

Persistent Homology for the Quantitative Prediction of Fullerene Stability

Kelin Xia,^[a] Xin Feng,^[b] Yiyong Tong,^{*,[b]} and Guo Wei Wei^{*,[a,c]}

Persistent homology is a relatively new tool often used for *qualitative* analysis of intrinsic topological features in images and data originated from scientific and engineering applications. In this article, we report novel *quantitative* predictions of the energy and stability of fullerene molecules, the very first attempt in using persistent homology in this context. The ground-state structures of a series of small fullerene molecules are first investigated with the standard Vietoris–Rips complex. We decipher all the barcodes, including both short-lived local bars and long-lived global bars arising from topological invariants, and associate them with fullerene structural details. Using accumulated bar lengths, we build quantitative models to correlate local and global Betti-2 bars,

respectively with the heat of formation and total curvature energies of fullerenes. It is found that the heat of formation energy is related to the local hexagonal cavities of small fullerenes, while the total curvature energies of fullerene isomers are associated with their sphericities, which are measured by the lengths of their long-lived Betti-2 bars. Excellent correlation coefficients (>0.94) between persistent homology predictions and those of quantum or curvature analysis have been observed. A correlation matrix based filtration is introduced to further verify our findings. © 2014 Wiley Periodicals, Inc.

DOI: 10.1002/jcc.23816

Introduction

Persistent homology, a method for studying topological features over changing scales, has received tremendous attention in the past decade.^[1,2] The basic idea is to measure the life cycle of topological features within a filtration, that is, a nested family of abstract simplicial complexes, such as Vietoris–Rips complexes, Čech complexes, or alpha complexes.^[3] Thus, long-lived topological characteristics, which are often the intrinsic invariants of the underlying system, can be extracted; while short-lived features are filtered out. The essential topological characteristics of three-dimensional (3D) objects typically include connected components, tunnels or rings, and cavities or voids, which are invariant under the nondegenerate deformation of the structure. Homology characterizes such structures as groups, whose generators can be considered independent components, tunnels, cavities, and so forth. Their times of “birth” and “death” can be measured by a function associated with the filtration, calculated with ever more efficient computational procedures,^[4–7] and further visualized through barcodes,^[8] a series of horizontal line segments with the horizontal x -axis representing the changing scale and the vertical y -axis representing the index of the homology generators. Numerous software packages, such as Perseus, Dionysus, and Javaplex,^[9] based on various algorithms have been developed and made available in the public domain. As an efficient tool to unveil topological invariants, persistent homology has been applied to various fields, such as image analysis,^[10–12] chaotic dynamics verification,^[13,14] sensor network,^[15] complex network,^[16,17] data analysis,^[18] geometric processing,^[19] and computational biology.^[20–23] Based on persistent homology analysis, we have proposed molecular topological fingerprints and used them to reveal the topology-function relationship of biomolecules.^[24] In general, persistent homology is devised as

a robust but *qualitative* topological tool and has been hardly used as a precise *quantitative* predictive tool.^[25,26]

To the best of our knowledge, persistent homology has not been applied to the study of fullerenes, special molecules comprised of only carbon atoms. The fullerene family shares the same closed carbon-cage structure, which contains only pentagonal and hexagonal rings. In 1985, Kroto et al.^[27] proposed the first structure of C_{60} , which was then confirmed in 1990 by Kratschmer et al.^[28] in synthesizing macroscopic quantities of C_{60} . Enormous interest has been aroused by these interesting discoveries. However, there are many challenges. Among them, finding the ground-state structure has been a primary target.

In general, two types of approaches are commonly used.^[29–34] The first method is based on the geometric and topological symmetries of fullerene.^[29–31] In this approach, one first constructs all possible isomers, and then chooses the best possible candidate based on the analysis of the highest-occupied molecular orbital (HOMO) energy and the lowest-unoccupied molecular orbital (LUMO) energy.^[30] In real applications, to generate all possible isomers for a fullerene with a

[a] K. Xia, G. W. Wei
Department of Mathematics, Michigan State University, Michigan 48824

[b] X. Feng, Y. Tong
Department of Computer Science and Engineering, Michigan State University, Michigan 48824
E-mail: ytong@msu.edu

[c] G. W. Wei
Department of Biochemistry and Molecular Biology, Michigan State University, Michigan 48824
E-mail: wei@math.msu.edu

Contract grant sponsor: NSF; Contract grant numbers: IIS-0953096, IIS-1302285, DMS-1160352; Contract grant sponsor: NIH; Contract grant number: R01GM-090208; Contract grant sponsor: MSU Center for Mathematical Molecular Biosciences initiative

© 2014 Wiley Periodicals, Inc.

given atom count is nontrivial until the introduction of Coxeter's construction method^[29,35] and the ring spiral method.^[30] In Coxeter's method, the icosahedral triangulations of the sphere are analyzed to evaluate the possible isomer structures. This method is mathematically rigorous. However, practical applications run into issues with low-symmetry structures. Conversely, based on the spiral conjecture,^[31] the ring spiral method simply lists all possible spiral sequences of pentagons and hexagons, and then winds them up into fullerenes. When a consistent structure is found, an isomer is generated; otherwise, the sequence is skipped. Although the conjecture breaks down for fullerenes with 380 or more atoms, the spiral method proves to be quite efficient.^[31]

For each isomer, its electronic structure can be modeled simply by the Hückel molecular orbital theory,^[36] which is known to work well for planar aromatic hydrocarbons using standard C—C and C—H σ bond energies. Similarly, the bonding connectivities in fullerene structures are used to evaluate orbital energies. The stability of the isomers, according to Manolopoulos,^[30] can then be directly related to the calculated HOMO-LUMO energy gap. However, this model falls short for large fullerene molecules. Even for small structures, its prediction tends to be inaccurate. One possible reason is fullerene's special cage structures. Instead of a planar shape, the structure usually has local curvatures, which jeopardizes the σ - π orbital separation.^[31,37] To account for curvature contributions, a strain energy is considered. It is found that the strain energy reaches its minimum when pentagonal faces are as far away as possible from each other. This is highly consistent with the isolated pentagon rule (IPR)—the most stable fullerenes are those in which all the pentagons are isolated.^[31]

Another approach to obtain ground-state structures for fullerene molecules is through simulated annealing.^[32–34] This global optimization method works well for some structures. However, if large energy barriers exist in the potential, the whole system is prone to be trapped into metastable high-energy state. This happens as breaking the carbon bonds and rearranging the structure need a huge amount of energy. A revised method is to start the system from a special face-dual network and then use the tight-binding potential model.^[34,38] This modified algorithm manages to generate the C_{60} structure of I_h symmetry that has the HOMO-LUMO energy gap of 1.61 eV, in contrast to 1.71 eV obtained using the ab initio local-density approximation.

In this article, persistent homology is, for the first time, used to quantitatively predict the stability of the fullerene molecules. The ground-state structures of a few small fullerene molecules are first studied using a distance based filtration process. Essentially, we associate each carbon atom of a fullerene with an ever-increasing radius and thus define a Vietoris–Rips complex. The calculated Betti numbers (i.e., ranks of homology groups), including β_0 , β_1 , and β_2 , are provided in the barcode representation. To further exploit the persistent homology, we carefully discriminate between the local short-lived and global long-lived bars in the barcodes. We define an average accumulated bar length as the negative arithmetic mean of β_2 bars. As the local β_2 bars represent the number of

cavities of the structure, when β_2 becomes larger, interconnectedness (and thus stability) tends to increase, and relative energy tends to drop. Therefore, the average accumulated bar length indicates the level of a relative energy. We validate this hypothesis with a series of ground-state structures of small fullerenes. It is found that our average accumulated bar length can capture the energy behavior remarkably well, including an anomaly in fullerene C_{60} energy. Additionally, we explore the relative stability of fullerene isomers. The persistence of the Betti numbers is calculated and analyzed. Our results are validated with the total curvature energies of two fullerene families. It is observed that the total curvature energies of fullerene isomers can be well represented with their lengths of the long-lived Betti-2 bars, which indicates the sphericity of fullerene isomers. For fullerenes C_{40} and C_{44} , correlation coefficients up to 0.956 and 0.948 are attained in the distance based filtration. Based on the flexibility-rigidity index (FRI),^[39–41] a correlation matrix based filtration process is proposed to validate our findings.

The rest of this article is organized as follows. In Section Rudiments of Persistent Homology, we discuss the basic persistent homology concepts, including simplices and simplicial complexes, chains, homology, and filtration. Section Algorithms for Persistent Homology is devoted to the description of algorithms. The alpha complex and Vietoris–Rips complex are discussed in some detail, including filtration construction, metric space design, and persistence evaluation. In Section Application to Fullerene Structure Analysis and Stability Prediction, persistent homology is used in the analysis of fullerene structure and stability. After a brief discussion of fullerene structural properties, we elaborate on their barcode representation. The average accumulated bar length is introduced and applied to the energy estimate of the small fullerene series. By validating with total curvature energies, our persistent homology based quantitative predictions are shown to be accurate. Fullerene isomer stability is also analyzed using the new correlation matrix based filtration. This article ends with a conclusion.

Rudiments of Persistent Homology

As representations of topological features, the homology groups are abstract abelian groups, which may not be robust or able to provide continuous measurements. Thus, practical treatments of noisy data require the theory of persistent homology, which provides continuous measurements for the persistence of topological structures, allowing both quantitative comparison and noise removal in topological analyses. The concept was introduced by Frosini and Landi^[42] and Robins,^[43] and in the general form by Zomorodian and Carlsson.^[2] Computationally, the first efficient algorithm for Z/2 coefficient situation was proposed by Edelsbrunner et al.^[1] in 2002.

Simplex and simplicial complex

For discrete surfaces, that is, meshes, the commonly used homology is called simplicial homology. To describe this notion, we first present a formal description of the meshes,

the common discrete representation of surfaces and volumes. Essentially, meshing is a process in which a geometric shape is decomposed into elementary pieces called cells, the simplest of which are called *simplices*.

Simplex. *Simplices* are the simplest polytopes in a given dimension, as described below. Let v_0, v_1, \dots, v_p be $p + 1$ affinely independent points in a linear space. A p -simplex σ_p is the convex hull of those $p + 1$ vertices, denoted as $\sigma_p = \text{convex} \langle v_0, v_1, \dots, v_p \rangle$ or shorten as $\sigma_p = \langle v_0, v_1, \dots, v_p \rangle$. A formal definition can be given as,

$$\sigma_p = \left\{ v \mid v = \sum_{i=0}^p \lambda_i v_i, \sum_{i=0}^p \lambda_i = 1, 0 \leq \lambda_i \leq 1, \forall i \right\}. \quad (1)$$

The most commonly used simplices in \mathbb{R}^3 are 0-simplex (vertex), 1-simplex (edge), 2-simplex (triangle), and 3-simplex (tetrahedron) as illustrated in Figure 1.

An m -face of σ_p is the m -dimensional subset of $m + 1$ vertices, where $0 \leq m \leq p$. For example, an edge has two vertices as its 0-faces and one edge as its 1-face. As the number of subsets of a set with $p + 1$ vertices is 2^{p+1} , there are a total of $2^{p+1} - 1$ faces in σ_p . All the faces are proper except for σ_p itself. Note that polytope shapes can be decomposed into cells other than simplices, such as hexahedron and pyramid. However, as non-simplicial cells can be further decomposed, we can, without loss of generality, restrict our discussion to shapes decomposed to simplices as we describe next.

Simplicial Complex. With simplices as the basic building blocks, we define a *simplicial complex* K as a finite collection of simplices that meet the following two requirements,

- Containment: Any face of a simplex from K also belongs to K .
- Proper intersection: the intersection of any two simplices σ_i and σ_j from K is either empty or a face of both σ_i and σ_j .

Two p -simplices σ^i and σ^j are *adjacent* to each other if they share a common face. The boundary of σ_p , denoted as $\partial\sigma_p$, is the union (which can be written as a formal sum) of its $(p-1)$ -faces. Its interior is defined as the set containing all nonboundary points, denoted as $\sigma - \partial\sigma_p$. We define a boundary operator for each p -simplex spanned by vertices v_0 through v_p as

$$\delta p \langle v_0, \dots, v_p \rangle = \sum_{i=0}^p \langle v_0, \dots, \hat{v}_i, \dots, v_p \rangle, \quad (2)$$

where \hat{v}_i indicates that v_i is omitted and $Z/2$ coefficient set is used. It is the boundary operator that creates the nested topological structures and the homomorphism among them as described in the next section.

If the vertex positions in the ambient linear space can be ignored or do not exist, the containment relation among the simplices (as finite point sets) defines an *abstract simplicial complex*.

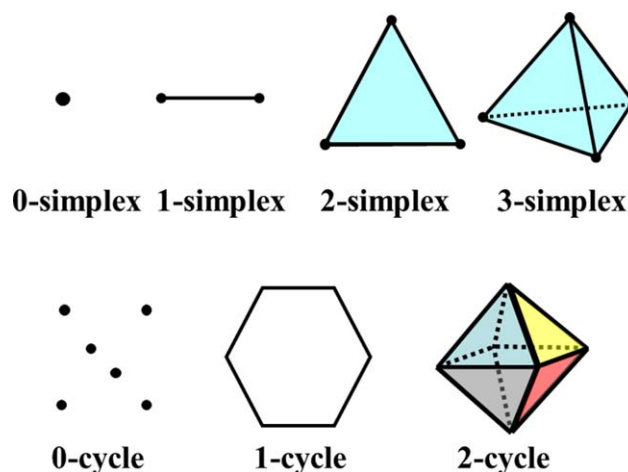


Figure 1. Illustration of 0-simplex, 1-simplex, and 2-simplex in the first row. The second row is simple 0-cycle, 1-cycle, and 2-cycle.

Homology

A powerful tool in topological analysis is homology, which represents certain structures in the meshes by algebraic groups to describe their topology. For regular objects in 3D space, essential topological features are connected components, tunnels and handles, and cavities, which are exactly described by the zeroth, first, and second homology groups, respectively.

Chains. The shapes to be mapped to homology groups are constructed from *chains* defined below. Given a simplicial complex (e.g., a tetrahedral mesh) K , which, roughly speaking, is a concatenation of p -simplices, we define a p -chain $c = \sum a_i \sigma_i$ as a formal linear combination of all p -simplices in K , where $a_i \in Z/2$ is 0 or 1 and σ_i is a p -simplex. Under such a definition, a 0-chain is a set of vertices, a 1-chain is a set of line segments which link vertices, a 2-chain is a set of triangles which are enclosed by line segments, and a 3-chain is a set of tetrahedrons which are enclosed by triangle surfaces.

We extend the boundary operator ∂_p for each p -simplex to a linear operator applied to chains, i.e., the extended operator meet following two conditions for linearity,

$$\begin{aligned} \partial_p(\lambda c) &= \lambda \partial_p(c), \\ \partial_p(c_i + c_j) &= \partial_p(c_i) + \partial_p(c_j), \end{aligned} \quad (3)$$

where c_i and c_j are both chains and λ is a constant, and all arithmetic is for modulo-2 integers, in which $1+1=0$.

An important property of the boundary operator is that the following composite operation is the zero map,

$$\partial_p \circ \partial_{p+1} = 0, \quad (4)$$

which immediately follows from the definition. Take the 2-chain $c = f_1 + f_2$ as an example, which represents a membrane formed by two triangles, $f_1 = \langle v_1, v_2, v_3 \rangle$ and $f_2 = \langle v_3, v_2, v_4 \rangle$. The boundary of c is a 1-chain, which turns out to be a loop,

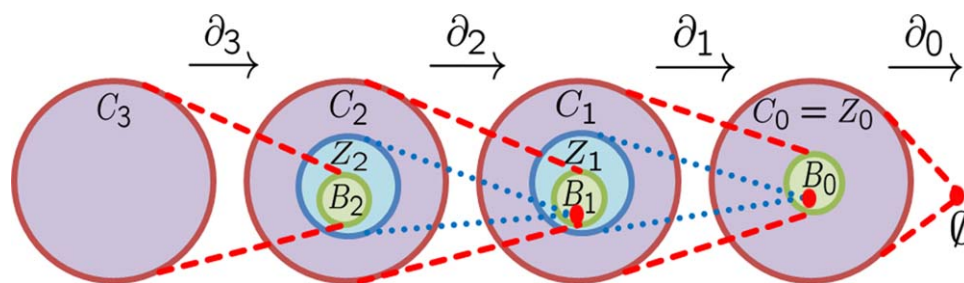


Figure 2. Illustration of boundary operators, and chain, cycle and boundary groups in \mathbb{R}^3 . Red dots stand for empty sets.

$$\begin{aligned} \partial_2(c) = & \langle v_1, v_2 \rangle + \langle v_2, v_3 \rangle + \langle v_3, v_1 \rangle + \langle v_3, v_2 \rangle \\ & + \langle v_2, v_4 \rangle + \langle v_4, v_3 \rangle = \langle v_1, v_2 \rangle \\ & + \langle v_3, v_1 \rangle + \langle v_2, v_4 \rangle + \langle v_4, v_3 \rangle . \end{aligned} \quad (5)$$

The boundary of this loop is thus

$$\begin{aligned} \partial_1 \circ \partial_2(c) = & \partial_1(\langle v_1, v_2 \rangle + \langle v_3, v_1 \rangle + \langle v_2, v_4 \rangle + \langle v_4, v_3 \rangle) \\ = & v_1 + v_2 + v_2 + v_4 + v_4 + v_3 + v_3 + v_1 = 0. \end{aligned} \quad (6)$$

Simplicial Homology. Simplicial homology is built on the *chain complex* associated to the simplicial complex. A chain complex is a sequence of abelian groups (C_1, C_2, \dots, C_n) connected by the homomorphism (linear operators) ∂_p , such that $\partial_p \circ \partial_{p+1} = 0$ as in eq.(4).

$$\dots \xrightarrow{\partial_{p+1}} C_p \xrightarrow{\partial_p} C_{p-1} \xrightarrow{\partial_{p-1}} \dots \xrightarrow{\partial_2} C_1 \xrightarrow{\partial_1} C_0 \xrightarrow{\partial_0} \emptyset. \quad (7)$$

The chain complex in the definition of simplicial homology is formed by C_p , the space of all p -chains, and ∂_p , the boundary operator on p -chains. As $\partial_p \circ \partial_{p+1} = 0$, the kernel of the boundary operator p -chains is a subset of the image of the boundary operator of $p+1$ -chains. The p -chains in the kernel of the boundary homomorphisms ∂_p are called p -cycles (p -chains without boundary) and the p -chains in the image of the boundary homomorphisms ∂_{p+1} are called p -boundaries. The p -cycles form an abelian group (with group operation being the addition of chains) called cycle group, denoted as $Z_p = \text{Ker } \partial_p$. The p -boundaries form another abelian group called boundary group, denoted as $B_p = \text{Im } \partial_{p+1}$.

Thus, p -boundaries are also p -cycles as shown in Figure 2. As p -boundaries form a subgroup of the cycles group, the quotient group can be constructed through cosets of p -cycles, that is, by equivalence classes of cycles. The p th homology, denoted as H_p , is defined as a quotient group,

$$\begin{aligned} H_p = & \text{Ker } \partial_p / \text{Im } \partial_{p+1} \\ = & Z_p / B_p, \end{aligned} \quad (8)$$

where $\text{Ker } \partial_p$ is the collection of p -chains with empty boundary and $\text{Im } \partial_{p+1}$ is the collection of p -chains that are boundaries of $p+1$ -chains.

Noticing that all groups with $p > 3$ cannot be generated from meshes in \mathbb{R}^3 , we only need chains, cycles and bounda-

ries of dimension p with $0 \leq p \leq 3$. See Figure 2 for an illustration.

We illustrate simplexes and cycles including 0-cycle, 1-cycle, and 2-cycle in Figure 1. Basically, an element in the p th homology group is an equivalence class of p -cycles. One of these cycles c can represent any other p -cycle that can be “deformed” through the mesh to c , because any other p -cycle in the same equivalence class differ with c by a p -boundary $b = \partial(\sigma_1 + \sigma_2 + \dots)$, where each σ_i is a $p+1$ -simplex. Adding the boundary of σ_i has the effect of deforming c to $c + \partial\sigma_i$ by sweeping through σ_i . For instance, a 0-cycle v_i is equivalent to v_j if there is a path $\langle v_i, v_{k1} \rangle + \langle v_{k1}, v_{k2} \rangle + \dots + \langle v_{kn}, v_j \rangle$. Thus each generator of zeroth-homology, (like a basis vector in a basis of the linear space of zeroth-homology) represents one connected component. Similarly, 1-cycles are loops, and first-homology generators represent independent nontrivial loops, that is, separate tunnels; 2-homology generators are independent membranes, each enclosing one cavity of the 3D object.

Define $\beta_p = \text{rank}(H_p)$ to be the p th Betti number. For a simplicial complex in 3D, β_0 is the number of connected components; β_1 is the number of tunnels; and β_2 is the number of cavities. As H_p is the quotient group between Z_p and B_p , we can also compute Betti numbers through,

$$\text{rank}(H_p) = \text{rank}(Z_p) - \text{rank}(B_p). \quad (9)$$

Note, however, H_p is usually of much lower rank than either Z_p or B_p .

Persistent homology

Homology generators identify the tunnels, cavities, and so forth, in the shape, but as topological invariants, they omit the metric measurements by definition. However, in practice, one often needs to compare the sizes of tunnels, for instance, to find the narrowest tunnel, or to remove tiny tunnels as topological noises. Persistent homology is a method of reintroducing metric measurements to the topological structures.^[1,2]

The measurement is introduced as an index i to a sequence of nested topological spaces $\{\mathbb{X}_i\}$. Such a sequence is called a *filtration*,

$$\emptyset = \mathbb{X}_0 \subseteq \mathbb{X}_1 \subseteq \mathbb{X}_2 \subseteq \dots \subseteq \mathbb{X}_m = \mathbb{X}. \quad (10)$$

As each inclusion induces a mapping of chains, it induces a linear map for homology,

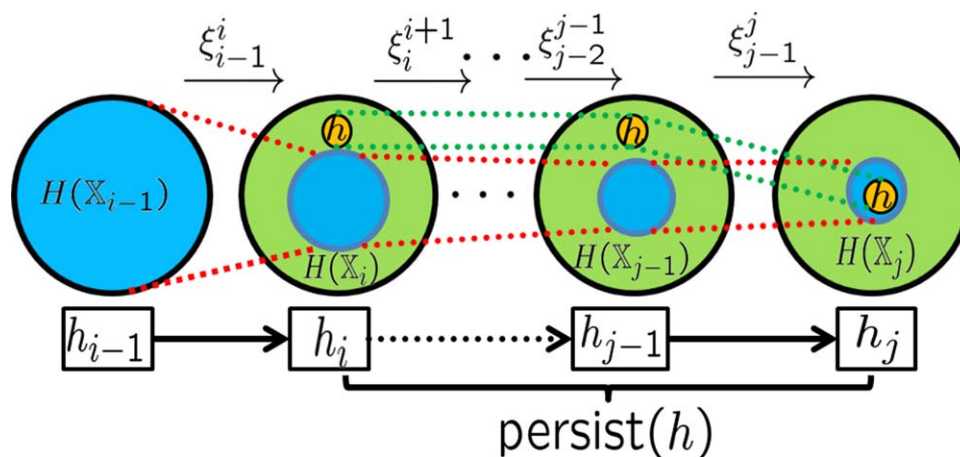


Figure 3. Illustration of the birth and death of a homology generator c . [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

$$\emptyset = H(\mathbb{X}_0) \rightarrow H(\mathbb{X}_1) \rightarrow H(\mathbb{X}_2) \rightarrow \dots \rightarrow H(\mathbb{X}_m) = H(\mathbb{X}). \quad (11)$$

The above sequence describes the evolution of the homology generators. We follow the exposition in Ref. [44] and define by a composition mapping from $H(\mathbb{X}_i)$ to $H(\mathbb{X}_j)$ as $\zeta_j^i : H(\mathbb{X}_i) \rightarrow H(\mathbb{X}_j)$. A new homology class c is created (born) in \mathbb{X}_i if it is not in the image of ζ_{i-1}^i . It dies in \mathbb{X}_j if it becomes trivial or is merged to an “older” (born before i) homology class, i.e., its image in $H(\mathbb{X}_j)$ is in the image of ζ_{i-1}^j , unlike its image under ζ_i^{j-1} .

As shown in Figure 3, if we associate with each space \mathbb{X}_i a value h_i denoting “time” or “length,” we can define the duration, or the persistence length of the each homology generator c as

$$\text{persist}(c) = h_j - h_i. \quad (12)$$

This measurement h_i is usually readily available when analyzing the topological feature changes. For instance, when the filtration arises from the level sets of a height function.

Algorithms for Persistent Homology

In computational topology, intrinsic features of point cloud data, that is, a point set $S \subset \mathbb{R}^n$ without additional structure, are common subjects of investigation. For such data, a standard way to construct the filtration is to grow a solid ball centered at each point with an ever-increasing radius. If the differences between points can generally be ignored, as is the case for fullerenes, a common radius r can be used for all points. In this setting, the radius r is used as the parameter for the family of spaces in the filtration. As the value of r increases, the solid balls will grow and simplices can be defined through the overlaps among the set of balls. In Figure 4, fullerene C_{60} is used to demonstrate this process. There are various ways of constructing abstract simplicial complexes from the intersection patterns of the set of expanding balls, such as Čech complex, Vietoris–Rips complex and alpha complex. The corresponding topological invariants, for example, the Betti numbers, are in general different due to different

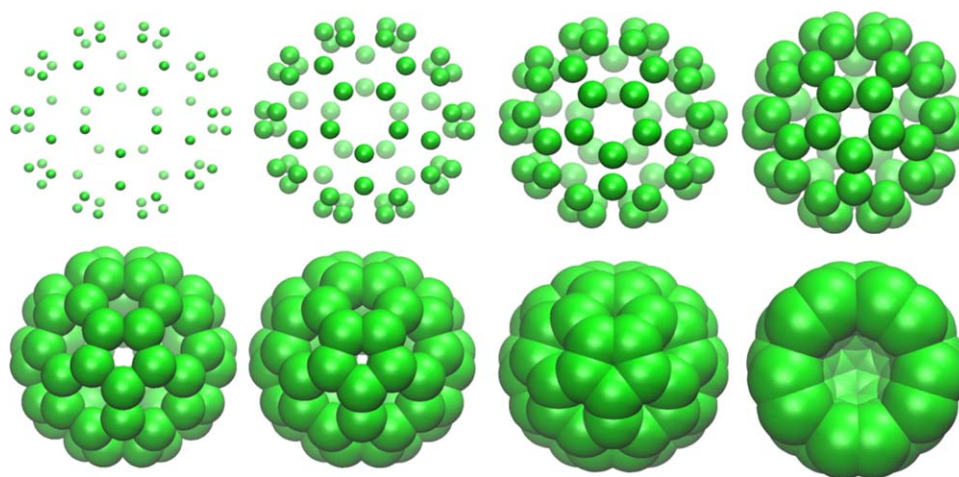


Figure 4. Illustration of filtrations built on fullerene C_{60} . Each point or atom in the point cloud data (i.e., coordinates) of the C_{60} is associated with a common radius r which increases gradually. As the value of r increases, the solid balls centered at given coordinates grow. These balls eventually overlap with their neighbors at certain r values. Simplices indicating such neighborhood information can be defined through abstract r -dependent simplicial complexes, e.g., alpha complexes and Rips complexes. Note that in the last chart, we have removed some atoms to reveal the central void. [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

definitions of simplicial complexes. In this section, we discuss computational algorithms for the design of filtrations, the construction of abstract simplicial complexes, and the calculation of Betti numbers.

Alpha complex

One possible filtration that can be derived from the unions of the balls with a given radius around the data points (as shown in Fig. 4) is the family of d -dependent Čech complexes, each of them is defined to be a simplicial complex, whose k -simplices are determined by $(k + 1)$ -tuples of points, such that the corresponding $d/2$ -balls have a nonempty intersection. However, it may contain many simplices for a large d . A variant called the alpha complex can be defined by replacing the $d/2$ -ball in the above definition by the intersection of the $d/2$ -ball with the Voronoi cells for these data points. In both cases, they are homotopic to the simple unions of balls, and thus produce the same persistent homology. Interested readers are referred to the nerve theorem for details.^[45]

Vietoris–Rips complex

The Vietoris–Rips complex, which is also known as Vietoris complex or Rips complex, is another type of abstract simplicial complex derived from the union of balls. In this case, for a k -simplex to be included, instead of requiring that the $(k + 1)$ $d/2$ -balls to have a common intersection, one only needs them to intersect pairwise. The Čech complex is a subcomplex of the Rips complex for any given d , however, the latter is much easier to compute and is also a subcomplex of the former at the filtration parameter of $\sqrt{2}d$.

Euclidean-distance based filtration

It is straightforward to use the metric defined by the Euclidean space in which the data points are embedded. The pairwise distance can be stored in a symmetric distance matrix (d_{ij}) , with each entry d_{ij} denoting the distance between point i and point j . Each diagonal term of the matrix is the distance from a certain point to itself, and thus is always 0. The family of Rips complexes is parameterized by d , a threshold on the distance. For a certain value of d , the Vietoris–Rips complex can be calculated. In 3D, more specifically, for a pair of points whose distance is below the threshold d , they form a 1-simplex in the Rips complex; for a triplet of points, if the distance between every pair is smaller than d , the 2-simplex formed by the triplet is in the Rips complex; whether a 3-simplex is in the Rips complex can be similarly determined. The Euclidean-distance based Vietoris–Rips complexes are widely used in persistent homology due to their simplicity and efficiency.

Correlation matrix based filtration

Another way to construct the metric space is through a certain correlation matrix, which can be built, for example, from theoretical predictions and experimental observations. From a previous study on protein stability, flexibility-rigidity index (FRI)

theory has been proven accurate and efficient.^[39] The reason for its success is that the geometric information is harnessed properly through the special transformation to a correlation matrix. The key to this transformation is the geometric to topological mapping. Instead of direct geometric information of the embedding in the Euclidean space, a mapping through certain kernel functions is able to interpret spatial locations of atoms in a particular way that reveals the atom stability quantitatively. We believe that this functional characterization is of importance to the study of not only proteins, but also other molecules.

Here, we present a special correlation matrix based Vietoris complex on the FRI method. To define the metric used, we briefly review the concepts of the FRI theory. First, we introduce the geometry to topology mapping.^[39–41] We denote the coordinates of atoms in the molecule we study as $\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_j, \dots, \mathbf{r}_N$, where $\mathbf{r}_j \in \mathbb{R}^3$ is the position vector of the j th atom. The Euclidean distance between i th and j th atoms r_{ij} can then be calculated. Based on these distances, topological connectivity matrix can be constructed with monotonically decreasing radial basis functions. A general form for a connectivity matrix is,

$$C_{ij} = w_j \Phi(r_{ij}, \eta_j), \quad (13)$$

where w_j is associated with atomic types, parameter $\eta_j > 0$ is the atom-type related characteristic distance, and $\Phi(r_{ij}; \eta_j)$ is a radial basis correlation kernel.

The choice of kernel is of significance to the FRI model. It has been shown that highly predictive results can be obtained by the exponential type and Lorentz type of kernels.^[39–41] Exponential type of kernels is

$$\Phi(r, \eta) = e^{-(r/\eta)^\kappa}, \quad \eta > 0, \kappa > 0 \quad (14)$$

and the Lorentz type of kernels is

$$\Phi(r, \eta) = \frac{1}{1 + (r/\eta)^v}, \quad \eta > 0, v > 0 \quad (15)$$

The parameters κ and v are adjustable.

We define the atomic rigidity index μ_i for i th atom as

$$\mu_i = \sum_{j=1}^N w_j \Phi(r_{ij}, \eta_j), \quad \forall i = 1, 2, \dots, N. \quad (16)$$

A related atomic flexibility index can be defined as the inverse of the atomic rigidity index.

$$f_i = \frac{1}{\mu_i}, \quad \forall i = 1, 2, \dots, N. \quad (17)$$

The FRI theory has been intensively validated by comparing with the experimental data, especially the Debye–Waller factor (commonly known as the B-factor).^[39] While simple to evaluate, their applications in B-factor prediction yield decent results. The predicted results are proved to be highly accurate while the procedure remains efficient. FRI is also used to analyze the protein folding behavior.^[41]

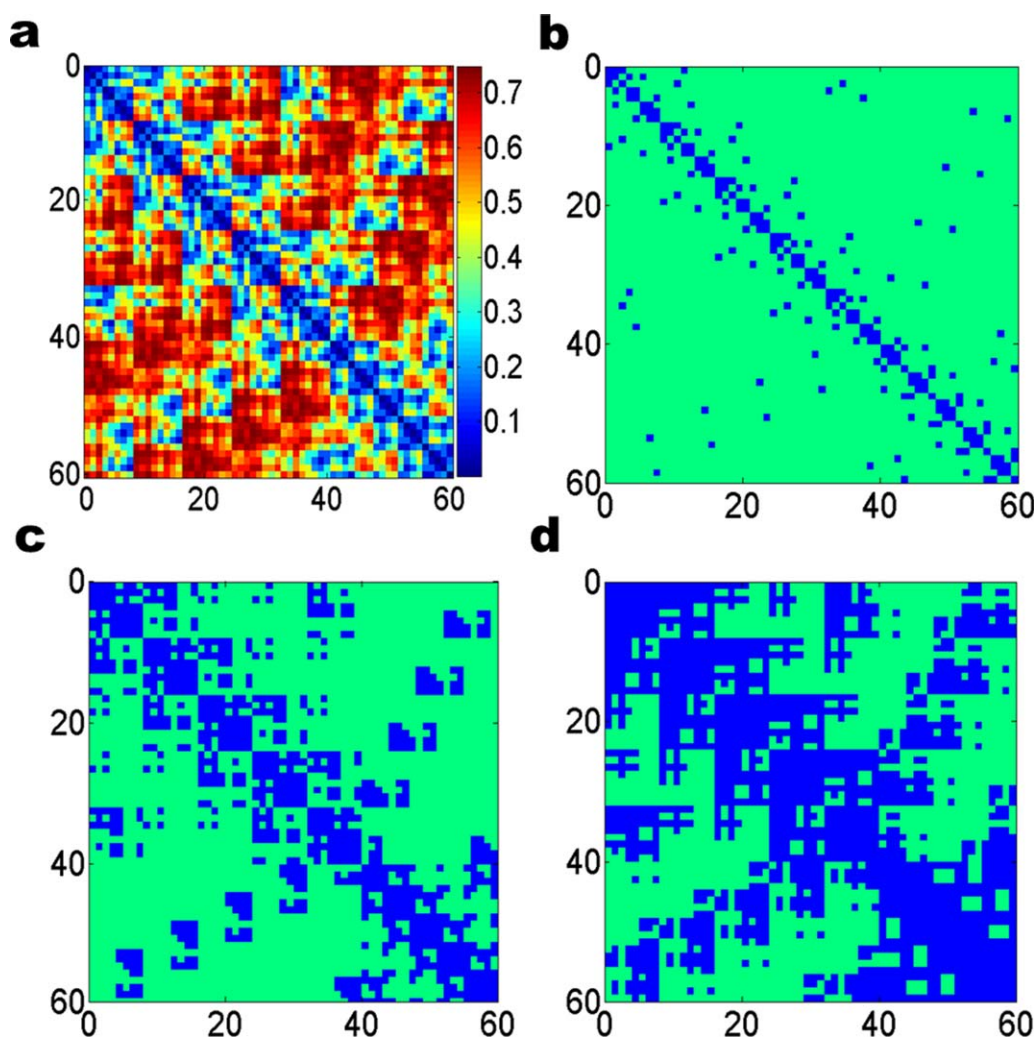


Figure 5. Correlation matrix based filtration of fullerene C_{60} (labels on both axes are atomic numbers). A correlation matrix is constructed from the FRI theory. As the filtration parameter increases, the Rips complex based on this matrix expands accordingly. a) The correlation based matrix for fullerene C_{60} ; b)–d) demonstrate the connectivity between atoms at the filtration threshold $d = 0.1, 0.3,$ and 0.5 \AA , respectively. The blue color entries represent the pairs already forming simplices. [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

To construct an FRI-based metric space, we need to design a special distance matrix, in which the functional correlation is measured. If we directly use the correlation matrix in eq. (13) for the filtration, atoms with less functional relation form more simplices, resulting in a counter-intuitive persistent homology. However, this problem can be easily remedied by defining a new correlation matrix as $M_{ij} = 1 - C_{ij}$, that is,

$$M_{ij} = 1 - w_j \Phi(r_{ij}, \eta_j). \quad (18)$$

Thus a kernel function induces a metric space under this definition. Figure 5a demonstrates such a metric space based filtration of fullerene C_{60} , in which we assume $w_j = 1$ as only one type of atom exists in this system. The generalized exponential kernel in eq. (14) is used with parameters $\kappa = 2.0$ and $\eta = 6.0 \text{ \AA}$.

With the correlation matrix based filtration, the corresponding Vietoris–Rips complexes can be straightforwardly constructed. Specifically, given a certain filtration parameter h_0 , if the matrix entry $M_{ij} \leq h_0$, an edge formed between i th and j th

atoms, and a simplex is formed if all of its edges are present. The complexes are built incrementally as the filtration parameter grows. Figures 5b–5d illustrate this process with three filtration threshold values $h = 0.1, 0.3,$ and 0.5 \AA , respectively. We use the blue color to indicate formed edges. It can be seen that simplicial complexes keep growing with the increase of filtration parameter h . The diagonal terms are always equal to zero, which means that N atom centers (0-simplices) form the first complex in the filtration.

Application to Fullerene Structure Analysis and Stability Prediction

In this section, the theory and algorithms of persistent homology are used to study the structure and stability of fullerene molecules. The ground-state structural data of fullerene molecules used in our tests are downloaded from the CCL webpage and fullerene isomer data and corresponding total curvature energies^[46] are adopted from David Tomanek's carbon

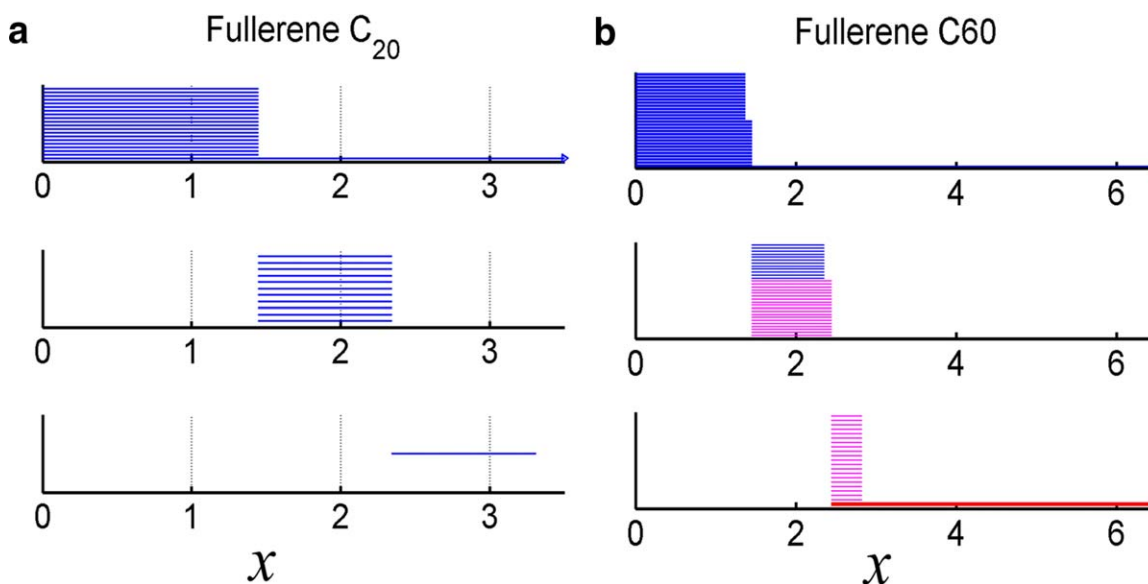


Figure 6. Illustration of the barcodes for fullerene C_{20} (left chart) and C_{60} (right chart) filtration on associated Rips complexes. Each chart contains three panels corresponding to the Betti number sequences β_0 , β_1 , and β_2 , from top to bottom. [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

fullerene webpage. In these structural data, coordinates of fullerene carbon atoms are given. The collection of atom center locations of each molecule forms a point cloud in \mathbb{R}^3 . With only one type of atom, the minor heterogeneity of atoms due to their chemical environments in these point clouds can be ignored in general. We examined both distance based and correlation matrix based metric spaces in our study. The filtration based on the FRI theory is shown to predict the stability very well.

Before we discuss the more informative persistent homology of fullerenes, we discuss the basic structural properties simply based on their Euler characteristics (vertex number minus edge number plus polygon number). The Euler characteristic, as a topological property, is invariant under nondegenerate shape deformation. For a fullerene cage composed of only pentagons and hexagons, the exact numbers of these two types of polygons can be derived from the Euler characteristic. For instance, if we have n_p pentagon and n_h hexagons in a C_N fullerene cage, the corresponding numbers of vertices, edges and faces are $(5n_p + 6n_h)/3$, $(5n_p + 6n_h)/2$ and $n_p + n_h$, respectively, as each vertex is shared by three faces, and each edge is shared by two faces. As the fullerene cage is treated as a two dimensional surface, we have the Euler characteristic $(5n_p + 6n_h)/3 - (5n_p + 6n_h)/2 + (n_p + n_h) = 2$, according to Euler's polyhedron formula, as it is a topological sphere. Thus, we have $n_p = 12$, which means a fullerene cage structure must have 12 pentagons and correspondingly $N/2 - 10$ hexagons. Therefore, for a C_N fullerene cage, we have N vertices, $3N/2$ edges and $N/2 + 2$ faces.

Barcode representation of fullerene structures and nanotube

Barcodes for Fullerene Molecule. In Figure 6, we demonstrate the persistent homology analysis of fullerene C_{20} and C_{60} using the barcode representation generated by Javaplex.^[9] The x-axis represents the filtration parameter h . If the distance between two vertices is below or equal to certain h_0 , they will

form an edge (1-simplex) at h_0 . Stated differently, the simplicial complex generated is equivalent to the radius filtration with radius parameter $h/2$. In the barcode, the persistence of a certain Betti number is represented by an interval (also known as bar), denoted as $L_i^{\beta_j}, j=0, 1, 2; i=1, 2, \dots$. Here, $j \in \{0, 1, 2\}$ as we only consider the first three Betti numbers in this work. From top to bottom, the behaviors of β_0 , β_1 , and β_2 are depicted in three individual panels. It is seen that as h grows, isolated atoms initialized as points will gradually grow into solid spheres with an ever-increasing radius. This phenomenon is represented by the bars in the β_0 panel. Once two spheres overlap with each other, one β_0 bar is terminated. Therefore, the bar length for the independent 0-th homology generator (connected component) c_i^0 , denoted as $L_i^{\beta_0} = \text{persist}(c_i^0)$, indicates the bond length information of the molecule. As can be seen from Figure 6, for fullerene C_{20} , all β_0 bar lengths are around 1.45 Å and the total number of components equals exactly to 20. Conversely, fullerene C_{60} has two different kinds of bars with lengths around 1.37 and 1.45 Å, respectively, indicating its two types of bond lengths.

More structure information is revealed as β_1 bars, which represent independent noncontractible 1-cycles (loops), emerge. It is seen in the fullerene C_{20} figure, that there are 11 equal-length β_1 bars persisting from 1.45 to 2.34 Å. As fullerene C_{20} has 12 pentagonal rings, the Euler characteristics for a 1D simplicial subcomplex (1-skeleton) can be evaluated from the Betti numbers,

$$n_{\text{vertex}} - n_{\text{edge}} = \beta_0 - \beta_1. \quad (19)$$

Here, β_0 , n_{vertex} , and n_{edge} are 1, 20, and 30, respectively. Therefore, it is easy to obtain that $\beta_1 = 11$ for fullerene C_{20} , as demonstrated in Figure 6. It should be noticed that all β_1 bars end at filtration value $h = 2.34$ Å, when five balls in each pentagon with their ever-increasing radii begin to overlap to form a pentagon surface.

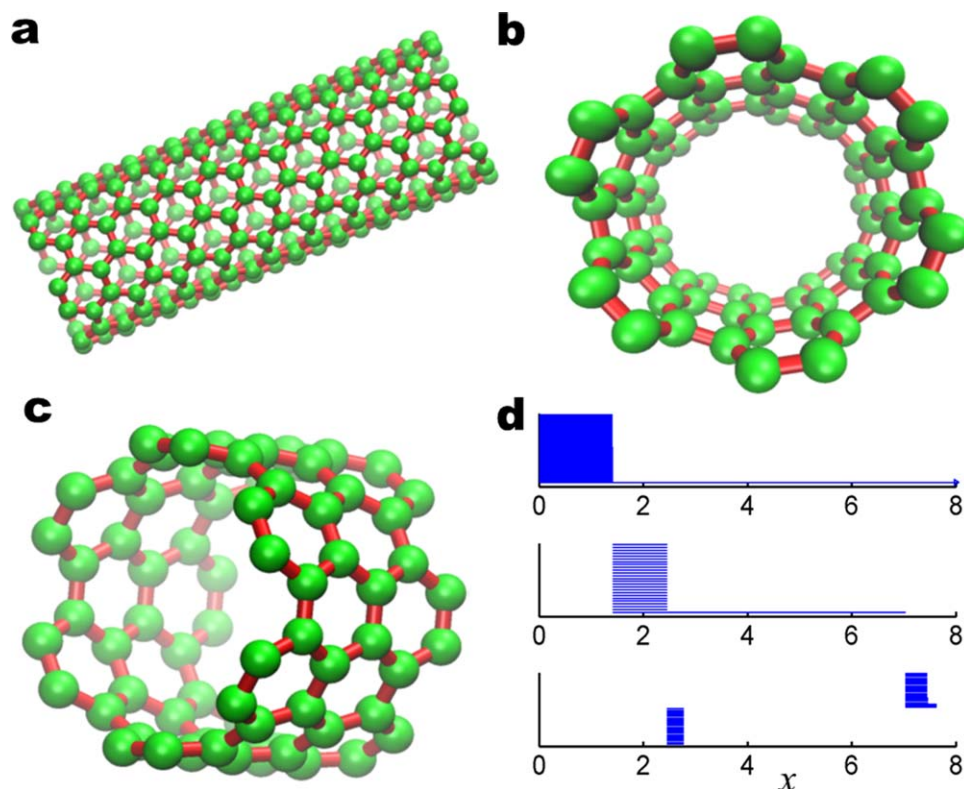


Figure 7. Illustration of persistent homology analysis for a nanotube. a) The generated nanotube structure with 10 unit layers. b) and c) A 3 unit layer segment extracted from the nanotube molecule in a). d) Barcodes representation of the topology of the nanotube segment. [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

Even more structural information can be derived from fullerene C_{60} 's β_1 barcodes. First, there are 31 bars for β_1 . This is consistent with the Euler characteristics in eq. (19), as we have 12 pentagons and 20 hexagons. Second, two kinds of bars correspond to the coexistence of pentagonal rings and hexagonal rings. They persist from 1.45 to 2.35 Å and from 1.45 to 2.44 Å, respectively.

As the filtration progresses, β_2 bars (membranes enclosing cavities) tend to appear. In fullerene C_{20} , there is only one β_2 bar, which corresponds to the void structure in the center of the cage. For fullerene C_{60} , we have 20 β_2 bars persisting from 2.44 to 2.82 Å, which corresponds to hexagonal cavities as indicated in the last chart of Figure 1. Basically, as the filtration goes, each node in the hexagon ring joins its four nearest neighbors, and fills in the relevant 2-simplices, yielding a simplicial complex whose geometric realization is exactly the octahedron. There is another β_2 bar due to the center void as indicated in the last chart of Figure 6, which persists until the complex forms a solid block. Note that two kinds of β_2 bars represent entirely different physical properties. The short-lived bars are related to local behaviors and fine structure details, while the long-lived bar is associated with the global feature, namely, the large cavity.

Barcodes for nanotube

Another example is a nanotube as demonstrated in Figure 7. The nanotube structure is constructed using the software TubeApplet webpage. We set tube indices to (6,6), the number

of unit cell to 10, tube radius to 4.05888 Å, and lattice constant to 2.454 Å. We extract a segment of three unit cells from the nanotube and use the persistent homology analysis to generate its barcodes. Our results are demonstrated in Figure 7. Different from fullerene molecules, the nanotube has a long β_1 bar representing the tube circle. It should also be noticed that β_2 barcodes are concentrated in two different regions. The first region is when x is around 2.5–2.7 Å. The β_2 barcodes in this domain are generated by hexagonal rings on the nanotube. The other region appears when x is slightly larger than 7.0 Å. The corresponding β_2 barcodes are representation of the void formed between different layer of carbons.

Unlike commonly used topological methods,^[31] persistent homology is able to provide a multiscale representation of the topological features. Usually, global behavior is of major concern. Therefore, the importance of the topological features is typically measured by their persistence length. In our analysis, we have observed that except for discretization errors, topological invariants of all scales can be equally important in revealing various structural features of the system of interest. In this work, we demonstrate that both local and global topological invariants play important roles in quantitative physical modeling.

Stability analysis of small fullerene molecules

From the above analysis, it can be seen that detailed structural information has been incorporated into the corresponding

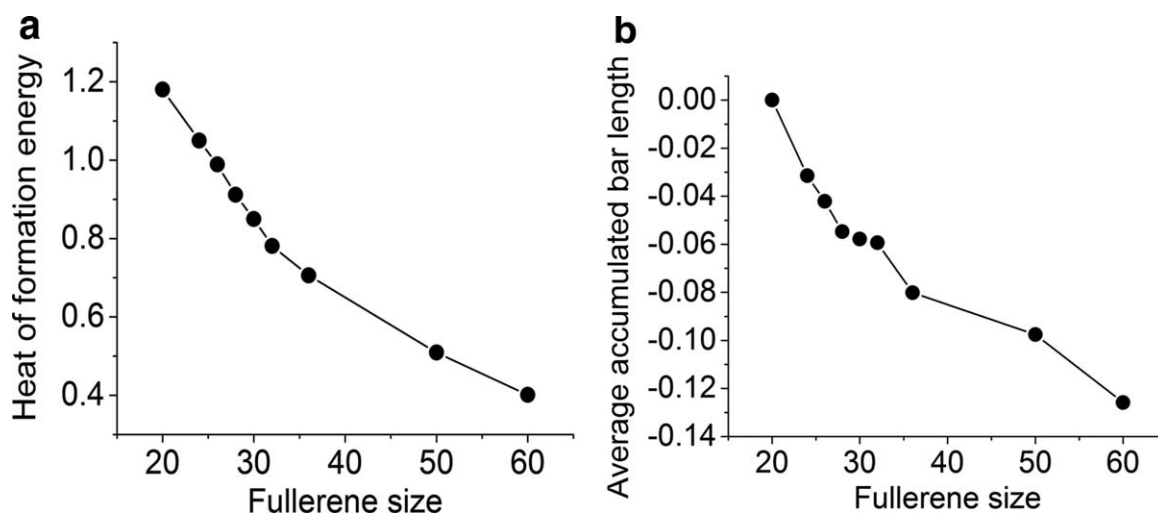


Figure 8. Comparison between the heat of formation energies computed using a quantum theory^[34] (left chart) and average accumulated bar length (right chart) for fullerenes. The units for the heat of formation energy and average accumulated bar length are eV/atom and Å/atom, respectively. Although the profile of average accumulated bar length of fullerenes does not perfectly match the fullerene energy profile, they bear a close resemblance in their basic characteristics.

barcodes. Conversely, molecular structures determine molecular functions.^[39–41] Therefore, persistent homology can be used to predict molecular functions of fullerenes. To this end, we analyze the barcode information. For each Betti number β_j , we define an accumulated bar length A_j as the summation of barcode lengths,

$$A_j = \sum_{i=1}^L L_i^j, j=0, 1, 2, \quad (20)$$

where L_i^j is the length of the i th bar in the j th-homology barcode. Sometimes, we may only sum over certain types of barcodes. We define an average accumulated bar length as $B_j = -\sum L_i^j/N$, where N is the total number of atoms in the molecule.

Zhang et al.^[34,38] found that for small fullerene molecule series C_{20} – C_{70} , their ground-state heat of formation energies gradually decrease with the increase of the number of atoms, except for C_{60} and C_{70} . The decreasing rate, however, slows down with the increase of the number of atoms. With data adopted from Ref. [34], Figure 8 demonstrates this phenomenon. This type of behavior is also found in the total energy (STO-3G/SCF at MM3) per atom,^[47] and in average binding energy of fullerene C_{2n} which can be broken down to n dimmers (C_2).^[48]

To understand this behavior, many theories have been proposed. Zhang et al.^[38] postulate that the fullerene stability is related to the ratio between the number of pentagons and the number of atoms for a fullerene molecule. Higher percentage of pentagon structures results in relatively higher levels of the heat of formation. Conversely, a rather straightforward isolated pentagon rule (IPR) states that the most stable fullerenes are those in which all the pentagons are isolated. The IPR explains why C_{60} and C_{70} are relatively stable as both have only isolated pentagons. Raghavachari's neighbour index^[49] provides another approach to quantitatively characterize the

relative stability. For example, in C_{60} of I_h symmetry, all 12 pentagons have neighbour index 0, thus the I_h C_{60} structure is very stable.

In this work, we hypothesize that fullerene stability depends on the average number of hexagons per atom. The larger number of hexagons is in a given fullerene structure, the more stable it is. We utilize persistent homology to verify our hypothesis. As stated in Section Barcode representation of fullerene structures and nanotube, there are two types of β_2 bars, namely, the one due to hexagon-structure-related holes and that due to the central void. Their contributions to the heat of formation energy are dramatically different. Based on our hypothesis, we only need to include those β_2 bars that are due to hexagon-structure-related holes in our calculation of the average accumulated bar length B_2 . As depicted in the right chart of Figure 8, the profile of the average accumulated bar length closely resembles that of the heat of formation energy. Instead of a linear decrease, both profiles exhibit a quick drop at first, then the decreasing rate slows down gradually. Although our predictions for C_{30} and C_{32} fullerenes do not match the corresponding energy profile precisely, which may be due to the fact that the data used in our calculation may not be exactly the same ground-state data as those in the literature,^[38] the basic characteristics and the relative relations in the energy profile are still well preserved. In fact, the jump at the C_{60} fullerene is captured and stands out more obviously than the energy profile. This may be due to the fact that our method distinguishes not only pentagon and hexagon structures, but also the size differences within each of them. We are not able to present the full set of energy data in Ref. [34] because we are limited by the availability of the ground-state structure data.

To quantitatively validate our prediction, the least squares method is used to fit our prediction with the heat of formation energy, and a correlation coefficient is defined,^[39]

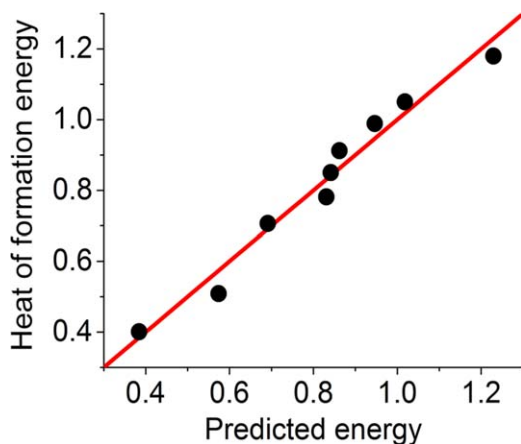


Figure 9. The comparison between quantum mechanical simulation results^[34] and persistent homology prediction of the heat of formation energy (eV/atom). Only local β_2 bars that are due to hexagon structures are included in our average accumulated bar length B_2 . The correlation coefficient from the least-squares fitting is near perfect ($C_c=0.985$). [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

$$C_c = \frac{\sum_{i=1}^N (B_i^e - \bar{B}^e)(B_i^t - \bar{B}^t)}{\left[\sum_{i=1}^N (B_i^e - \bar{B}^e)^2 \sum_{i=1}^N (B_i^t - \bar{B}^t)^2 \right]^{1/2}}, \quad (21)$$

where B_i^e represents the heat of formation energy of the i th fullerene molecule, and B_i^t is our theoretical prediction. The parameter \bar{B}^e and \bar{B}^t are the corresponding mean values. The fitting result is demonstrated in Figure 9. The correlation coefficient is close to unity (0.985), which indicates the soundness of our model and the power of persistent homology for quantitative predictions.

Total curvature energy analysis of fullerene isomers

Having demonstrated the ability of persistent homology for the prediction of the relative stability of fullerene molecules, we further illustrate the effectiveness of persistent homology for analyzing the total curvature energies of fullerene isomers. Fullerene molecules C_N are well known to admit various isomers,^[50] especially when the number (N) of atoms is large. To identify all of the possible isomers for a given N , many elegant mathematical algorithms have been proposed. Coxeter's construction method^[29,35] and the ring spiral method^[30] are two popular choices. Before discussing the details of these two methods, we need to introduce the concept of fullerene dual. Mathematically, a dual means dimension-reversing dual. From Euler's polyhedron theorem, if a spherical polyhedron is composed of n_{vertex} vertices, n_{edge} edges and n_{face} faces, we have the relation $n_{\text{vertex}} - n_{\text{edges}} + n_{\text{face}} = 2$. Keeping the n_{edge} unchanged while swapping the other two counts, we have its dual, which has n_{vertex} faces and n_{face} vertices. For example, the cube and the octahedron form a dual pair, the dodecahedron and the icosahedron form another dual pair, and the tetrahedron is its self-dual. This duality is akin to the duality

between the Delaunay triangulation and the corresponding Voronoi diagram in computational geometry.

In fullerenes, each vertex is shared by three faces (each of them is either a pentagon or a hexagon). Therefore, fullerene dual can be represented as a triangulation of the topological sphere. Based on this fact, Coxeter is able to analyze the icosahedral triangulations of the sphere and predict the associated isomers. This method, although mathematically rigorous, is difficult to implement for structures with low symmetry, thus is inefficient in practical applications.^[31] Conversely, in the Schlegel diagram,^[51] each fullerene structure can be projected into a planar graph made of pentagons and hexagons. The ring spiral method is developed based on the spiral conjecture,^[31] which states "The surface of a fullerene polyhedron may be unwound in a continuous spiral strip of edge-sharing pentagons and hexagons such that each new face in the spiral after the second shares an edge with both (a) its immediate predecessor in the spiral and (b) the first face in the preceding spiral that still has an open edge." Basically, for fullerenes of N atoms, one can list all possible spiral sequences of pentagons and hexagons, and then wind them up into fullerenes. If no conflict happens during the process, an isomer is generated. Otherwise, we neglect the spiral sequence. Table 1 lists the numbers of isomers for different fullerenes,^[31] when enantiomers are regarded as equivalent 1. It is seen that the number of isomers increases dramatically as N increases. Total curvature energies of many fullerene isomers are available at the carbon fullerene webpage.

In 1935, Hakon defined sphericity as a measure of how spherical (round) an object is.^[52] By assuming particles having the same volume but differing in surface areas, Hakon came up with a sphericity function,^[52]

$$\Psi = \frac{\pi^{1/3}(6V_p)^{2/3}}{A_p}, \quad (22)$$

where V_p and A_p are the volume and the surface area of the particle. Obviously, a sphere has sphericity 1, while the sphericity of nonspherical particles is less than 1. Let us assume that fullerene isomers have the same surface area as the perfect sphere $A_p = 4\pi R^2$, we define a sphericity measure as

$$\Psi_c = \frac{V_p}{V_s} = \frac{6\pi^{1/2}V_p}{A_p^{3/2}}, \quad (23)$$

where V_s is the volume of a sphere with radius R . By the isoperimetric inequality, among all simple closed surfaces with given surface area A_p , the sphere encloses a region of maximal volume. Thus, the sphericity of nonspherical fullerene isomers is less than 1. Consequently, in a distance based filtration

Table 1. Numbers of isomers for small fullerenes.

N_{atom}	20	24	26	28	30	32	34	36	38	40	50	60
N_{isomer}	1	1	1	2	3	6	6	15	17	40	271	1812

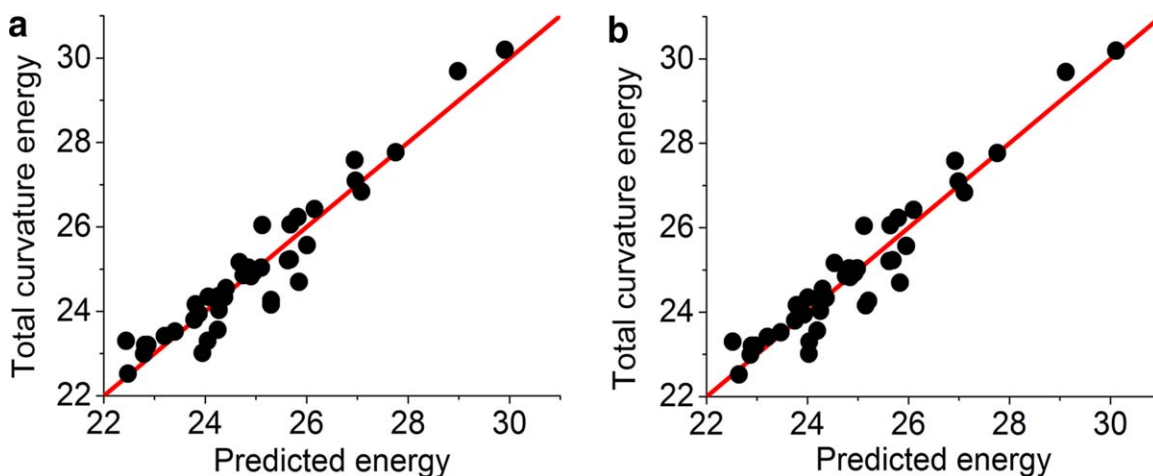


Figure 10. Comparison between the distance filtration (left chart) and the correlation matrix filtration (right chart) in fullerene C_{40} stability analysis. Fullerene C_{40} has 40 isomers. Each of them has an associated total curvature energy (eV). We calculate our average accumulated bar lengths from both distance filtration and the correlation matrix based filtration, and further fit them with total curvature energies. The correlation coefficients for our fitting are 0.956 and 0.959, respectively. It should be noticed that only the central void related β_2 bars (i.e., the long-lived bars) are considered. The exponential kernel is used in matrix filtration with parameter $\eta = 4$ and $\kappa = 2$. [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

process, the smaller sphericity a fullerene isomer is, the shorter its global β_2 bar will be.

On fullerene surface, the local curvature characterizes the bond bending away from the plane structure required by the $sp^{[2]}$ hybrid orbitals.^[53] Therefore, the relation between fullerene curvature and stability can be established and confirmed using *ab initio* density functional calculations.^[46] However, such an analysis favors fullerenes with infinitely many atoms. Let us keep the assumption that for a given fullerene C_N , all its isomers have the same surface area. We also assume that the most stable fullerene isomer C_N is the one that has a near perfect spherical shape. Therefore, each fullerene isomer is subject to a (relative) total curvature energy E_c per unit area due to its accumulated deviations from a perfect sphere,

$$E_c = \int_{\Gamma} \mu \left[(\kappa_1 - \kappa_0)^2 + (\kappa_2 - \kappa_0)^2 \right] dS \quad (24)$$

$$= \int_{\Gamma} 2\mu \left[\frac{1}{2} (2\mathbf{H} - \kappa_0)^2 + \mathbf{K} \right] dS, \quad (25)$$

where Γ is the surface, μ is bending rigidity, κ_1 and κ_2 are the two principal curvatures, and $\kappa_0 = 1/R$ is the constant curvature of the sphere with radius R . Here, \mathbf{H} and \mathbf{K} are the mean and Gaussian curvature of the fullerene surface, respectively. Therefore, a fullerene isomer with a smaller sphericity will have a higher total curvature energy. Based on the above discussions, we establish the inverse correlation between fullerene isomer global β_2 bar lengths and fullerene isomer total curvature energies.

Obviously, the present fullerene curvature energy (24) is a special case of the Helfrich energy functional for elasticity of cell membranes^[54]

$$E_c = \int_{\Gamma} \left[\frac{1}{2} \mathcal{K}_C (2\mathbf{H} - C_0)^2 + \mathcal{K}_G \mathbf{K} \right] dS, \quad (26)$$

where, C_0 is the spontaneous curvature, and \mathcal{K}_C and \mathcal{K}_G are the bending modulus and Gaussian saddle-splay modulus,

respectively. The Gauss–Bonnet theorem states that for a compact two-dimensional Riemannian manifold without boundary, the surface integral of the Gaussian curvature is $2\pi\chi$, where χ is the Euler characteristic. Therefore, the curvature energy admits a jump whenever there is a change in topology which leads to a change in the Euler characteristic. A problem with this discontinuity in the curvature energy is that the topological change may be induced by an infinitesimal change in the geometry associated with just an infinitesimal physical energy, which implies that the Gaussian curvature energy functional is unphysical. Similarly, Hadwiger type of energy functionals, which make use of a linear combination of the surface area, surfaced enclosed volume, and surface integral of mean curvature and surface integral of Gaussian curvature,^[55] may be unphysical as well for systems involving topological changes. However, this is not a problem for differential geometry based multiscale models which utilize only surface area and surface enclosed volume terms,^[56–59] as we use the Eulerian representation and the proposed generalized mean curvature terms but not Gaussian curvature terms. Moreover, in the present model for fullerene isomers, there is no topological change.

To verify our assumptions, we consider a family of isomers for fullerene C_{40} . It has a total of 40 isomers. We compute the global β_2 bar lengths of all isomers by Euclidean distance filtration and fit their values with their total curvature energies with a negative sign. Figure 10 (right chart) shows an excellent correlation between the fullerene total curvature energies and our persistent homology based predictions. The correlation coefficient is 0.956, which indicates that the proposed persistent homology analysis of nonsphericity and our assumption of a constant surface area for all fullerene isomers are sound. In reality, fullerene isomers may not have an exactly constant surface area because some distorted bonds may have a longer bond length. However, the high correlation coefficient found in our persistent homology analysis implies that either the average bond lengths for all isomers are similar or the error due to nonconstant surface area is offset by other errors.

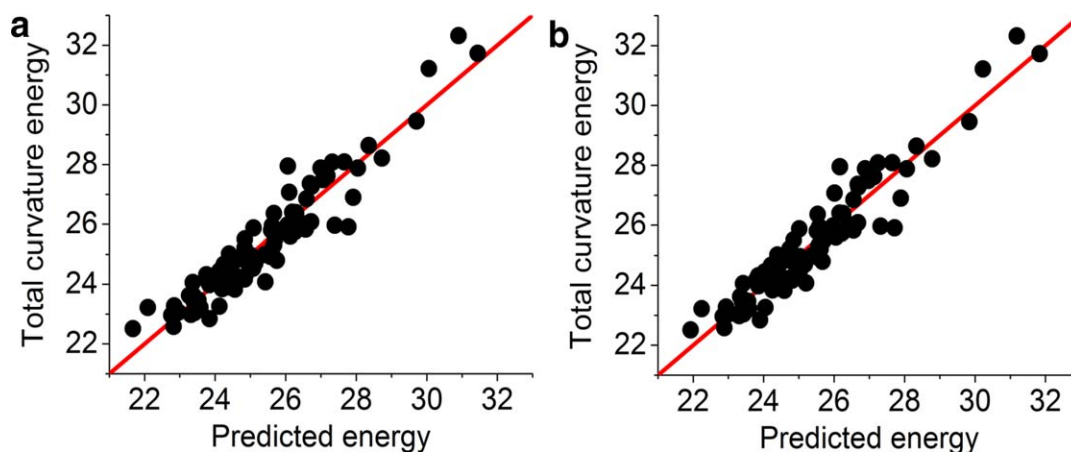


Figure 11. Further validation of our method with 89 isomers for fullerene C_{44} . The correlation coefficients for distance filtration (left chart) and correlation matrix based filtration (right chart) are 0.948 and 0.952, respectively. In the latter method, the exponential kernel is used with parameter $\eta = 4$ and $\kappa = 2$. [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

To further validate our persistent homology based method for the prediction of fullerene total curvature energies, we consider a case with significantly more isomers, namely, fullerene C_{44} , which has 89 isomers. In this study, we have again found an excellent correlation between the fullerene total curvature energies and our persistent homology based predictions as depicted in the right chart of Figure 11. The correlation coefficient for this case is 0.948. In fact, we have checked more fullerene isomer systems and obtained similar predictions.

Finally, we explore the utility of our correlation matrix based filtration process for analysis of fullerene total curvature energies. In place of Euclidean distance based filtration, the correla-

tion matrix based filtration is used. To demonstrate the basic principle, eq. (18) with the generalized exponential kernel in eq. (14) is used in the filtration. We assume $w_j = 1$ as fullerene molecules have only carbon atoms. To understand the correlation matrix based filtration method, the fullerene C_{60} is used again. We fixed the power $\kappa = 2$, and adjust the value of characteristic distance η . Figure 12 gives the calculated barcodes with $\eta = 2$ and $\eta = 20$. It can be seen that these barcodes share a great similarity with the Euclidean distance based filtration results depicted in the right chart of Figure 6. All of topological features, namely, two kinds of bonds in β_0 , the pentagonal rings and the hexagonal rings in β_1 , and also the hexagonal cavities and the central void in β_2 are clearly

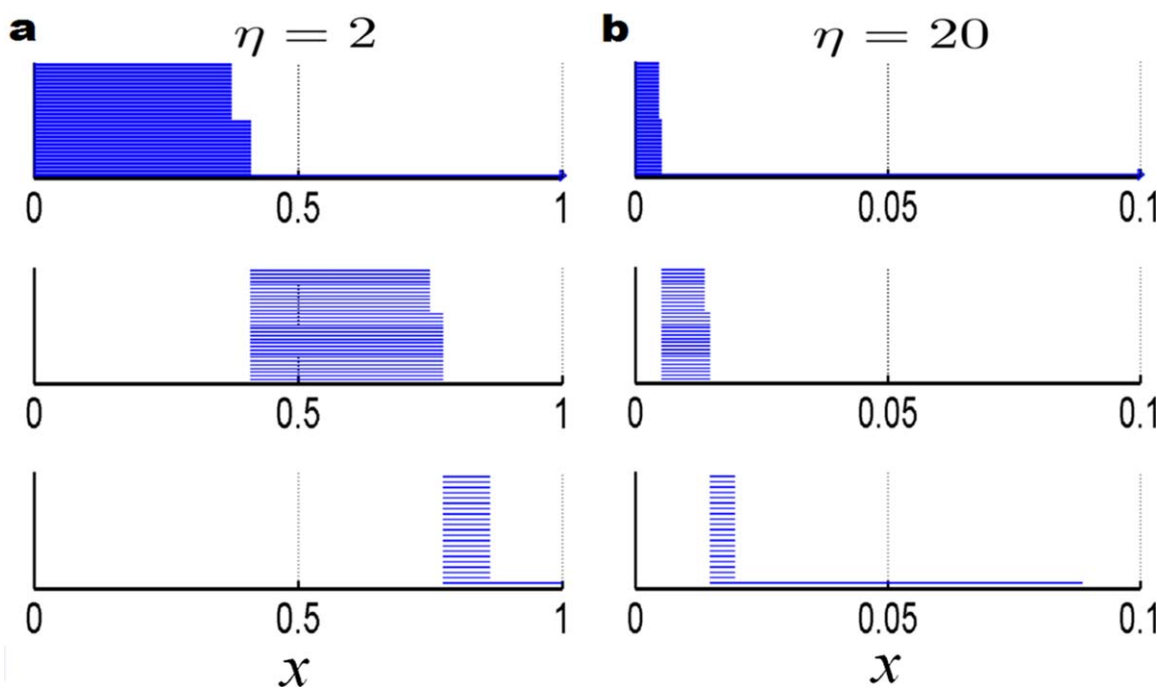


Figure 12. Illustration of the persistent barcodes generated using correlation matrix based filtrations with different characteristic distances. The exponential kernel model with power $\kappa = 2$ is used. The characteristic distances in the left and right charts are respectively $\eta = 2$ and $\eta = 20$. [Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.]

demonstrated. However, it should be noticed that, unlike the distance based filtration, the matrix filtration does not generate linear Euclidean distance relations. However, relative correspondences within the structure are kept. For instances, in β_2 bars, the bar length ratio between the central void part and the hexagonal hole part in Figure 12 is drastically different from its counterpart in Figure 6. From our previous experience in flexibility and rigidity analysis,^[39–41] these rescaled distance relations have a great potential in capturing the essential physical properties, such as, flexibility, rigidity, stability, and compressibility of the underlying system.

Similarly, the global β_2 bar lengths obtained from the correlation matrix based filtration are utilized to fit with the total curvature energies of fullerene isomers. The correlation coefficients for the correlation distance matrix filtration are 0.959 and 0.952, respectively for C_{40} and C_{44} fullerene isomers. The corresponding results are demonstrated in the right charts of Figures 10 and 11, respectively. It can be seen that the correlation matrix filtration is able to capture the essential stability behavior of fullerene isomers. In fact, results from correlation matrix based filtrations are slightly better than those of Euclidean distance based filtrations. In correlation matrix based filtrations, the generalized exponential kernel is used with parameter $\eta = 4$ and $\kappa = 2$. These parameters are chosen based on our previous flexibility and rigidity analysis of protein molecules. Overall, best prediction is obtained when the characteristic distance is about 2–3 times of the bond length and power index κ is around 2–3. Fine tuning of the parameters for each single case may yield even better result. However, this aspect is beyond the scope of the present work.

Conclusion

Persistent homology is an efficient tool for the qualitative analysis of topological features that last over scales. In the present work, for the first time, persistent homology is introduced for the quantitative prediction of fullerene energy and stability. We briefly review the principal concepts and algorithms in persistent homology, including simplex, simplicial complex, chain, filtration, persistence, and pairing algorithms. Euler characteristics analysis is used to decipher the barcode representation of fullerene C_{20} and C_{60} . A thorough understanding of fullerene barcode origins enables us to construct physical models based on local and/or global topological invariants and their accumulated persistent lengths. By means of an average accumulated bar length of the second Betti number that corresponds to fullerene hexagons, we are able to accurately predict the relative energies of a series of small fullerenes. To analyze the total curvature energies of fullerene isomers, we propose to use sphericity to quantify the nonspherical fullerene isomers and correlate the sphericity with fullerene isomer total curvature energies, which are defined as a special case of the Helfrich energy functional for elasticity. Topologically, the sphericity of a fullerene isomer is measured by its global second homology bar length in the barcode, which in turn gives rise to the prediction of fullerene isomer total curvature energies. We demonstrate an excellent agreement between total curvature

energies and our persistent homology predictions for the isomers of fullerene C_4 and C_{44} . Finally, a new filtration based on the correlation matrix of the flexibility and rigidity index is proposed and found to provide even more accurate predictions of fullerene isomer total curvature energies.

Acknowledgment

G.W.W. acknowledges the Mathematical Biosciences Institute for hosting valuable workshops. K.L.X. thanks Bao Wang for useful discussions.

Keywords: persistent homology · filtration · fullerene · isomer · nanotube · stability · curvature

How to cite this article: K. Xia, X. Feng, Y. Tong, G. W. Wei. *J. Comput. Chem.* **2015**, *36*, 408–422. DOI: 10.1002/jcc.23816

- [1] H. Edelsbrunner, D. Letscher, A. Zomorodian, *Discrete Comput. Geom.* **2002**, *28*, 511.
- [2] A. Zomorodian, G. Carlsson, *Discrete Comput. Geom.* **2005**, *33*:249.
- [3] H. Edelsbrunner, E. P. Mucke, *Phys. Rev. Lett.* **1994**, *13*, 43.
- [4] H. Edelsbrunner, J. Harer, *Computational Topology: An Introduction*; American Mathematical Society, Providence, RI, **2010**.
- [5] T. K. Dey, K. Y. Li, J. Sun, C. S. David, *ACM Trans. Graph.*, **2008**, *45*, 1.
- [6] T. K. Dey, Y. S. Wang, *Discrete Comput. Geom.* **2013**, *49*, 46.
- [7] K. Mischaikow, V. Nanda, *Discrete Comput. Geom.* **2013**, *50*, 330.
- [8] R. Ghrist, *Bull. Am. Math. Soc.* **2008**, *45*, 61.
- [9] A. Tausz, M. Vejdemo-Johansson, H. Adams, *Javaplex: A Research Software Package for Persistent (co)Homology*. Software Available at: <http://code.google.com/p/javaplex>, 2011.
- [10] G. Carlsson, T. Ishkhanov, V. Silva, A. Zomorodian, *Int. J. Comput. Vision* **2008**, *76*, 1.
- [11] D. Pachauri, C. Hinrichs, M. Chung, S. Johnson, V. Singh, *IEEE Trans. Med. Imaging* **2011**, *30*, 1760.
- [12] G. Singh, F. Memoli, T. Ishkhanov, G. Sapiro, G. Carlsson, D. L. Ringach, *J. Vision* **2008**, *11*, 1.
- [13] K. Mischaikow, M. Mrozek, J. Reiss, A. Szymczak, *Phys. Rev. Lett.* **1999**, *82*, 1144.
- [14] T. Kaczynski, K. Mischaikow, M. Mrozek, *Computational Homology*; Springer-Verlag, New York, NY, **2004**.
- [15] V. D. Silva, R. Ghrist, In *Proceedings of Robotics: Science and Systems*, Cambridge, USA, 2005, p. 1.
- [16] H. Lee, H. Kang, M. K. Chung, B. Kim, D. S. Lee, *IEEE Trans. Med. Imaging* **2012**, *31*, 2267.
- [17] D. Horak, S. Maletic, M. Rajkovic, *J. Stat. Mech.: Theory Experiment* **2009**, *2009*, P03034.
- [18] G. Carlsson, *Am. Math. Soc.* **2009**, *46*, 255.
- [19] X. Feng, Y. Tong, *IEEE Trans. Visual. Comput. Graphics* **2013**, *19*, 1298.
- [20] P. M. Kasson, A. Zomorodian, S. Park, N. Singhal, L. J. Guibas, V. S. Pande, *Bioinformatics* **2007**, *23*, 1753.
- [21] M. Gameiro, Y. Hiraoka, S. Izumi, M. Kramar, K. Mischaikow, V. Nanda, *Jpn. J. Ind. Appl. Math.*, **2014**, DOI: 10.1007/s13160-014-0153-5.
- [22] Y. Dabaghian, F. Memoli, L. Frank, G. Carlsson, *PLoS Comput. Biol.* **2012**, *8*, e1002581.
- [23] Y. Yao, J. Sun, X. H. Huang, G. R. Bowman, G. Singh, M. Lesnick, L. J. Guibas, V. S. Pande, G. Carlsson, *J. Chem. Phys.* **2009**, *130*, 144115.
- [24] K. L. Xia, G. W. Wei, *Int. J. Num. Methods Biomed. Eng.* **2014**, *30*, 814.
- [25] A. Adcock, E. Carlsson, G. Carlsson, *The Ring of Algebraic Functions on Persistence Bar Codes*, arXiv: 1304.0530, 2013.
- [26] P. Bendich, S. Chin, J. Clarke, J. deSena, J. Harer, E. Munch, A. Newman, D. Porter, D. Rouse, N. Strawn, A. Watkins, *Topological and Statistical Behavior Classifiers for Tracking Applications*, arXiv:1406.0214, 2014.
- [27] H. W. Kroto, J. R. Heath, S. C. O'Brien, R. F. Curl, R. E. Smalley, *Nature* **1985**, *318*, 162.

- [28] W. Kratschmer, L. D. Lamb, K. Fostiropoulos, D. Huffman, *Nature* **1990**, 347, 354.
- [29] P. W. Fowler, J. E. Cremona, J. I. Steer, *Theor. Chim. Acta*, **1988**, 73, 1.
- [30] D. E. Manolopoulos, J. C. May, S. E. Down, *Chem. Phys. Lett.* **1991**, 181, 105.
- [31] P. W. Fowler, D. E. Manolopoulos, *An Atlas of Fullerenes*; Clarendon Press, Oxford, **1995**.
- [32] P. Ballone, P. Milani, *Phys. Rev. B* **1990**, 42, 3201.
- [33] J. R. Chelikowsky, *Phys. Rev. Lett.* **1991**, 67, 2970.
- [34] B. L. Zhang, C. Z. Wang, K. M. Ho, C. H. Xu, C. T. Chan, *J. Chem. Phys.* **1992**, 97, 5007.
- [35] H. S. Coxeter, *Virus Macromolecules and Geodesic Domes*. In H.S. Coxeter and J.C. Butcher, Eds. *A Spectrum of Mathematics*; Auckland University Press, Auckland, **1971**, pp. 98–107.
- [36] A. Streitwieser, *Molecular Orbital Theory for Organic Chemists*; Wiley, New York, **1961**.
- [37] P. W. Fowler, S. J. Austin, D. E. Manolopoulos, *Phys. Chem. Fullerenes* **1994**, 443, 41.
- [38] B. L. Zhang, C. H. Xu, C. Z. Wang, C. T. Chan, K. M. Ho, *Phys. Rev. B* **1992**, 46, 11.
- [39] K. L. Xia, K. Opron, G. W. Wei, *J. Chem. Phys.* **2013**, 139, 194109.
- [40] K. L. Xia, G. W. Wei, *Phys. Rev. E* **2013**, 88, 062709.
- [41] K. L. Xia, G. W. Wei, *Chaos* **2014**, 24, 013103.
- [42] P. Frosini, C. Landi, *Pattern Recognit. Image Anal.* **1999**, 9, 596.
- [43] V. Robins, In *Topology Proceedings*, Alabama, USA, Vol 24; 1999, pp. 503–532.
- [44] E. Munch, *Applications of Persistent Homology to Time Varying Systems*, Dissertation of Duke University, Durham, NC, **2013**.
- [45] A. Zomorodian, *Topology for Computing*; Cambridge Monographs on Applied and Computational Mathematics, 2009.
- [46] J. Guan, Z. Q. Jin, Z. Zhu, D. Tománek, *Phys. Rev. B.*, **2014**, 90, 245403.
- [47] R. L. Murry, D. L. Strout, G. E. Scuseria, *Int. J. Mass Spectrom. Ion Processes* **1994**, 138, 113.
- [48] Y. F. Chang, J. P. Zhang, H. Sun, B. Hong, Z. An, R. S. Wang, *Int. J. Quantum Chem.* **2005**, 105, 142.
- [49] K. Raghavachari, *Chem. Phys. Lett.* **1992**, 190, 397.
- [50] P. W. Fowler, T. Heine, D. Mitchell, G. Orlandi, R. Schmidt, G. Seifert, F. Zerbetto, *J. Chem. Soc. Faraday Trans.* **1996**, 92, 12.
- [51] V. Schlegel, *Theorie der Homogen Zusammengesetzten Raumgebilde*; Dresden, Druck von E. Blochmann und Sohn, 1883.
- [52] W. Hakon, *J. Geol.* **1935**, 43, 250.
- [53] W. Helfrich, *Zeitschrift für Naturforschung Teil C* **1973**, 28, 693.
- [54] D. Holec, M. A. Hartmann, F. D. Fischer, F. G. Rammerstorfer, P. H. Mayrhofer, O. Paris, *Phys. Rev. B* **2010**, 81, 235403.
- [55] H. Hadwiger, *Vorlesungen Über Inhalt, Oberfläche und Isoperimetrie*; Springer, Berlin/Goettingen/Heidelberg, **1975**.
- [56] G. W. Wei, *Bull. Math. Biol.* **2010**, 72, 1562.
- [57] G.-W. Wei, Q. Zheng, Z. Chen, K. Xia, *SIAM Rev.* **2012**, 54, 699.
- [58] G.-W. Wei, *J. Theor. Comput. Chem.* **2013**, 12, 1341006.
- [59] Z. Chen, N. A. Baker, G. W. Wei, *J. Comput. Phys.* **2010**, 229, 8231.

Received: 3 September 2014

Revised: 25 October 2014

Accepted: 23 November 2014

Published online on 19 December 2014