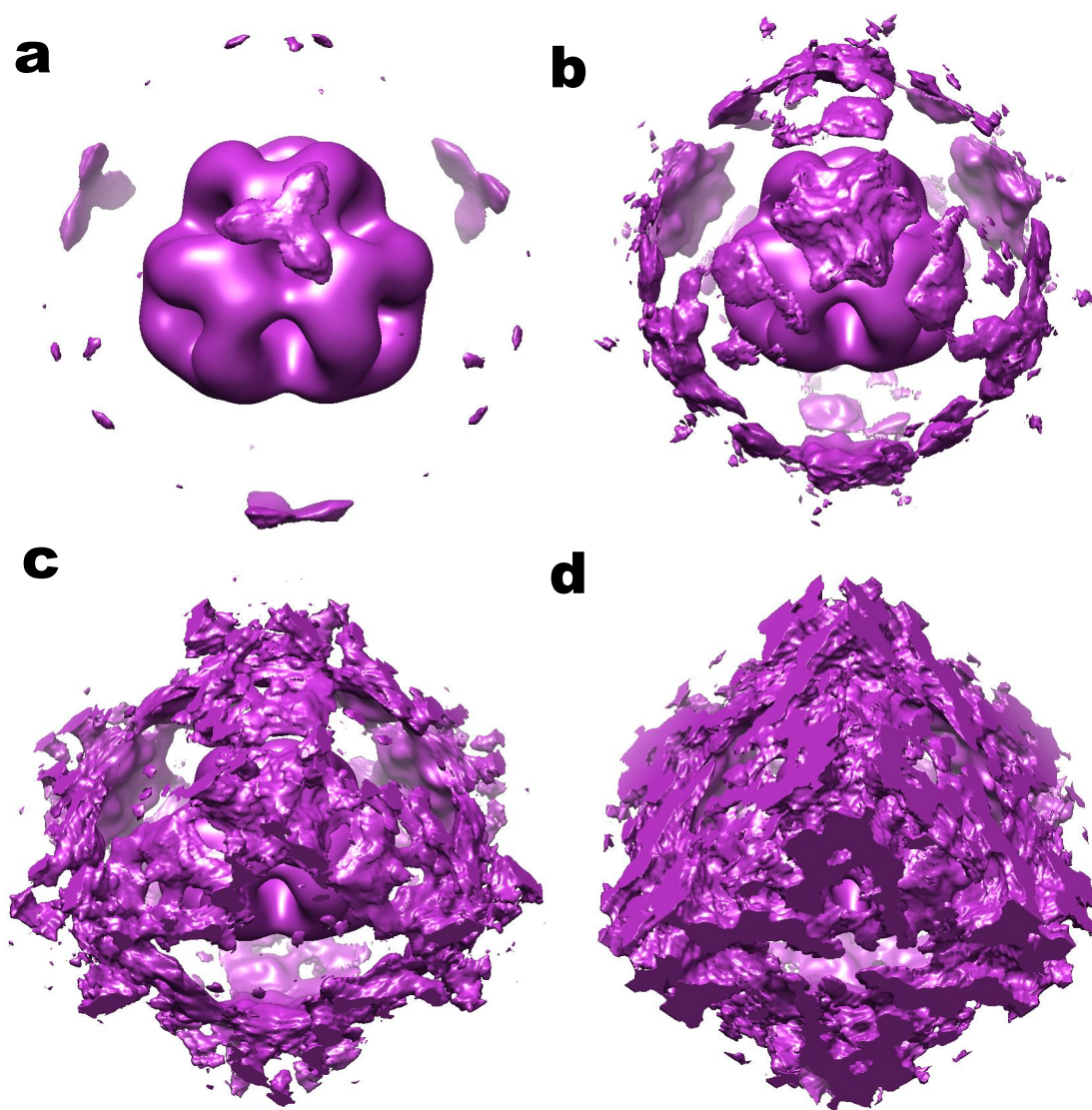
Figure 10: The original noise in EMD 1776 data. The isovalues for **a**, **b**, **c** and **d** are 0.020, 0.010, 0.005, and 0.000, respectively.
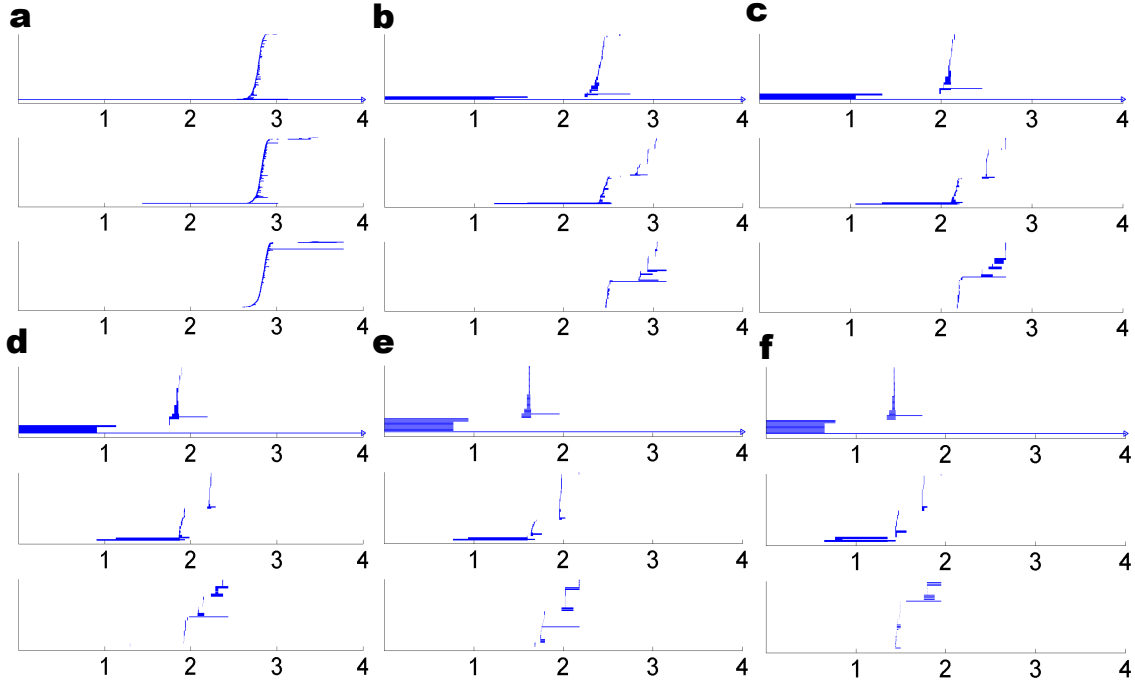
Figure 11: The barcodes representation for the noise removal of contaminated EMD 1776 with SNR 1.0. The denoising steps for **a**, **b**, **c**, **d**, **e** and **f** are 10, 40, 80, 120, 160 and 200, respectively. In this case, the noise induced topological invariants have been gradually weakened but not eradicated. Compared with the original noise-polluted barcodes in Fig. 7 **b**, the noise effect has been enormously reduced after 10 steps of denoising as indicated in **a**. From **b** to **f**, there is a clear separation between the intrinsic topological features of protein segment and noise induced topological invariants. The noise effect continuously wanes, and we are able to identify a persistent barcode pattern, which is an indication of the intrinsic topological invariants of the protein segment. However, unlike the fullerene $C_{20}$, further denoising will remove both the noise and intrinsic structure related topological information.

indicated in Fig. 8**b**. In Fig. 8**c**, there is a clear separation between the intrinsic topological persistence of fullerene $C_{20}$ and noise induced topological persistence. From Fig. 8**d** to Fig. 8**f**, the noise effect is further weakened, and thus we are able to identify a consistent barcode pattern, which is an indication of the intrinsic topological invariants of fullerene $C_{20}$. It should also be noticed that since noise intensity diminishes during the denoising process, noise induced topological patterns gradually shift to the right of the filtration parameter and eventually disappear.

The second case is a noise contaminated protein segment from 2GR8 with SNR 1.0. Its topological behavior under the denoising process is illustrated in Fig. 9. In this case, the noise induced topological invariants have been gradually weakened but not eradicated. Compared with the original noise-polluted barcodes in Fig. 9**a**, the noise persistence has been enormously reduced after 10 steps of denoising as indicated in Fig. 9**b**. From Fig. 9**b** to Fig. 9**f**, there is a clear separation between the intrinsic topological invariants of protein segment and noise induced topological invariants. As the noise effect is continuously weakened, we are able to identify a persistent barcode pattern, which is an indication of the intrinsic topological invariants of the protein segment. However, unlike the situation for fullerene $C_{20}$, further denoising will remove both the noise and intrinsic structure related topological information.

Finally, we consider a more realistic example, i.e., EMD 1776, obtained from cryo-EM. For EMD 1776, when the isovalue goes to around 0.020, noise begins to emerge. Figure 10 depicts noise in EMD 1776 data. The isovalues for Figs. 10 **a**, **b**, **c**, and **d** are 0.020, 0.010, 0.005 and 0.000, respectively. The geometric flow based denoising method is employed. Persistent homology results are demonstrated in Fig. 11. The numbers of denoising steps in Figs. 11 **a**, **b**, **c**, **d**, **e** and **f** are 10, 40, 80, 120, 160 and 200, respectively. In this case, the noise induced topological invariants have been gradually weakened but have not been cleaned up.
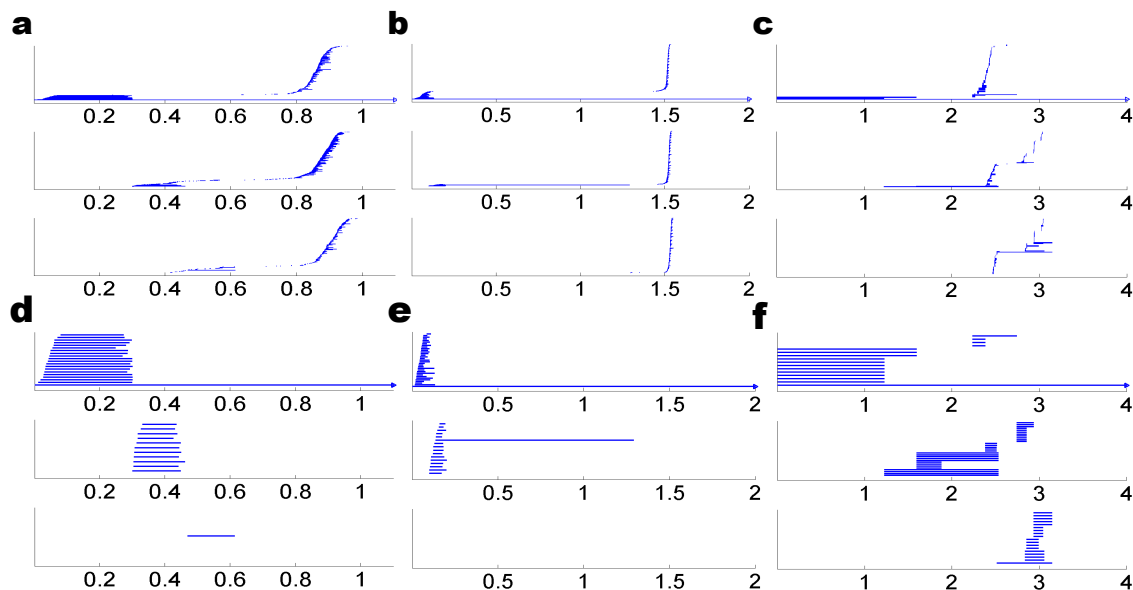
Figure 12: Retrieve barcode patterns through noise thresholding for three systems. Subfigure **a** is the barcode for denoising fullerene $C_{20}$ data with 20 iterations as in Fig. 8 **b**. Subfigure **b** is the barcode for denoising protein 2GR8 beta segment data with 20 iterations as in Fig. 9 **c**. Subfigure **c** is the barcode for denoising EMD 1776 data with 40 iterations as in Fig. 11 **b**. Barcodes in **d**, **e** and **f** are all obtained respectively from **a**, **b** and **c** by setting noise threshold as 0.1, i.e., removing all barcodes with their lengths shorter than 0.1.

From the above analysis, some common features can be unveiled. First, signal related topological features tend to be buried near the left end of the filtration parameter during the denoising process. Second, topological features corresponding to signal and noise begin to separate as the denoising procedure advances. Third, intrinsic topological invariants associated with the "signal" of the data are essentially preserved over the denoising process. These features can be used to guide the evaluation and thresholding of the denoising process.

When the persistent features in the barcode representation is identified, one can retrieve the intrinsic barcodes due to the signal by simply setting up a noise threshold and removing all barcodes with length less than it.[104] Figure 12 demonstrates this technique. Figure 12**a** shows the barcodes after 20 iterations of the noisy fullerene $C_{20}$ data as given in Fig. 8 (**b**). Figure 12 **b** illustrates the barcodes after 20 iterations of the noisy protein 2GR8 beta segment data as given in Fig. 9 (**c**). Figure 12 **c** depicts the barcodes after 40 iterations of the noisy EMD1776 data as in presented Fig. 11 **b**. Figures 12 **d**, **e** and **f** display the barcodes of the above three cases with a noise threshold of 0.1, i.e., removing all barcodes with their lengths shorter than 0.1.

Another importance aspect is that for cryo-EM data, we just need to consider voxels with isovalue larger than a certain threshold. For instance, in EMD 1776 data, noise begins to emerge when the isovalue goes down to about 0.020 as indicated in Fig. 10. As the isovalue decreases, more noise emerges. More importantly, in most cryo-EM data, isovalue can even go below 0.0. These special voxels, as far as we know, do not really represent the desirable biomolecular structure or more specifically do not directly reflect the desirable biomolecular structure. Usually, a recommended isovalue is specified for each cryo-EM data. The meaningful structure information can be only derived from voxels with value near this specified isovalue. From our persistent analysis, we believe that all the voxels with isovalue larger than certain threshold can be related to their inner structure properties. In order not to overlook certain potential structure pattern, in our persistent analysis, we just ignore all voxels with isovalue smaller than 0.0. This is usually done by assigning all negative isovalues to 0.0.
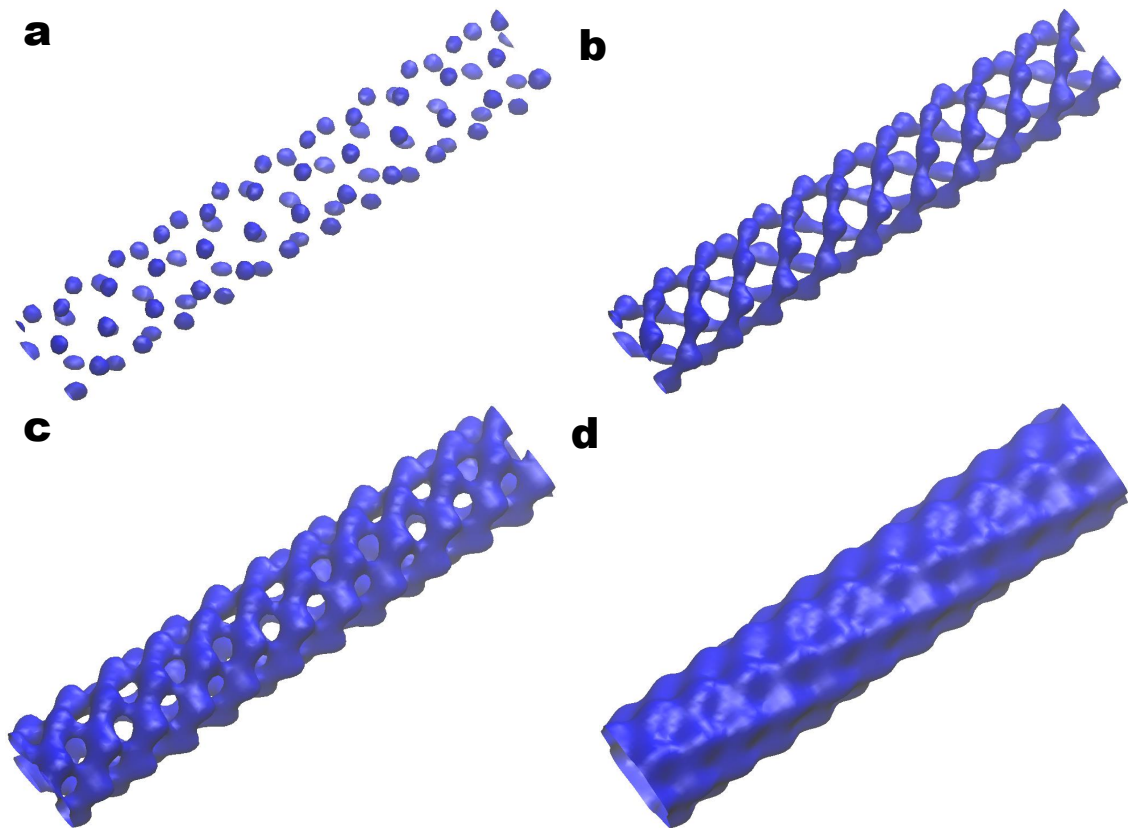
Figure 13: Isosurfaces of EMD 5729 data after 20 steps of noise reduction. Images in **a**, **b**, **c** , and **d** are extracted from isovalues 0.35, 0.30, 0.25, and 0.10, respectively.

### 3.3   Case study: EMD 5729

Finally, we consider EMD 5729, a supramolecular filamentous complex,[105] to demonstrate the application of topological denoising in cryo-EM data analysis. We first process the EMD 5729 data with 20 steps of noise reduction using our geometric flow method. The resulting data are illustrated in Fig. 13 with four different isovalues.

Figure 14 depicts barcodes computed for EMD 5729. The $\beta_0$ pattern in Fig. 14**a** shows a large number of bars of similar lengths, which indicates only one type of protein monomers. Four relatively long persistent bars in the highlighted circle indicate that there are four major pieces in the structure. The $\beta_1$ panel appears to be heavily contaminated by noise, which suggests the necessity for a denoising process. The denoised $\beta_0$ topological persistent patterns in Figs. 14 **b**, **c** and **d** confirm that only one type of bars can be found, which means there is only one type of polymer monomers. Additionally, four relatively long $\beta_0$ bars confirm there are four polymer chains. The $\beta_1$ bars in Figs. 14 **b**, **c** and **d** are relatively consistent over the denoising process and have very similar lengths, which suggest these polymers are evenly distributed and form certain tunnel type of global structures with high symmetry. The long $\beta_1$ bar in Fig. 14 **d** indicates a large cylinder structure. Indeed, Fig. 13 shows that protein monomers form four helix polymers and then bind together to result in a hollow cylinder structure.[105]

### 4   Persistent homology analysis of microtubule

In this section, persistent homology and topological denoising approaches are applied to the microtubule cryo-EM data analysis.
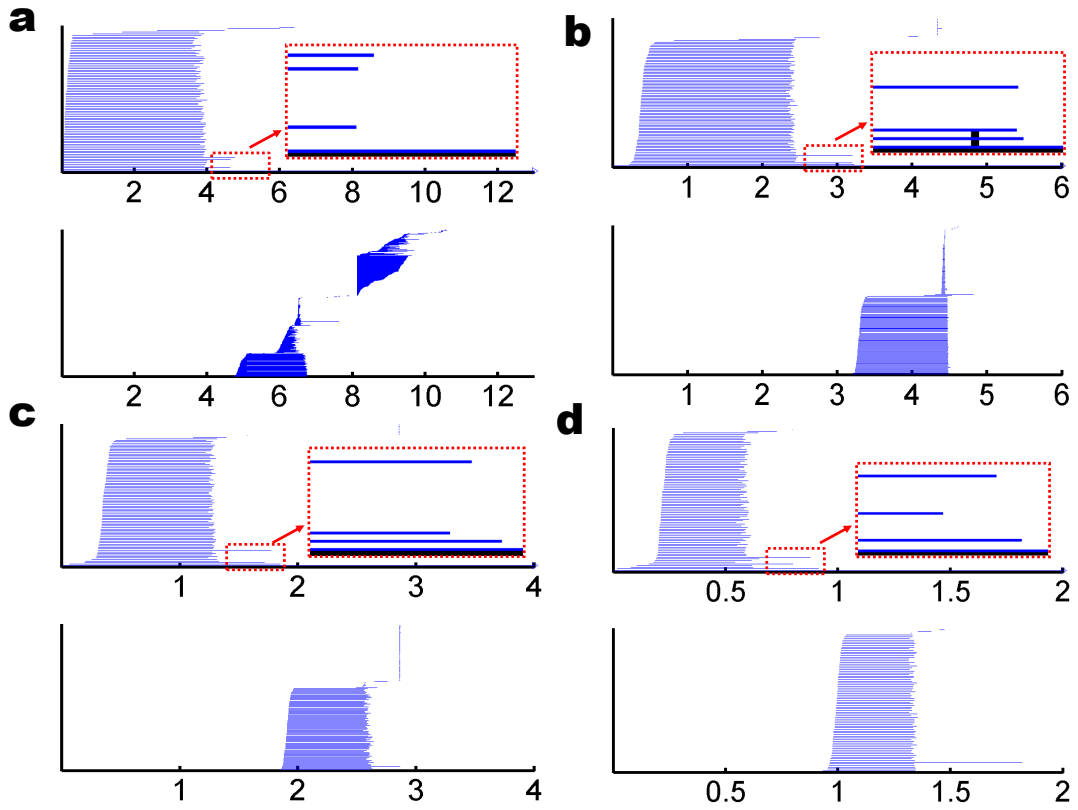
Figure 14: The $\beta_0$ and $\beta_1$ barcodes for EMD 5729 in the noise reduction process. The barcodes for the original data are shown in **a**. Here **b**, **c** , and **d** are barcodes after 10, 20 and 30 steps of noise reduction, respectively. Four relatively long $\beta_0$ bars in the highlighted circles indicate that there are four polymer strands in the structure.
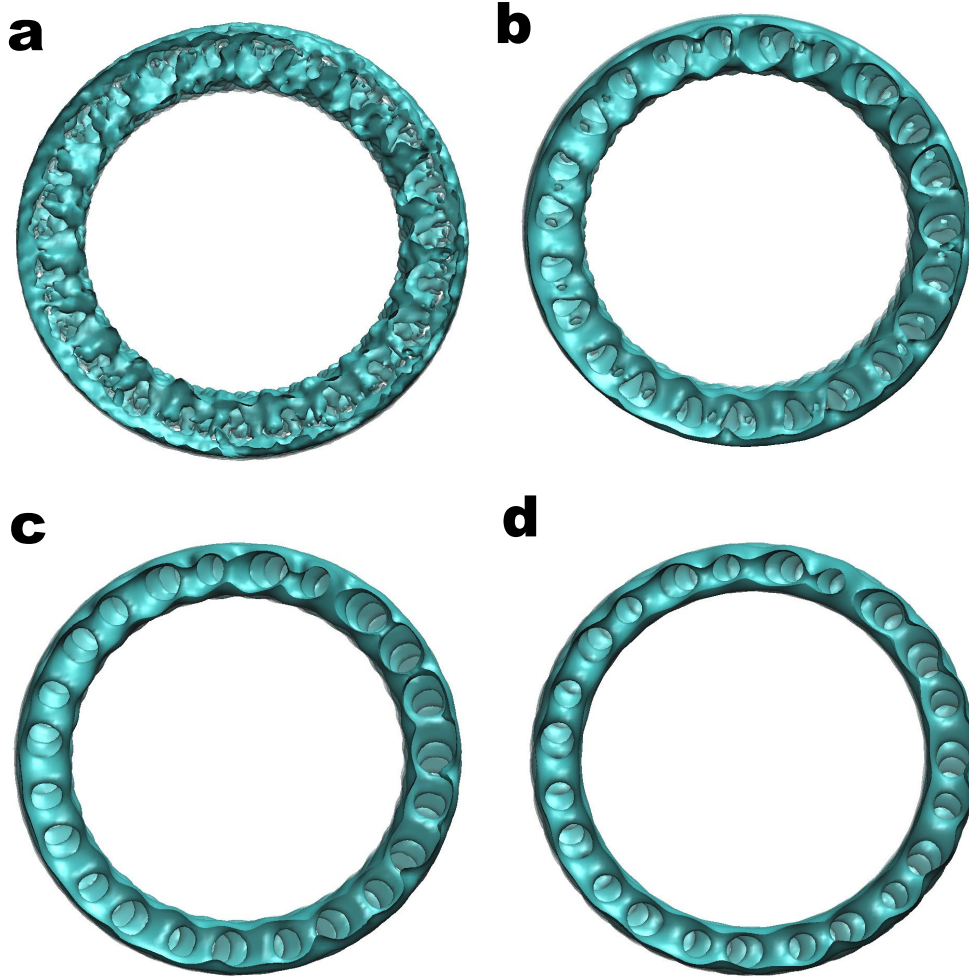
Figure 15: EMD 1129 data preprocessing. **a**: Surface extracted from original data with isovalue 16. **b**, **c** and **d** are surfaces extracted from denoising data with 10, 20, and 40 iterations, respectively.

## 4.1 Microtubule structure EMD-1129

Microtubule is a cytoskeleton component of eukaryotic cells. It plays important roles in maintaining the structure or shape of the cell, supporting intracellular transport and facilitating cell division (mitosis and meiosis).[106] Microtubule has a long hollow cylinder structure made up of polymerized $\alpha$- and $\beta$- tubulin dimers. These hetero-dimers bind head to tail into protofilaments, which further combine with each other in a parallel manner. The hollow cylinder structure of microtubule is finalized by attach about 13 protofilaments with each other side by side. Although the crystal structures of $\alpha-$ and $\beta-$ tubulin are available, experimental microtubule structure data are usually in low resolution and inadequate to separate between neither $\alpha-$ tubulin and $\beta-$ tubulin, nor intra-dimer interface and inter-dimer interface. Recently, a microtubule intermediate structure data (EMD-1129) with a 12 Angstrom resolution is obtained.[107] Based on these data, we demonstrate how can we make use of our persistent homology and topological denoising methods to aid the modeling of cryo-EM structures.

## 4.2 Coarse-grained models for microtubule

For cryo-EM data of low resolution or intermediate resolution, it is well-known that atomic scale models are unreliable. As such, coarse-grained models in terms of residues or even proteins can be useful. In this work,
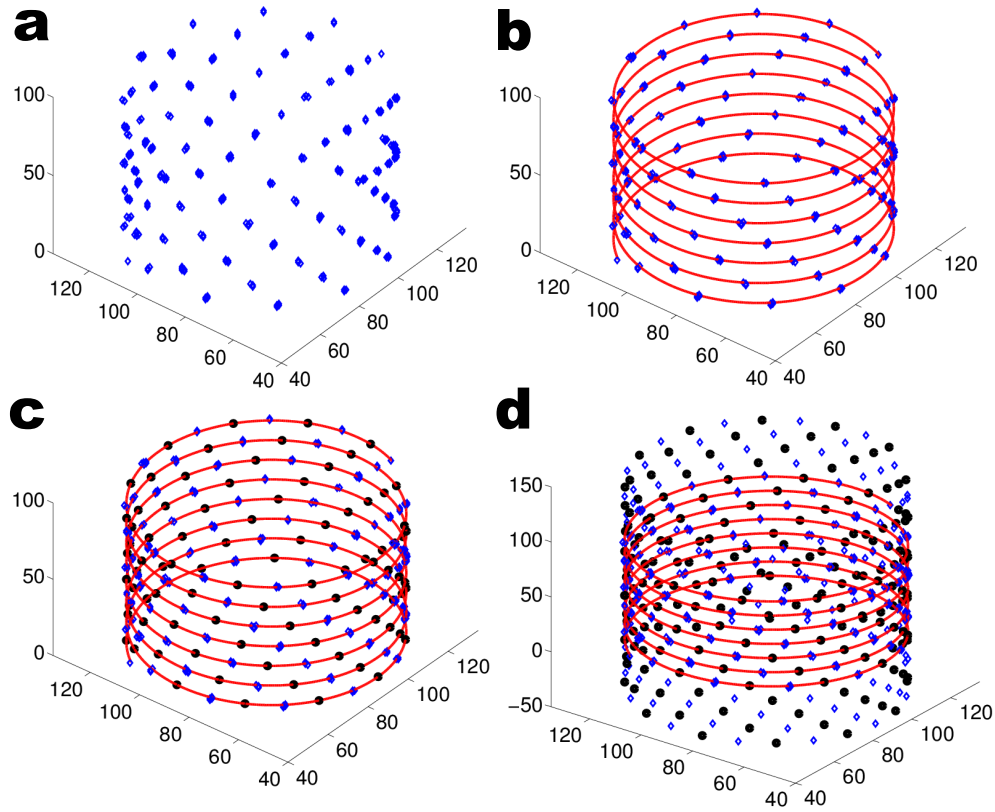
Figure 16: The helix backbone reconstruction for the theoretical model with two types of protein monomers. In **a**, positions of voxels with isovalue larger than 31.2 are marked with solid blue diamonds for EMD 1129 denoising data after 10 iterations. A helix backbone structure can be clearly identified. In **b**, a helix function is parametrized based on fitting with marked positions in **a**. In each circle, there are 12 independent blue color nodes, which represent the same type protein monomers (we call them Type "I" monomers). In **c**, 12 black color nodes are added evenly on helix curve in each circle, representing Type "II" protein monomers. Together with the blue nodes, they are used to compare with the EMD1129 experimental data, in which there are 24 proteins in each circle. In **d**, three and four more layers of proteins are added to two ends of the helix curve to eliminate the boundary effect in our density function.

we propose a coarse-grained model for microtubule.

The EMD-1129 data seriously suffer from noise as demonstrated in Fig. 15 (**a**). To build up a coarse-grained model, a topological denoising process as discussed in the previous section is employed. Surfaces extracted from denoising data are illustrated in Fig. 15. It can be seen that after ten iterations, the noise intensity is dramatically reduced and the basic geometry of the structure is preserved. More iterations are considered due to the requirement of persistent homology analysis, which will be discussed later in Section 4.3.

Based on the denoising data processed with ten iterations, we analyze the structure of this microtubule intermediate and build up coarse-grained models. Through the observation of different isosurfaces from the data it can be found that this microtubule intermediate has a unique helix backbone configuration. We assume that the center of each component protein has the largest electron density value. Through a threshold value of 31.2, we are able to identify centers of these component proteins and the helix backbone is thus constructed. Figure 16 **a** demonstrates the construction of our theoretical model with two types of protein monomers. A helix function as illustrated in Fig. 16 **b** is parametrized based on fitting with these marked positions. It is seen that in each circle, there are about 12 blue color nodes, representing 12 protein monomers of the same type which we denote as type "I". However, there are 12 type "II" monomers missing in this model as they has a slightly lower electron density value. These type "II" protein monomers are further accounted by adding 12 black nodes evenly distributed on the helix curve so that they can pair up with type "I" protein monomers as demonstrated in Fig. 16 **c**. Finally, to avoid the boundary effect in our density function evaluation, about three layers are added to the top and bottom parts of the helix structure.

To avoid the complexities and present our persistent homology analysis directly and clearly, a simple coarse-grained model is considered. Basically, we use ellipsoids to represent protein monomers. The density function of our microtubule intermediate model can be expressed as,

$$\rho(x, y, z) = \sum_i W_i e^{-\left[(\frac{x-x_i}{(\sigma_i^x)})^2 + (\frac{y-y_i}{(\sigma_i^y)})^2 + (\frac{z-z_i}{(\sigma_i^z)})^2\right]}, \tag{13}$$

where $\rho(x, y, z)$ is the density function of the model, parameter $W_i$ is the weight coefficient, parameters $(\sigma_i^x)$, $(\sigma_i^y)$ and $(\sigma_i^z)$ are ellipsoid radii. Coordinates $(x_i, y_i, z_i)$ denote the positions of protein monomer centers. To eliminate the boundary effect, the simulated models incorporate extra protein elements as illustrated in Fig. 16 **d**. All the above parameters are optimized by the least-square fitting using the denoising data. It is found that in any $xy$-cross section, the electron density of microtubule tightly concentrates in a highly symmetric ring-band region as shown in Fig. 17. This ring-band region can be characterized by an inner circles and an outer circle. These circles share the same center at grid position (86, 86) and their radii are 37 and 48 voxels, respectively. As discussed in the literature,[108] only the regions that have sufficiently large density values should be included in the fitting. In the present work, the fitting region is limited to the region within two dash-line circles in the $xy$-cross section as illustrated in Fig. 17. In 3D, a thick-layer-cylinder region that encompasses the structure is considered as the fitting region and is denoted as $V$.

To obtain the optimized fitting parameters, two evaluation coefficients, i.e., cross-correlation coefficients (CCF)

$$\text{CCF} = \frac{\sum_{j \in V} \rho_j^e \sum_{j \in V} \rho_j}{\sqrt{\sum_{j \in V} (\rho_j^e)^2 \sum_{j \in V} (\rho_j)^2}}, \tag{14}$$

and correlation coefficients (CF)

$$\text{CF} = \frac{\sum_{j \in V} (\rho_j^e - \bar{\rho}_j^e) \sum_{j \in V} (\rho_j - \bar{\rho}_j)}{\sqrt{\sum_{j \in V} (\rho_j^e - \bar{\rho}_j^e)^2 \sum_{j \in V} (\rho_j - \bar{\rho}_j)^2}}. \tag{15}$$

are used.[108] Here $\rho_i = \rho(x_j, y_j, z_j)$, $\rho_j^e$ is the experimental electron-density value at $(x_j, y_j, z_j)$ after 10 denoising iterations, $\bar{\rho}$ denotes the average of $\rho$ and parameter $V$ is the fitting region as stated above.

Based on the helix backbone configuration, three theoretical models are constructed for microtubule structure. Using the least square fitting, we determine the fitting parameters, i.e, $W_i$ and $\sigma$ in Eq. (13) for these
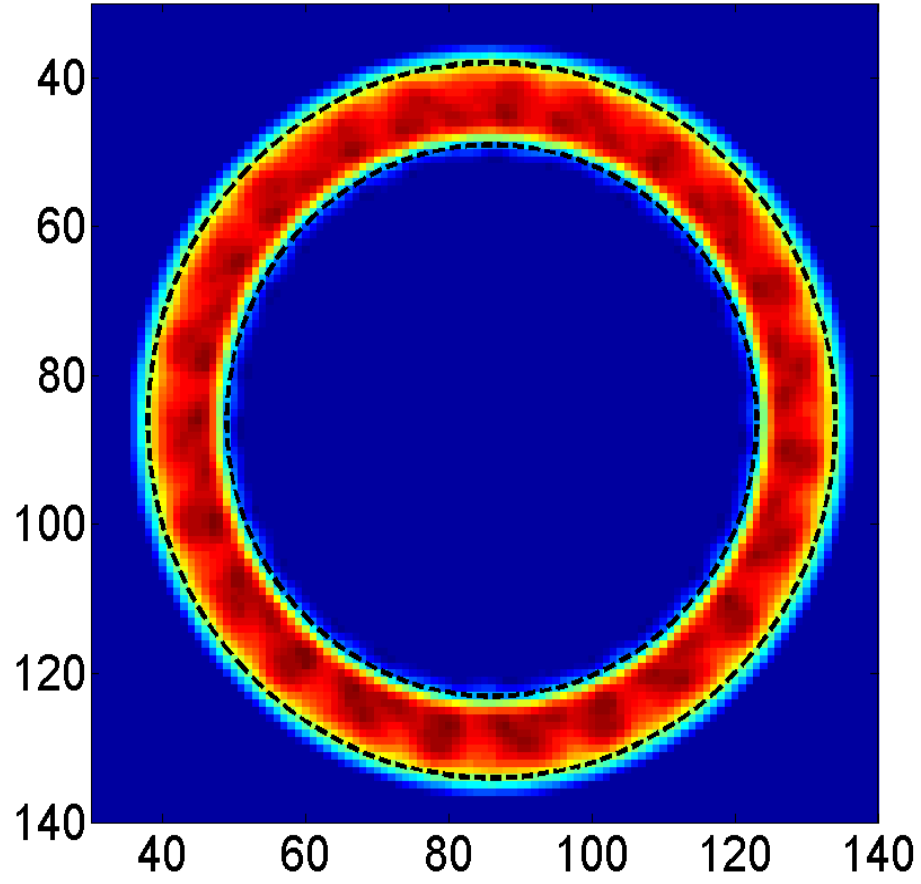
Figure 17: Illustration of the fitting region. The average isovalues over the $z$-axis, i.e., $\frac{1}{n_z} \sum_k \rho^e(x_i, y_j, z_k)$ are given. The fitted region is indicated by two dash lines, i.e., the inner circle and the outer circle. Two circles share the same center point coordinate (86, 86) and their radius are 37 and 48 voxels, respectively.
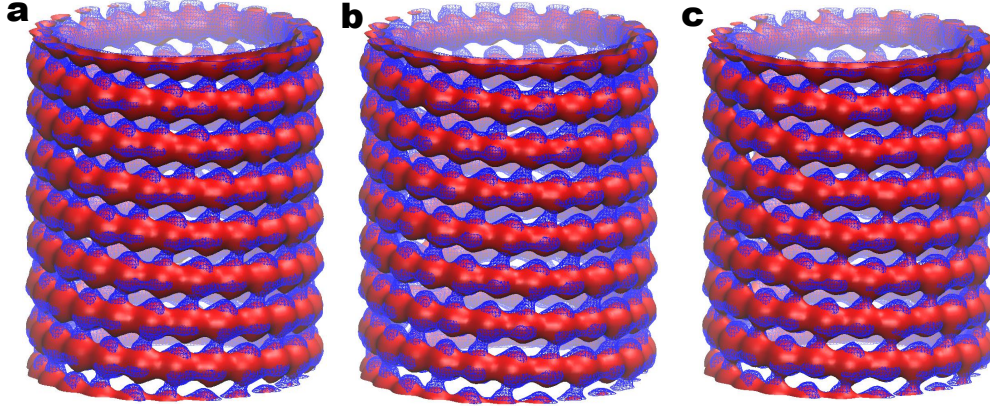
Figure 18: Three theoretical models for microtubule structures constructed from fitting the experimental data using proteins on helix backbone curve. Based on the coarse-grained representation, we use an ellipsoid to represent a protein monomer. In **a**, only one type of ellipsoids is used. In **b**, two types of ellipsoids with different weight functions are used. In **c**, two types of ellipsoids with two different weight functions and modified locations are considered. The isosurfaces in **a** and **b** look similar to each other. However, they have dramatically different topological behaviors. In **c**, we systematically shift type "II" monomers to form dimmers with type "I" monomers. The blue meshed surfaces are obtained from the denoising data, and red solid surfaces are computed from the corresponding three theoretical models.

models. Results are evaluated by aforementioned CCF and CF criteria. In the first model, only one type of ellipsoids is used, i.e., one type of monomers with $W_i = 42$ for all protein monomers. The second model has two types of monomers, see Figure 16. Two types of ellipsoids with $w_1 = 42$ and $w_2 = 38$ are considered by setting $\{W_i; W_i = w_1 \text{ or } w_2\}$. The third model is dimer one. In this model, location modification is considered to generate dimers by shifting the type "II" protein monomers simultaneously along the helix backbone closer to type "I" protein monomers slightly. In this manner, we discriminate between intra-dimer and inter-dimer distances. All ellipsoids are parametrized uniformly by assuming $\sigma_i^x = 24$, $\sigma_i^y = 24$ and $\sigma_i^z = 22$. The unit for these parameters is Angstrom (Å), and voxel spacing is 4 Å. The results are illustrated in Fig. 18. All three models look similarly. Actually, the CCFs for the three models are 0.9601, 0.9607 and 0.9604, respectively. The CFs for them are 0.7392, 0.7436 and 0.7662, respectively. The difference in these coefficients is very small. Therefore, we cannot determine with a good confidence that one model is definitely better than others. We therefore encounter a standard ill-posed inverse problem by using the structural optimization approach.

### 4.3 Persistent homology based microtubule model evaluation

However, if we pay attention to the $\beta_1$ patterns of holes formed between the upper and lower helix circles as illustrated in Fig. 18 (**c**), it can been seen that only the third model preserves these features. These linkage properties are directly related to biomolecular flexibility and functional properties. Therefore it is important for us to characterize and capture them in our models. Unfortunately, least square optimization approach is insensitive to these little structural characteristics. Therefore, techniques that are sensitive to geometric variations are required to guide the structure determination. We show that persistent homology is a desirable technique for detecting geometric changes in the rest of this section.

To understand how the persistent homology can be employed to guide our model construction and evaluation, we investigate the topological fingerprint of the microtubule intermediate structure. As described earlier, due to the noise, a denoising process is required. The geometric flow based denoising algorithm is used with different numbers of iterations. The denoising data are carefully analyzed and the topological persistence results are demonstrated in Fig. 19. It can be seen that a special pattern begin to emerge when the number of iterations approaches 20. More specifically, groups of bars in both $\beta_0$ and $\beta_1$ panels appear with distinctive persisting patterns.

In the $\beta_0$ panel of Fig. 19, bars can be roughly groups into three parts from the top to the bottom, i.e., an irregular "hair-like" part on the top, a narrow regular "body" part in the middle and a large regular "base"
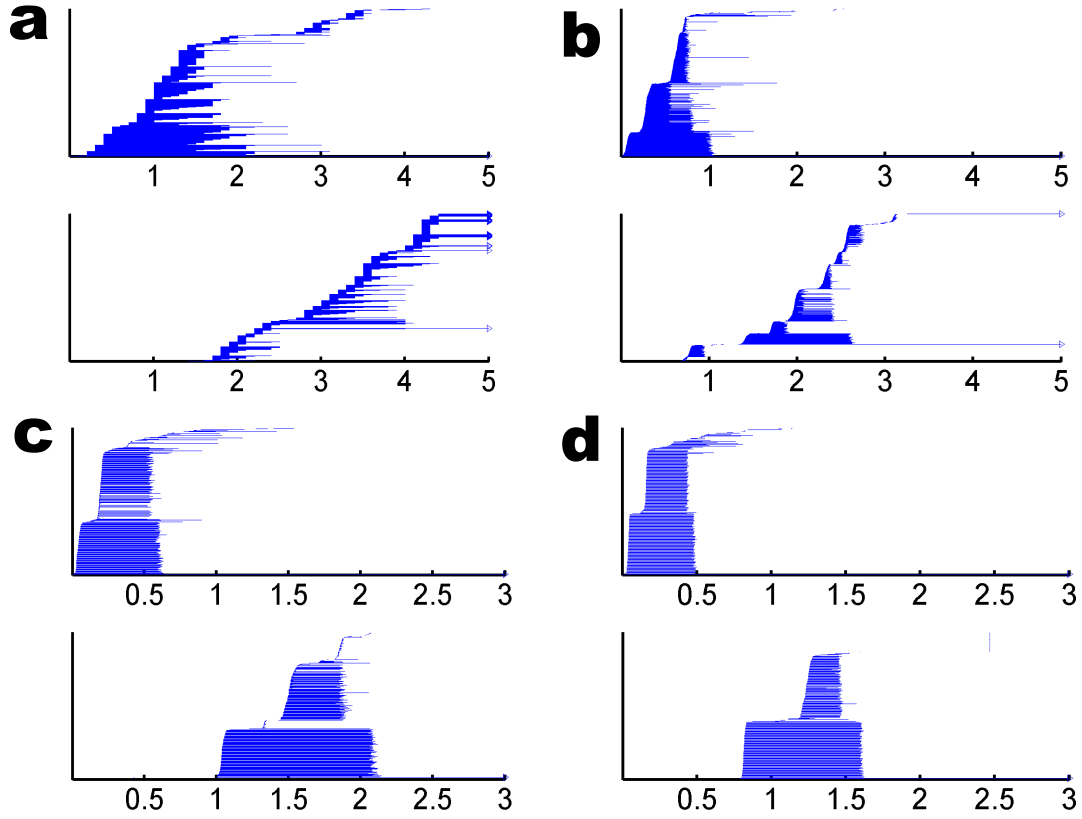
Figure 19: Topological persistence of $\beta_0$ (top row) and $\beta_1$ (bottom row) generated from original and preprocessed EMD 1129 data. **a** is the barcodes of the original EMD1129 data. **b**, **c** and **d** are barcodes for EMD-1129 data after 10, 20 and 40 denoising iterations. A special pattern, i.e., two individual bands of bars in both $\beta_0$ and $\beta_1$, persists in **c** and **d**.



Figure 20: Topological fingerprints of $\beta_0$ (top row) and $\beta_1$ (bottom row) generated from three theoretical models as depicted in Fig. 18. **a** is the barcodes of the first model with only one type of monomers. **b** is the barcodes of the second fitted model with only two types of monomers using different weight functions. **c** is barcodes of the third model with two types of monomers using different weight functions and modified locations. It can be seen that only the final model is able to capture the topological properties of the intrinsic topological fingerprints of cryo-EM data in Fig. 19 **d**.

Figure 21: The topological transitions of microtubule geometry. Here **a**, **b** and **c** are the original data, denoising data (40 iterations) and theoretical model (the third model), respectively. The subscripts 1 to 4 represent four topological transitions during the filtration process, i.e., hetero-dimmer formation, large circles formation, evolution of each large circles into two circles, and finally death of one of two circles. It can be seen that, these topological transitions are well-preserved in our denoising data and theoretical model, which explains the excellent topological consistency between them.

part in the bottom. Topologically, these parts represent different components in the microtubule intermediate structure. The irregular "hair-like" part corresponds to the partial monomer structures located on the top and the bottom boundaries of the structure. As can be seen in Fig. 18, each monomer has lost part of the structure at the boundary regions. The regular "body" and "base" parts are basically related to two types of monomers in the middle region where the structure is free of boundary effect. From the barcodes, it can be seen that "body" part has a later "birth" time and earlier "death" time compared with the "base" barcode part. This is due to the reason that this type of monomers has relative lower electron density. As the filtration is defined to go from highest electron density values to lowest ones, their corresponding barcodes appear later. Their earlier death time, however, is due to the reason that they form dimers with the other type of monomers represented by the "base" barcode part. It can be derived from these nonuniform behavior that monomers are not equally distributed along the helix backbone structure. Instead, two adjacent different types of monomers form a dimer first and then all these dimers simultaneously connect with each other as the filtration goes on. Moreover, from the analysis in the previous section, it is obvious to see that the "body" and "base" parts are topological representations of type "II" monomer and type "I" monomer, respectively.

For the $\beta_1$ panel of Fig. 19, there also exists a consistent pattern when the denoising process passing a certain stage. Two distinctive types of barcodes can be identified in the fingerprint, i.e., a shorter band of barcodes on the top and a longer band of bars on the bottom. Topologically, these $\beta_1$ bars correspond to the rings formed between two adjacent helix circles of monomers or dimers. During the filtration, dimers are formed between type "I" and type "II" monomers and soon after that, all dimers connect with each other and form the helix backbone. As the filtration goes on, type "I" monomers from the upper helix circle first connect with type "II" monomers at the lower circle. Geometrically, this means six monomers, three ("I-II-I") from the upper layer and three ("II-I-II") from the lower layer, form a circle. As the filtration goes further, this circle evolves into two circles when two middle monomers on two layers also connect. However, these two circles do not disappear simultaneously. Instead, one persists longer than the other. This entire process generates the unique topological fingerprint in $\beta_1$ barcodes.

The topological fingerprint we extracted from the denoising process can be used to guide the construction and evaluation of our microtubule models. To this end, we analyze the topological features of three theoretical models. Our persistent homology results for three models are demonstrated in Figs. 20 **a**, **b** and **c**, respectively. It can be seen that all the three models are able to capture the irregular "hair" region in their $\beta_0$ barcode chart. From the topological point of view, the first model is the poorest one. It fails to capture the regular fingerprint patterns in both $\beta_0$ and $\beta_1$ panels of the original cryo-EM structure in Fig. 19 **d**. With two different weight functions to represent two types of monomers, the second model delivers a relatively better topological result. It is able to preserve part of the difference between type "I" and type "II" barcodes in the $\beta_0$ panel. In $\beta_1$ panel, some nonuniform barcodes emerges. The persistent homology results are further improved in the third model when the intra-dimer and inter-dimer interactions are considered. In our third model, fingerprint patterns of the cryo-EM structure in both $\beta_0$ and $\beta_1$ panels of Fig. 19 **d** are essentially recovered by those of Fig. 20 **c**. Even though their scales are different, their shapes are strikingly similar.

## 4.4 Discussion

The essential topological features that are associated with major topological transitions of the original cryo-EM structure are illustrated in Figs. 21 $\mathbf{a}_1$, $\mathbf{a}_2$, $\mathbf{a}_3$ and $\mathbf{a}_4$. As shown in Figs. 21 $\mathbf{b}_1$, $\mathbf{b}_2$, $\mathbf{b}_3$ and $\mathbf{b}_4$, these features have been well-preserved during the denoising process. Our best predicted model is depicted in Figs. 21 $\mathbf{c}_1$, $\mathbf{c}_2$, $\mathbf{c}_3$ and $\mathbf{c}_4$. In these figure labels, subscripts $1, 2, 3$ and $4$ denote four topological transition stages in the filtration process, namely hetero-dimmer formation, large circles formation, evolution of one large circle into two circles, and finally death of one of two circles. By the comparison of denoising results (Figs. $\mathbf{b}_1$, $\mathbf{b}_2$, $\mathbf{b}_3$ and $\mathbf{a}_4$) with original structures (Figs. $\mathbf{a}_1$, $\mathbf{a}_2$, $\mathbf{a}_3$ and $\mathbf{a}_4$), it is seen that, in the noise reduction process, although some local geometric and topological details are removed, fundamental topological characteristics are well preserved. As illustrated in Fig. 19, using the persistent homology description, these fundamental topological characteristics are well-preserved in topological persistence patterns, which are further identified as fingerprints of the microtubule intermediate structure. We believe that topological fingerprints are crucial to the characterization, identification, and analysis of the biological structure. As demonstrated in Figs. 21 $\mathbf{c}_1$, $\mathbf{c}_2$, $\mathbf{c}_3$ and $\mathbf{c}_4$, once our model successfully reproduces the topological fingerprints, the simulated structure is able to capture the essential topological characteristics of the original one. Moreover, through the analysis

in Section 4.3, it can be seen that to reproduce the topological fingerprint of EMD 1129, two conditions are essential. The first is the creation of two types of monomers. The second is the differentiation of intra-dimers and inter-dimers. Biologically, these requirements means: 1) there are two types of monomers, i.e., $\alpha$-tubulin monomers and $\beta$-tubulin monomers; and 2) intra-dimers and inter-dimers should behave differently from hetero-dimers.

It also should be noticed that a higher correlation coefficient may not guarantee the success of the model, especially when the original data is of low resolution and low SNR. As can be seen in Section 4.2, our three theoretical models have very similar fitting coefficients. The second model even has a slightly higher cross-correlation coefficients. However, only the third model is able to reproduce the essential topological features of the original cryo-EM data. This happens as topological invariants, i.e., connected components, circles, loops, holes or void, tend to be very sensitive to "tiny" linkage parts, which are almost negligible in the density fitting process, compared to the major body part. We believe these linage parts play important roles in biological system especially in macroproteins and protein-protein complexes. Different linkage parts generate different connectivity, thus can directly influence biomolecular flexibility, rigidity, and even its functions. By associates topological features with geometric measurements, our persistent homology analysis is able to distinguish these connectivity parts. Therefore, persistent homology is able to play a unique role in protein design, model evaluation and structure determination.

## 5    Conclusion

Cryo-electron microscopy (cryo-EM) is a major workhorse for the investigation of subcellular structures, organelles and large multiprotein complexes. However, cryo-EM techniques and algorithms are far from mature, due to limited sample quality and or stability, low signal to noise (SNR), low resolution, and the high complexity of the underlying molecular structures. Persistent homology is a new branch of topology that is known for its potential in the characterization, identification and analysis (CIA) of big data. In this work, persistent homology is, for the first time, employed for the cryo-EM data CIA.

Methods and algorithms for the geometric and topological modeling are presented. Here, geometric modeling, such the generation of density maps for proteins or other molecules, is employed to create known data sets for validating topological modeling algorithms. We demonstrate that cryo-EM density maps and fullerene density data can be effectively analyzed by using persistent homology.

Since topology is very sensitive to noise, the understanding of the topological signature of noise is a must in cryo-EM CIA. We first investigate the topological fingerprint of Gaussian noise. We reveal that for the Gaussian noise, its topological invariants, i.e., $\beta_0$, $\beta_1$ and $\beta_2$ numbers, all exhibit the Gaussian distribution in the filtration space, i.e., the space of volumetric density isovalues. At a low SNR, signal and noise are inseparable in the filtration space. However, after denoising with the geometric flow method, there is clear separation between signal and noise for various topological invariants. As such, a simple threshold can be prescribed to effectively remove noise. For the case of low SNR, the understanding of noise characteristic in the filtration space enable us to use persistent homology as an efficient means to monitor and control the noise removal process. This new strategy for noise reduction is called topological denoising.

Persistent homology has been applied to the theoretical modeling of a microtubule structure (EMD 1129). The backbone of the microtubule has a helix structure. Based on the helix structure, we propose three theoretical models. The first model assumes that protein monomers form the helix structure. The second model adopts two types of protein monomers evenly distributed along helix chain. The last model utilizes a series of protein dimers along the helix chain. These models are fitted with experimental data by the least square optimization method. It is found that all the three models give rise to similar high correlation coefficients with the experimental data, which indicates that the structural optimization is ill-posed. However, the topological fingerprints of three models are dramatically different. In the denoising process, the cryo-EM data of the microtubule structure demonstrate a consistent pattern which can be recognized as the intrinsic topological fingerprint of the microtubule structure. By careful examination of the fingerprint, we reveal two essential topological characteristics which discriminate the protein dimers from the monomers. As such, we conclude that only the third model, i.e., the protein dimer model, is able to capture the intrinsic topological characteristics of the cryo-EM structure and must be the best model for the experimental data. It is believed that the present work offers a novel topology based strategy for resolving ill-posed inverse problems.

**References**

[1] S. Nickell, C. Kofler, A. P. Leis, and W. Baumeister. A visual approach to proteomics. *Nat. Rev. Mol. Cell Biol.*, 7:225–230, 2006.

[2] C. V. Robinson, A. Sali, and W. Baumeister. The molecular sociology of the cell. *Nature*, 450:973C982, 2007.

[3] A. Leis, B. Rockel, L. Andrees, and W. Baumeister. Visualizing cells at the nanoscale. *Trends Biochem. Sci.*, 34:60C70, 2009.

[4] E. I. Tocheva, Z. Li, and G. J. Jensen. Electron cryotomography. In *Cold Spring Harb. Perspect. Biol.*, volume 2, page A003442, 2010.

[5] N. Volkmann. Methods for segmentation and interpretation of electron tomographic reconstructions. In *Methods Enzymol*, volume 483, pages 31–46, 2010.

[6] S. S. Abeysinghe, M. Baker, W. Chiu, and J. Tao. Segmentation-free skeletonization of grayscale volumes for shape understanding. *IEEE International Conference on Shape Modeling and Applications*, SMI Stony Brook, N.Y.:63–71, 2012.

[7] A. Biswas, D. Si, K. Al Nasr, D. Ranjan, M. Zubair, and J. He. Improved effeciency in Cryo-EM secondary structure topology determination from inaccurate data. *Journal of Bioinformatics and Computational Biology*, 10:1242006, 2012.

[8] T. Ju, M. L. Baker, and W. Chiu. Computing a family of skeletons of volumetric models for shape description. *Computer-Aided Design*, 39:352–360, 2007.

[9] M. L. Baker, W. Jiang, W. J. Wedemeyer, F. J. Rixon, D. Baker, and W. Chiu. Ab initio modeling of the herpesvirus vp26 core domain assessed by cryoem density. *PLoS Computational Biology*, 2:e146, 2006.

[10] W. Sun and Jing He. Native secondary structure topology has near minimum contact energy among all possible geometrically constrained topolgies. *Proteins: Structure, Function and Bioinformatics*, 77:159–173, 2009.

[11] Y. Lu and J. He. Deriving topology and sequence alignment for helix skeleton in low resolution protein density maps. *Journal of Bionformatics and Computational Biology*, 8:183–201, 2008.

[12] E. F. Pettersen, T. D. Goddard, C. C. Huang, G. S. Couch, D. M. Greenblatt, E. C. Meng, and T. E. Ferrin. UCSF Chimera–a visualization system for exploratory research and analysis. *J. Comput. Chem.*, 25:1605–1612, 2004.

[13] Shawn Q. Zhenga, Bettina Keszthelyi, Eric Branlunda, John M. Lyleb, Michael B. Braunfelda, John W. Sedatb, and David A. Agarda. UCSF tomography: an integrated software suite for real-time electron microscopic tomographic data collection, alignment, and reconstruction. *J. Struct. Biol.*, 157:138–147, 2007.

[14] F. Amat, F. Moussavi, L. R. Comolli, G. Elidan, K. H. Downing, and M. Horowitz. Markov random field based automatic image alignment for electron tomography. *J. Struct. Biol.*, 161:260–275, 2008.

[15] T. Hrabe, Y. Chen, S. Pfeffer, L.K. Cuellar, A.V. Mangold, and F. Forster. PyTom: a python-based toolbox for localization of macromolecules in cryo-electron tomograms and subtomogram analysis. *J. Struct. Biol.*, 178:178–188, 2012.

[16] J. R. Kremer, D. N. Mastronarde, and J. R. McIntosh. Computer visualization of three-dimensional image data using IMOD. *J. Struct. Biol.*, 116:71–76, 1996.

[17] D. Ress, M. L. Harlow, M. Schwarz, R. M. Marshall, and U. J. McMahan. Automatic acquisition of fiducial markers and alignment of images in tilt series for electron tomography. *J Electron Microsc (Tokyo)*, 48:277–287, 1999.

[18] C. A. Rohl, C. E. Strauss, K. M. Misura, and D. Baker. Protein structure prediction using Rosetta. *Methods Enzymol*, 383:66–93, 2004.

[19] A. Stoschek and R. Hegerl. Denoising of electron tomographic reconstructions using multiscale transformations. *J. Struct. Biol.*, 120:257C265, 1997.

[20] A. S. Frangakis and R. Hegerl. Noise reduction in electron tomographic reconstructions using nonlinear anisotropic diffusion. *J. Struct. Biol.*, 135:239C250, 2001.

[21] S. Fernandez, J. J.and Li. An improved algorithm for anisotropic nonlinear diffusion for denoising cryo-tomograms. *J. Struct. Biol.*, 144:152C161, 2003.

[22] J.J. Fernandez. Tomobflow: feature-preserving noise filtering for electron tomography. *BMC Bioinformatics*, 178:1–10, 2009.

[23] C. Tomasi and R. Manduchi. Bilateral filtering for gray and color images. *Proc. ICCV*, 98:839C846, 1998.

[24] W. Jiang, Q. Baker, M. L.and Wu, C. Bajaj, and W. Chiu. Applications of a bilateral denoising filter in biological electron microscopy. *J. Struct. Biol.*, 144:114 – 122, 2003.

[25] R. S. Pantelic, C. Y. Rothnagel, R.and Huang, D. Muller, D. Woolford, M. J. Landsberg, A. McDowall, B. Pailthorpe, P. R. Young, J. Banks, B. Hankamer, and G. Ericksson. The discriminative bilateral filter: An enhanced denoising filter for electron microscopy data. *J. Struct. Biol.*, 155:395C408, 2006.

[26] P. van der Heide, X. P. Xu, B. J. Marsh, D. Hanein, and N. Volkmann. Efficient automatic noise reduction of electron tomographic reconstructions based on iterative median filtering. *J. Struct. Biol.*, 158:196C204, 2007.

[27] Kunyu Tsai, Jianwei Ma, Datian Ye, and Jian Wu. Curvelet processing of mri for local image enhancement. *International Journal for Numerical Methods in Biomedical Engineering*, 28:661–677, 2012.

[28] Mei-sen Pan, Jing-tian Tang, Qiu-sheng Rong, and Fen Zhang. Medical image registration using modified iterative closest points. *International Journal for Numerical Methods in Biomedical Engineering*, 27:1150–1166, 2011.

[29] A. G. Radaelli and J. Peiro. On the segmentation of vascular geometries from medical images. *International Journal for Numerical Methods in Biomedical Engineering*, 26:3–34, 2010.

[30] Issei Fujishiro, Yuriko Takeshima, Taeko Azuma, and Shigeo Takahashi. Volume data mining using 3d field topology analysis. *IEEE Computer Graphics and Applications*, 20(5):46–51, 2000.

[31] G. Carlsson. Topology and data. *Am. Math. Soc*, 46(2):255–308, 2009.

[32] R. Ghrist. Barcodes: The persistent topology of data. *Bull. Amer. Math. Soc.*, 45:61–75, 2008.

[33] Harish Doraiswamy and Vijay Natarajan. Computing reeb graphs as a union of contour trees. *IEEE Transactions on Visualization and Computer Graphics*, 19:249–262, 2013.

[34] Patrizio Frosini and Claudia Landi. Size theory as a topological tool for computer vision. *Pattern Recognition and Image Analysis*, 9(4):596–603, 1999.

[35] Vanessa Robins. Towards computing homology from finite approximations. In *Topology Proceedings*, volume 24, pages 503–532, 1999.

[36] H. Edelsbrunner, D. Letscher, and A. Zomorodian. Topological persistence and simplification. *Discrete Comput. Geom.*, 28:511–533, 2002.

[37] A. Zomorodian and G. Carlsson. Computing persistent homology. *Discrete Comput. Geom.*, 33:249–274, 2005.

[38] Peter Bubenik and Peter T. Kim. A statistical approach to persistent homology. *Homology, Homotopy and Applications*, 19:337–362, 2007.

[39] Herbert Edelsbrunner and John Harer. *Computational topology: an introduction.* American Mathematical Soc., 2010.

[40] T. K. Dey, K. Y. Li, J. Sun, and C. S. David. Computing geometry aware handle and tunnel loops in 3d models. *ACM Trans. Graph.*, 27, 2008.

[41] Tamal K. Dey and Y. S. Wang. Reeb graphs: Approximation and persistence. *Discrete and Computational Geometry*, 49(1):46–73, 2013.

[42] K. Mischaikow and V. Nanda. Morse theory for filtrations and efficient computation of persistent homology. *Discrete and Computational Geometry*, 50(2):330–353, 2013.

[43] G. Carlsson, T. Ishkhanov, V. Silva, and A. Zomorodian. On the local behavior of spaces of natural images. *International Journal of Computer Vision*, 76(1):1–12, 2008.

[44] D. Pachauri, C. Hinrichs, M.K. Chung, S.C. Johnson, and V. Singh. Topology-based kernels with application to inference problems in alzheimer's disease. *Medical Imaging, IEEE Transactions on*, 30(10):1760–1770, Oct 2011.

[45] G. Singh, F. Memoli, T. Ishkhanov, G. Sapiro, G. Carlsson, and D. L. Ringach. Topological analysis of population activity in visual cortex. *Journal of Vision*, 8(8), 2008.

[46] Paul Bendich, Herbert Edelsbrunner, and Michael Kerber. Computing robustness and persistence for images. *IEEE Transactions on Visualization and Computer Graphics*, 16:1251–1260, 2010.

[47] Patrizio Frosini and Claudia Landi. Persistent betti numbers for a noise tolerant shape-based approach to image retrieval. *Pattern Recognition Letters*, 34:863–872, 2013.

[48] K. Mischaikow, M Mrozek, J. Reiss, and A. Szymczak. Construction of symbolic dynamics from experimental time series. *Physical Review Letters*, 82:1144–1147, 1999.

[49] T. Kaczynski, K. Mischaikow, and M. Mrozek. *Computational homology.* Springer-Verlag, 2004.

[50] V. D. Silva and R Ghrist. Blind swarms for coverage in 2-d. In *In Proceedings of Robotics: Science and Systems*, page 01, 2005.

[51] H Lee, H. Kang, M. K. Chung, B. Kim, and D. S. Lee. Persistent brain network homology from the perspective of dendrogram. *Medical Imaging, IEEE Transactions on*, 31(12):2267–2277, Dec 2012.

[52] D. Horak, S Maletic, and M. Rajkovic. Persistent homology of complex networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2009(03):P03034, 2009.

[53] P. Niyogi, S. Smale, and S. Weinberger. A topological view of unsupervised learning from noisy data. *SIAM Journal on Computing*, 40:646–663, 2011.

[54] Bei Wang, Brian Summa, Valerio Pascucci, and M. Vejdemo-Johansson. Branching and circular features in high dimensional data. *IEEE Transactions on Visualization and Computer Graphics*, 17:1902–1911, 2011.

[55] Bastian Rieck, Hubert Mara, and Heike Leitte. Multivariate data analysis using persistence-based filtering and topological signatures. *IEEE Transactions on Visualization and Computer Graphics*, 18:2382–2391, 2012.

[56] Xu Liu, Zheng Xie, and Dongyun Yi. A fast algorithm for constructing topological structure in large data. *Homology, Homotopy and Applications*, 14:221–238, 2012.

[57] Barbara Di Fabio and Claudia Landi. A mayer-vietoris formula for persistent homology with an application to shape recognition in the presence of occlusions. *Foundations of Computational Mathematics*, 11:499–527, 2011.

[58] P. M. Kasson, A. Zomorodian, S. Park, N. Singhal, L. J. Guibas, and V. S. Pande. Persistent voids a new structural metric for membrane fusion. *Bioinformatics*, 23:1753–1759, 2007.

[59] M. Gameiro, Y. Hiraoka, S. Izumi, M. Kramar, K. Mischaikow, and V. Nanda. Topological measurement of protein compressibility via persistence diagrams. *preprint*, 2013.

[60] Y. Dabaghian, F. Memoli, L. Frank, and G. Carlsson. A topological paradigm for hippocampal spatial map formation using persistent homology. *PLoS Comput Biol*, 8(8):e1002581, 08 2012.

[61] J. Kloke and G. Carlsson. Topological de-noising: Strengthening the topological signal. *arXiv preprint arXiv:0910.5947*, 2009.

[62] K. L. Xia and G. W. Wei. Persistent homology analysis of protein structure, flexibility and folding. *International Journal for Numerical Methods in Biomedical Engineerings*, 30:814–844, 2014.

[63] V. De Silva and G. Carlsson. Topological estimation using witness complexes. In *Proceedings of the First Eurographics conference on Point-Based Graphics*, pages 157–166. Eurographics Association, 2004.

[64] D. Günther, A. Jacobson, J. Reininghaus, H. P. Seidel, O. Sorkine-Hornung, and T. Weinkauf. Fast and memory-efficient topological denoising of 2D and 3D scalar fields. *IEEE Transactions on Visualization and Computer Graphics*, 20:12, 2014.

[65] U. Bauer, C. B. Schönlieb, and M. Wardetzky. Total variation meets topological persistence: A first encounter. *AIP Conference Proceedings*, 1281(1):1022, 2010.

[66] R. J. Adler, O. Bobrowski, M. S. Borman, E. Subag, and S. Weinberger. Persistent homology for random fields and complexes. In *Borrowing Strength: Theory Powering Applications–A Festschrift for Lawrence D. Brown*, volume 6, pages 124–143. Institute of Mathematical Statistics, 2010.

[67] E. F. Pettersen, T. D. Goddard, C. C. Huang, G. S. Couch, D. M. Greenblatt, E. C. Meng, and T. E. Ferrin. UCSF Chimera–A visualization system for exploratory research and analysis. *Journal of computational chemistry*, 25(13):1605–1612, 2004.

[68] K. L. Xia, X. Feng, Y. Y. Tong, and G. W. Wei. Persistent homology for the quantitative prediction of fullerene stability. *Journal of Computational Chemistry*, in press 2014.

[69] X. Feng, K. L. Xia, Y. Y. Tong, and G. W. Wei. Multiscale geometric modeling of macromolecules II: lagrangian representation. *Journal of Computational Chemistry*, 34:2100–2120, 2013.

[70] K. L. Xia, X. Feng, Y. Y. Tong, and G. W. Wei. Multiscale geometric modeling of macromolecules i: Cartesian representation. *Journal of Computational Physics*, 275:912–936, 2014.

[71] R. B. Corey and L. Pauling. Molecular models of amino acids, peptides and proteins. *Rev. Sci. Instr.*, 24:621–627, 1953.

[72] B. Lee and F. M. Richards. The interpretation of protein structures: estimation of static accessibility. *J Mol Biol*, 55(3):379–400, 1971.

[73] F. M. Richards. Areas, volumes, packing, and protein structure. *Annual Review of Biophysics and Bioengineering*, 6(1):151–176, 1977.

[74] W. Rocchia, S. Sridharan, A. Nicholls, E Alexov, A Chiabrera, and B. Honig. Rapid grid-based construction of the molecular surface and the use of induced surface charge to calculate reaction field energies: Applications to the molecular systems and geometric objects. *Journal of Computational Chemistry*, 23:128 – 137, 2002.

[75] M. L. Connolly. Depth buffer algorithms for molecular modeling. *J. Mol. Graphics*, 3:19–24, 1985.

[76] F. Eisenhaber and P. Argos. Improved strategy in analytic surface calculation for molecular systems: Handling of singularities and computational efficiency. *J. Comput. Chem.*, 14:1272–1280, 1993.

[77] V. Gogonea and E. E. Osawa. Implementation of solvent effect in molecular mechanics. 1. model development and analytical algorithm for the solvent -accessible surface area. *Supramol. Chem.*, 3:303–317, 1994.

[78] M. F. Sanner, A. J. Olson, and J. C. Spehner. Reduced surface: An efficient way to compute molecular surfaces. *Biopolymers*, 38:305–320, 1996.

[79] G. W. Wei. Differential geometry based multiscale models. *Bulletin of Mathematical Biology*, 72:1562 – 1622, 2010.

[80] Z. Chen, N. A. Baker, and G. W. Wei. Differential geometry based solvation models I: Eulerian formulation. *J. Comput. Phys.*, 229:8231–8258, 2010.

[81] G. W. Wei, Y. H. Sun, Y. C. Zhou, and M. Feig. Molecular multiresolution surfaces. *arXiv:math-ph/0511001v1*, pages 1 – 11, 2005.

[82] P. W. Bates, G. W. Wei, and S. Zhao. The minimal molecular surface. *arXiv:q-bio/0610038v1*, [q-bio.BM], 2006.

[83] P. W. Bates, G. W. Wei, and Shan Zhao. Minimal molecular surfaces and their applications. *Journal of Computational Chemistry*, 29(3):380–91, 2008.

[84] P. W. Bates, Z. Chen, Y. H. Sun, G. W. Wei, and S. Zhao. Geometric and potential driving formation and evolution of biomolecular surfaces. *J. Math. Biol.*, 59:193–231, 2009.

[85] Guo-Wei Wei, Qiong Zheng, Zhan Chen, and Kelin Xia. Variational multiscale models for charge transport. *SIAM Review*, 54(4):699 – 754, 2012.

[86] Guo-Wei Wei. Multiscale, multiphysics and multidomain models I: Basic theory. *Journal of Theoretical and Computational Chemistry*, 12(8):1341006, 2013.

[87] Guoliang Xu, Qing Pan, and Chandrajit L. Bajaj. Discrete surface modeling using partial differential equations. *Computer Aided Geometric Design*, 23(2):125–145, 2006.

[88] L. T. Cheng, Joachim Dzubiella, Andrew J. McCammon, and B. Li. Application of the level-set method to the implicit solvation of nonpolar molecules. *Journal of Chemical Physics*, 127(8), 2007.

[89] Shan Zhao. Pseudo-time-coupled nonlinear models for biomolecular surface representation and solvation analysis. *International Journal for Numerical Methods in Biomedical Engineering*, 27:1964–1981, 2011.

[90] Shan Zhao. Operator splitting adi schemes for pseudo-time coupled nonlinear solvation simulations. *Journal of Computational Physics*, 257:1000 – 1021, 2014.

[91] K. L. Xia, K. Opron, and G. W. Wei. Multiscale multiphysics and multidomain models — Flexibility and rigidity. *Journal of Chemical Physics*, 139:194109, 2013.

[92] K. Opron, K. L. Xia, and G. W. Wei. Fast and anisotropic flexibility-rigidity index for protein flexibility and fluctuation analysis. *Journal of Chemical Physics*, 140:234105, 2014.

[93] G. W. Wei. Wavelets generated by using discrete singular convolution kernels. *Journal of Physics A: Mathematical and General*, 33:8577 – 8596, 2000.

[94] Y. Zhang, I. A. Hubner, A. K. Arakaki, E. Shakhnovich, and J. Skolnick. On the origin and completeness of highly likely single domain protein structures. *Proc. Natl. Acad. Sci. USA*, 103:2605–2610, 2006.

[95] Z. Y. Yu, M. Holst, Y. Cheng, and J. A. McCammon. Feature-preserving adaptive mesh generation for molecular shape modeling and simulation. *Journal of Molecular Graphics and Modeling*, 26:1370–1380, 2008.

[96] Q. Zheng, S. Y. Yang, and G. W. Wei. Molecular surface generation using PDE transform. *International Journal for Numerical Methods in Biomedical Engineering*, 28:291–316, 2012.

[97] K. L. Xia and G. W. Wei. A galerkin formulation of the mib method for three dimensional elliptic interface problems. *Computers and Mathematics with Applications*, submitted 2013.

[98] K. L. Xia and G. W. Wei. A stochastic model for protein flexibility analysis. *Physical Review E*, 88:062709, 2013.

[99] Andrew Tausz, Mikael Vejdemo-Johansson, and Henry Adams. Javaplex: A research software package for persistent (co)homology. Software available at http://code.google.com/p/javaplex, 2011.

[100] Vidit Nanda. Perseus: the persistent homology software. Software available at http://www.sas.upenn.edu/~vnanda/perseus.

[101] G. W. Wei. Generalized Perona-Malik equation for image restoration. *IEEE Signal Processing Lett.*, 6:165–167, 1999.

[102] M. Lysaker, A. Lundervold, and X. C. Tai. Noise removal using fourth-order partial differential equation with application to medical magnetic resonance images in space and time. *IEEE Transactions on Image Processing*, 12(12):1579–1590, 2003.

[103] G. Gilboa, N. Sochen, and Y. Y. Zeevi. Image sharpening by flows based on triple well potentials. *Journal of Mathematical Imaging and Vision*, 20(1-2):121–131, 2004.

[104] B. T. Fasy, F. Lecci, A. Rinaldo, L. Wasserman, S. Balakrishnan, and A. Singh. Statistical inference for persistent homology: Confidence sets for persistence diagrams. *arXiv preprint arXiv:1303.7117*, 2013.

[105] Q. Qiao, C. H. Yang, C. Zheng, L. Fontán, L. David, X. Yu, C. Bracken, M. Rosen, A. Melnick, E. H. Egelman, and H. Wu. Structural architecture of the CARMA1/Bcl10/MALT1 signalosome: nucleation-induced filamentous assembly. *Molecular cell*, 51(6):766–779, 2013.

[106] E. Nogales and H. W. Wang. Structural intermediates in microtubule assembly and disassembly: how and why? *Current opinion in cell biology*, 18(2):179–184, 2006.

[107] W. H. Wang and E. Nogales. Nucleotide-dependent bending flexibility of tubulin regulates microtubule assembly. *Nature*, 435:911–915, 2005.

[108] M. Topf, K. Lasker, B. Webb, H. Wolfson, W. Chiu, and A. Sali. Protein structure fitting and refinement guided by Cryo-EM density. *Structure*, 16(2):295–307, 2008.