

Topological Methods for Exploring Low-density States in Biomolecular Folding Pathways

Yuan Yao* Jian Sun[†] Xuhui Huang[‡] Gregory R. Bowman[§] Gurjeet Singh[¶]
Michael Lesnick^{||} Leonidas J. Guibas^{**} Vijay S. Pande^{††} Gunnar Carlsson^{‡‡}

August 25, 2008

*Department of Mathematics, Stanford University, Stanford, CA 94305. Email: yuany@stanford.edu.

[†]Department of Computer Science, Stanford University, Stanford, CA 94305. Email: sunjian@stanford.edu.

[‡]Department of Bioengineering, Stanford University, Stanford, CA 94305. Email: huangx@stanford.edu.

[§]Biophysics Program, Stanford University, Stanford, CA 94305. Email: gbowman@stanford.edu.

[¶]Department of Mathematics, Stanford University, Stanford, CA 94305. Email: gurjeet@stanford.edu.

^{||}Institute of Computational and Mathematical Engineering. Email: mlesnick@stanford.edu.

^{**}Department of Computer Science, Stanford University, Stanford, CA 94305. Email: guibas@cs.stanford.edu.

^{††}Department of Chemistry, Stanford University, Stanford, CA 94305. Email: pande@stanford.edu.

^{‡‡}Department of Mathematics, Stanford University, Stanford, CA 94305. Email: gunnar@math.stanford.edu.

Abstract

Characterization of transient intermediate or transition states is crucial for the description of biomolecular folding pathways, which is however difficult in both experiments and computer simulations. Such transient states are typically of low population in simulation samples. Even for simple systems such as RNA hairpins, recently there are mounting debates over the existence of multiple intermediate states. In this paper, we develop a computational approach to explore the relatively low populated transition or intermediate states in biomolecular folding pathways, based on a topological data analysis tool, Mapper, with simulation data from large-scale distributed computing. The method is inspired by the classical Morse theory in mathematics which characterizes the topology of high dimensional shapes via some functional level sets. In this paper we exploit a conditional density filter which enables us to focus on the structures on pathways, followed by clustering analysis on its level sets, which helps separate low populated intermediates from high populated uninteresting structures. A successful application of this method is given on a motivating example, a RNA hairpin with GCAA tetraloop, where we are able to provide structural evidence from computer simulations on the multiple intermediate states and exhibit different pictures about unfolding and refolding pathways. The method is effective in dealing with high degree of heterogeneity in distribution, capturing structural features in multiple pathways, and being less sensitive to the distance metric than nonlinear dimensionality reduction or geometric embedding methods. The methodology described in this paper admits various implementations or extensions to incorporate more information and adapt to different settings, which thus provides a systematic tool to explore the low density intermediate states in complex biomolecular folding systems.

Key words: biomolecular folding, RNA hairpin, SREMD, Morse theory, Mapper

Author Summary

In biomolecular folding problems, identification of intermediate or transition states connecting folded or unfolded states is crucial for understanding the folding pathways, but is difficult in both experiments due to the transient nature and computer simulations due to the low sampling population. We are motivated by recent debates over whether RNA hairpins fold in a two-state or multi-state manner with intermediates on pathways. To address this issue, we develop a computational approach to explore the relatively low populated transition or intermediate states in biomolecular folding pathways, based on a topological data analysis tool, Mapper, with simulation data from large-scale distributed computing. The method in this paper exploits a conditional density filter which effectively focuses on pathways, followed by clustering analysis on filter level sets which separates low populated intermediate states

from others of high population. Such a scheme, when applied to a RNA hairpin system, successfully provides important structural evidence on the multiple intermediate states on folding pathways which may guide further experimental investigations. The methodology adds to the communities of simulation and informatics an exploratory tool to data mine large data sets and propose experimentally testable properties, such as folding pathways for biological macromolecules.

1 Introduction

The folding of biomolecules is a classic biophysical problem. Proteins and nucleic acids are synthesized as linear polymer chains. They must then spontaneously and rapidly fold into their three-dimensional native states. The folding process is determined by the underlying free energy landscape. These landscapes are rugged, having many local minima corresponding to intermediates and misfolded states. Characterizing these states is critical for a full understanding of biomolecular folding. Experimental studies may point to the existence of such states but are usually unable to provide high resolution structural information due to the transience and/or heterogeneity of such states. Computer simulations have proved useful for sampling this complex high-dimensional space while yielding structures at full-atom resolution. However, these simulations tend to generate millions of configurations. The volume and high-dimensional nature of the output make it extremely difficult to discern the structure of the data.

One common approach to dealing with computer simulation results is to apply K-means clustering to the entire data set. However, K-means clustering suffers from a number of important limitations. First, it is limited by the need to specify the number of states from the beginning. Second, it tends to create spherical states. The relevant states of the free energy landscape, on the other hand, may be non-convex. In this case, K-means clustering will tend to lump unrelated configurations together or split related configurations into separate states. This limitation may be overcome by splitting the configurations into many small states and grouping them together using various metrics that allow non-convex states, such as in [1]. There is another widely used clustering method, single-linkage, which may overcome these issues in K-means. Unfortunately, identifying sparsely populated intermediate states is still difficult. Simulation data tends to be heavily dominated by the most stable states, such as the folded and unfolded states, and single-linkage clustering of the entire data set tends to pick up densest states only and hardly distinguish the intermediates from noise.

Recently, geometric embedding techniques, such as nonlinear dimensionality reduction [2-7], have been explored as a means to overcome the dimensionality hurdle in complex biomolecular systems. For example, ISOMAP [2] has been applied to protein folding [8] and Laplacian eigenmap [4] has been applied to the dynamics of biological networks [9]. This class of techniques maps the data in high dimensional spaces to a low dimensional space by preserving some local/global metric relationship among neighboring data points. In this way, one can easily visualize data and possibly gain important insights. For instance, the new embedding coordinates may be biologically relevant reaction coordinates [8]. However,

the performance of these geometric embedding techniques will suffer from the high degree of heterogeneity in distribution and be sensitive to the choice of the distance metric.

One efficient strategy to address these issues is to stratify the data into density level sets and study its topological features such as clustering which are less sensitive to the metric than geometric methods. High density levels will contain the dominant states, such as the folded and unfolded states, while less populated states, such as intermediates, will occupy the low density levels. Clustering on level sets of similar density, will be less affected by the distributional heterogeneity and thus effectively disclose structural information about intermediates. This idea of stratification is reminiscent of Morse theory, which provides a general machinery for studying the topology of high dimensional manifolds by looking at level sets of some nicely-behaved function [10]. Inspired by Morse theory, Singh *et al.* recently introduced Mapper [11], a topological data analysis tool for high dimensional data sets.

Mapper is a way to visualize and cluster high dimensional data. In its simple form, a filter function is used to decompose the data into overlapping level sets and clustering is then carried out in each of them. A graph is then generated by connecting clusters in neighboring level sets with an edge if they have non-empty overlapping. If an energy function is taken as the filter, such a graph in the limiting case will summarize the same kind of topological information as a disconnectivity graph of the energy landscape [12]. But the flexibility of the filter design in Mapper makes it applicable in non-equilibrium data as well. In its extended form, Mapper can return a simplicial complex with high dimensional topological information about the data. The method is computationally efficient and amenable to parallelization.

In this work we demonstrate the applicability of Mapper in its simple form to the biomolecular folding problem. We begin with a discussion of Mapper itself from a perspective of Morse Theory and then present the details of a filtering function that is well-suited for biomolecular folding problems, the conditional density filter. This filter weights configurations close to a state of interest more heavily, thus facilitating the identification of any intermediate state leading up to it. We then describe the use of single-linkage clustering within level sets to allow the identification of an unspecified number of non-convex states. Finally, we discuss the application of Mapper to the folding of a small RNA hairpin, which gives some structural evidence from computer simulations in support of the multi-state hypothesis [13]. The biological implications of the Mapper results are discussed by [14] elsewhere. We also briefly discuss the advantages of Mapper over nonlinear dimensionality reduction techniques. In the future we hope to explore the combination of those geometric embedding techniques with Mapper in order to take advantage of the strengths of both approaches.

2 Materials and Methods

2.1 Mapper: A Tool for Topological Data Analysis

One way to reduce the computational complexity in the study of massive data sets is to decompose the data by classifying the data into groups and doing analysis on each of the

group individually instead of performing analysis on the whole. This strategy is amenable for parallel computation, which is particularly important for studies of biomolecular folding, where a great amount of configurations are normally generated.

Here we pursue this idea in the particular case where the decomposition is induced by the choice of some filter function on the data set, $h : \mathcal{X} \rightarrow \Omega$. In this paper, we will only consider filters that take values in the real line, though the Mapper methodology is equally applicable for filter functions taking values in higher dimensional space, or even spheres, tori, or any other topological space. With this choice, we introduce Mapper from a perspective of Morse Theory, which differs from the original paper [11] but discloses a deeper inspiration.

Morse Theory [10] tells us that when $h : \mathcal{X} \rightarrow \mathbb{R}$ is some nicely-behaved function¹, topological information of \mathcal{X} can be inferred from the level sets $h^{-1}(\omega)$. Morse theory is an extremely powerful tool to analyze the topology of high dimensional manifolds, which lies in the heart of proving the celebrated Poincare Conjecture of dimension no less than five [15].

The simplest example in this spirit may be *Reeb graph* [16], by contracting to points the connected components within level sets $h^{-1}(\omega)$, illustrated as Figure 1 (a). This simple scheme turns out to be useful in various fields under different names, *e.g.* contour trees in computational geometry [17] and cluster trees in statistics [18, 19, 20].

Mapper [11] extends this construction to incorporate the discrete setting where \mathcal{X} is a finite set of data points in high dimensional spaces or metric spaces. First, instead of working with the level set of a single value which is difficult to capture in discrete settings, Mapper considers the preimage of subinterval $h^{-1}([a, b])$. Second, it replaces by clustering the contraction of connected components in continuous settings. Specifically, the procedure of Mapper used in this paper is as follows.

1. **Level-set formation.** Cover the range of $h : \mathcal{X} \rightarrow \mathbb{R}$ by a set of subintervals which overlap in neighbors, *i.e.* $U_i = [a_i, b_i]$ with $U_i \cap U_{i+1} \neq \emptyset$ and $U_i \cap U_j \cap U_k = \emptyset$, and stratify \mathcal{X} into level sets by taking inverse images $h^{-1}([a_i, b_i])$;
2. **Clustering.** On each level set $h^{-1}([a_i, b_i])$, construct the connected components or point clusters;
3. **Graph representation.** Represent each component or cluster by a node. Add an edge between a node pair whenever they have nonempty intersection.

Mapper thus returns an undirected graph representing the connectivity information between data clusters across level sets $h^{-1}([a_i, b_i])$. See the example in Figure 1 (b). Note that those degree-one nodes lie in the intervals containing local minima/maxima and the branching (degree-three) nodes lie in the intervals with saddle points, a sort of critical points.

More generally, if the filter value range Ω takes some higher dimensional space or other topological spaces, Mapper may return a simplicial complex which is however not pursued

¹They are called Morse functions, *i.e.* those smooth functions with only nondegenerate critical points; in other words, the Hessian at each critical point where the gradient vanishes has full rank. Morse functions are generic in the sense that they are dense in the space of smooth functions, as well of continuous functions. Hence every continuous function can be approximated arbitrarily well by Morse functions.

in this paper. This construction can easily yield a multiresolution structure by choosing subintervals of different granularities, which helps handle noise.

The key choice in Mapper will be the filter map $h : \mathcal{X} \rightarrow \Omega$. In fact, the name, *Mapper*, was coined to emphasize the importance of choosing such a map. There is no universal scheme for this choice, which may vary from application to application. In [11] some examples are presented with the choice of density function and a certain eccentricity function measuring data depth as filters. In the following, we will discuss it in detail in the setting of biomolecular folding problems, with a particular example in RNA hairpin folding.

2.2 Mapper Design in Biomolecular Folding

Simulation data in biomolecular systems produces massive data in high dimensional space, and exhibits heterogeneity in distributions. The general procedure of Mapper above is adapted toward such challenges. The first crucial design is to construct filters based on conditional density functions estimated from data, which effectively enables us to focus on important local regions in configuration spaces and separate less populated pathways from the overwhelmed uninterested states. In clustering we choose the single-linkage method to capture possibly non-convex clusters. Below we give a detailed description on these particular implementations.

2.2.1 Conditional Density Filters for Mapper

Our key construction of filters here is based on conditional density functions estimated from data, conditioning on the states of interests. For example, in the study of folding process we extract configurations from folding events and focus on the region close to folded states, while in unfolding process we draw samples from unfolding events and pay more attention to the zone around extended states. Simulation trajectories of those processes are often dominated by stochastic fluctuations around the initial states. It is near the target states that one may observe interesting structural information about pathways. The conditional density filters are chosen to reflect the localized free energy landscape around the states on pathways without being disturbed by off-pathway structures which are quite noisy and of high population in samples.

Although the simulation data of biomolecular systems often lie in a high dimensional configuration space, the degree of freedom are much less due to the constraints and cooperation among atoms in folding process. It is often expected that the pathway samples are concentrated around some low dimensional manifolds which can be described by a relatively small number of intrinsic reaction coordinates [8]. The existence of multiple pathways as in the example of this paper may lead to holes in such manifolds with nontrivial topology. Note that in the continuous case, the Reeb graph of a (unconditional) density function defined on the Euclidean space \mathbb{R}^n turns out to be trivially a tree. However conditional density functions adopted here may restrict on interesting regions where the loops in the Reeb graph might shed light on the hole structures. Reconstructing the low-dimensional topology of

densely sampled regions, thus may disclose the nature of multiple pathways. In theory, it is possible to efficiently recover the topology from samples of such low dimensional manifolds [21]. In this paper, through conditional density filters we approach such manifolds via data level sets and extract some low-dimensional topological features which provide structural evidence on the existence of multiple pathways.

Here we describe a general approach to construct conditional density filters, which will be specialized in the next section with the application to RNA hairpin folding.

- Draw random samples $S \subseteq \mathcal{X}$ from the folding events. Choose importance weights on S , $w(x) \geq 0$, with higher values on interested states.
- Define the filter function by

$$h(x) = -\log \frac{\sum_{y \in S} w(y) K(x, y)}{\sum_{y \in S} w(y)}, \quad (1)$$

where the kernel function is defined by

$$K(x, y) = e^{-\frac{d^\beta(x, y)}{\alpha}}, \quad (2)$$

where $d(x, y)$ is some distance function between configuration x and y , $\alpha > 0$ is the band-width, and $\beta > 0$ is the exponent. For example the Euclidean distance with $\beta = 2$ is used in case of Gaussian kernels, or the Hamming distance between structural contact maps with $\beta = 1$ is used later in this paper.

- Resample from S according to the new distribution,

$$p(x) = \frac{w(x)}{\sum_x w(x)}.$$

To avoid the normalization in large data sets, we can use the rejection method or extended [22], *e.g.* a sequential Bernoulli experiments where a new configuration x is accepted with probability $q(x) = \frac{w(x)}{\max\{w(x)\}}$.

These configurations, together with the filter function (1), will be the inputs of Mapper procedure shown in the last section.

Filter (1) assumes a density function in Boltzman form $f(x) = \frac{1}{Z} e^{-h(x)}$, with partition function $Z = \sum_x e^{-h(x)}$. Thus up to a constant filter (1) approximates the free energy near the folded state. Since only order information of $h(x)$ will be used below, it leads to the same result choosing any monotone transform on h , *e.g.* $\sum_{y \in S} w(y) K(x, y)$. Our construction is equivalent to a kernel density estimator which can be replaced by other methods [23].

2.2.2 Level-set Formation in Mapper

To increase the robustness of Mapper allowing more errors in density estimation, we only use the order information of filter (1) to construct level sets.

- **Level-set formation.** Order the samples according to values of $h(x)$, and classify the samples into m consecutive overlapping groups of equal or similar size, whose filter value ranges $[a_i, b_i]$ cover the range of h .

Up to an arbitrary small perturbation, a real valued function $h : \mathcal{X} \rightarrow \mathbb{R}$ induces a linear order on samples. Therefore any monotone transform on $h(x)$, such as $c_1 \exp c_2 h(x)$, leads to the same level sets.

2.2.3 Clustering in Mapper

The graphical representation of Mapper depends on the choice of clustering methods. Mapper itself does not place any prerequisite on the clustering algorithm. In the study of biomolecular folding such as RNA hairpins, our purpose is to identify those connected components in free energy or density level sets, which might be of non-convex shapes and whose numbers are unknown to us beforehand. Single-linkage clustering is the simplest choice to meet those two features.

- On each level set, construct a weighted graph, with nodes for configurations and edge weights as pairwise distances.
- Find a Minimal Spanning Tree (MST) of such a graph.
- Find a threshold value for edges. We construct a histogram of MST edge weights with k bins. Once some empty bins are found from top bins containing p longest edges, we set the threshold to be the center of the first empty bin. Otherwise, set the threshold the maximal edge (diameter).
- Truncate the graph by breaking those edges greater than the threshold, dividing the graph into connected components.
- Prune those components of size no more than q .

Single-linkage will separate those clusters where within each cluster two points can be joined by a path consisting of short edges, but relatively longer edges are required to merge the clusters. When we draw random samples from compact connected components in an Euclidean space, the distances between configurations within the same components will drop down to zero as the sample size grows. Hence the distances across components will be kept in the longest edges and can be separated from a large amount of short edges. Thresholding above tries to capture such a gap. Truncation may create several components/clusters of different sizes, where pruning helps reduce the noise and identify those dominant components.

In the continuous setting, single-linkage clustering will consistently locate those connected components when the samples are dense enough [18]. Such a feature makes it a desirable choice for Mapper [11], as well as density cluster trees [19, 20]. However, in the latter part of this paper, we will meet a discrete configuration space, *i.e.* the space of contact maps as undirected graphs. Thus we need to explain in what sense we extend the “connected components” in such a discrete setting.

Equipped with a metric, *e.g.* Hamming distance, the discrete configuration set can be viewed as a weighted complete graph, where each node represents a structure and the weight of an edge is the distance between its endpoints. Single-linkage clustering firstly builds up a Minimal Spanning Tree (MST) of this graph and then truncate the MST by keeping the edges with the length less than a given threshold, which breaks the MST into several connected components or clusters. In this way, single-linkage compute the components where two nodes within a component are joined by a path consisting of the short edges, but relatively longer edges are required to merge different components.

One may also consider other clustering schemes, such as k -means, which is widely used in clustering the configurations in the biomolecular folding simulations. In contrast to single-linkage, k -means attempts to find the clusters such that within cluster *any* two nodes are connected by a short edge, rather than by a path made up of short edges. Therefore, roughly speaking, k -means attempts to find *spherical-shape* clusters while single-linkage can discover *snake-shape* clusters. Both may provide useful but different kinds of information in biomolecular folding problems. However k -means needs one to specify the number of clusters *a priori* while the single-linkage does not. This is a short-coming for k -means since we don’t have such information in advance. So in this paper, single-linkage is chosen as the basic scheme and k -means is only used in comparative studies, when we already know the number of clusters from single-linkage. Other choice of clustering methods includes average-linkage, complete-linkage, spectral clustering [24, 25], *etc.*, which are however not pursued in this paper.

3 Results and Discussions

Recently, [14] performed Serial Replica Exchange Molecular Dynamics (SREMD)[26, 27] simulations of the GCAA tetraloop (5’-GGGCGCAAGCCU-3’) on the Folding@home distributed computing platform. The hairpin motif consists of a primarily Watson-Crick base-paired stem capped with a loop of unpaired or non-Watson-Crick base-paired nucleotides, as shown in Figure 2 (a). Despite their simple structures, there is some debate over whether or not there are intermediate states in the folding of hairpins, *e.g.* see [13].

With the technique developed in this paper, we are able to disclose the structures of multiple intermediate states on the folding pathways, which in the first time provides structural evidence from computer simulations about RNA hairpin folding pathways. The biological implications of this discovery are discussed in detail by Bowman *et. al.* [14]. Here, we only focus on details of data analysis.

The RNA molecule examined here has 389 atoms. Including the solvent there are about

$N = 12,000$ atoms in the system, yielding $3N = 36,000$ parameters. To reduce the dimensionality of this large space we chose to represent each configuration with a contact map. Contact maps can faithfully describe the base-pair interactions in the stem, which provides important structural information of RNA hairpin folding. A contact map is a bit string specifying pairs of contacting residues that are not immediately adjacent in the sequence. Thus, even this coarse-grained space is \mathbb{R}^{55} . Following Bowman *et. al.* [14], we define the *native state* as any conformation with all four stem base-pair contacts formed. Each of these base-pair contacts is referred to as a native contact. For example, Figure 2 (a) shows a native state whose contact map model is illustrated in (b). An *unfolding* event is defined as the set of conformations between the first point with no contacts between any two residues on opposite sides of the stem and the first preceding point with four native contacts. A *refolding* event is defined as the set of conformations between the first point with no contacts between any two residues on opposite sides of the stem and the first subsequent point where the number of native contacts is four.

3.1 Structural Analysis by Mapper

Mapper is an ideal tool for such a problem due to the enormous size of the simulation dataset, the high probability of non-convex states, and the need to identify folding intermediates with low populations relative to the folded and unfolded states. Application of Mapper to this data set revealed a number of intermediate states.

The data generated from SREMD simulations is normally dominated by the folded and unfolded structures. For example, a typical refolding trajectory starts from an unfolded state, undergoing a significant period of stochastic fluctuation around that, then proceeds gradually to the folded state. It is in the neighborhood of folded states that interesting structural information about folding pathways are exhibited. Therefore, in the construction of the conditional density filters, we treat folding and unfolding separately. In the study of folding pathways, we take configurations from refolding events, and then weight heavily a neighborhood around the native states. However in the study of unfolding pathways, we sample from unfolding events, and focus on a neighborhood of the unfolded states.

The following parameters are used to produce the results in Figure 3. We use the Hamming distance $d_H(x, y)$ between a pair of contact maps in the conditional density function (Equation (1)) and choose $\alpha = \beta = 1$ in kernel (2). For simplicity, the importance weights are set to one within the neighborhood of the state of interest and zero otherwise. In refolding events, we choose a neighborhood within 7-bit Hamming distance from the native state in Figure 2. In unfolding events, a neighborhood of the extended state is chosen as the set of configurations with no more than 6 non-adjacent contacts formed. In the level-set formation, the filter is divided into 8 levels of equal size with 25% overlap. In the clustering, a histogram with 5 bins is used, with thresholding from top bins consisting largest $p = 20\%$ edges and the cluster pruning size $q = 2\%$ of the level sample size. More details on parameter tuning will be provided in supplementary material.

The graphical output of Mapper with such parameters shows distinct pictures about

folding and unfolding pathways. Unfolding has a single dominant pathway characterized by unzipping from the end base-pair (Figure 3 (a)), while folding process has two dominant pathways, passing through either the formation of the closing base-pair or the end base-pair (Figure 3 (b)). Such an observation reveals a number of intermediate states in the folding process, which supports the multi-state hypothesis. It is interesting to notice in Figure 3 that conditional density filters seem good indicators of reaction coordinates, suggesting that the folding/unfolding processes start from the densest zone and become sparser as the reactions proceed.

3.2 Kinetic Verifications

Are the two pathways in refolding (Figure 3 (b)) are truly separate pathways or just the artifact of noise? This question can be answered from the kinetic information of simulation trajectories by computing the transition probability. Note that our purpose here is not to create a Markov model [1] for metastable states, but investigate how the two intermediate states in refolding pathways are kinetically connected. Therefore the shortest lag time, 2 ps, is chosen which provides the finest resolution in simulation trajectories.

To simplify the result, we merge the four nodes with extended structures as a single unfolded state, U, and collapse the three blue nodes with folded structures as folded state, F, leaving alone the two intermediates, I1 and I2. This does not change the topology of Mapper graph, but highlights the dynamics associated with intermediate states. Configurations in simulations are mapped to such four-node states by nearest neighbor method. One-step (2 ps) transition probability are then computed among the four states.

The result is shown in Figure 4. It can be seen that the two intermediates, I1 and I2, are kinetically well separated on folding pathways. Once the simulation climb up the energy barrier I1 and I2, the majority will either proceed to F or withdraw to U, while an ignorable minority will cross the intermediates from I1 to I2.

3.3 Importance of Conditional Density Filters

Conditional density filters play a crucial role here, without which clustering methods like K-means or single-linkage tend to split the sparse intermediates and lump them with densest clusters.

To see this, we make a comparison between Mapper clusters found in Figure 3 (b) and K -means clustering on the same data set. Since the number of K -means clusters is not unknown *a priori*, we performed a series of experiments with k varying from 1 to 80, each of which has 20 repeated experiments. Our first purpose is to locate the value of k around which the Mapper cluster with end base-pair formed becomes identifiable. Hence for each K -means experiment, we count the number of the end base-pair clusters, defined as the clusters containing more than 75% configurations with native contact 4 (Figure 2 (b)) formed, and less than 25% for any other native contact. Figure 5 plots a rough distribution of the numbers of end base-pair clusters against the growth of k . It can be seen that around $k = 25$ this

intermediate state becomes identifiable, in the sense that with more than 1/2 probability such clusters are found indicated by nonzero medians. Notice that as k grows, the variation range (10% \sim 90%) of such cluster numbers expands, showing a trend of increasing instability. Particularly around $k = 55$, such a state begins to split into several K -means clusters.

We can further see how K -means clusters might split the intermediate states and lump them toward densest clusters. Figure 6 illustrates this when $k = 30$ for K -means clustering, on the same data set for the construction of Mapper clusters on refolding pathways.

3.4 Comparative Studies on Single-linkage vs. k -means

Single-linkage clustering is motivated by its ability to identify possibly non-convex clusters of unknown number. It is also interesting to explore other clustering methods such as K -means which tries to group data in *spherical* clusters and is widely used in the studies of biomolecular folding simulations. Given the cluster number returned by single-linkage, comparisons with K -means of similar number of clusters on the same level sets might disclose how far the intermediate states deviate from spherical shapes. For this purpose, we perform K -means clustering on the same data set for refolding pathways in Figure 3 (b). We use the same number of clusters returned by single-linkage and especially on level five we set $k = 2$. It turns out that K -means finds two clusters on level five with similar structural features to single-linkage, *i.e.* one with closing base-pair formed and the other with the end base-pair. However K -means has different partition: 48% vs. 52%, in contrast to 23% vs. 44% in single-linkage. Clearly to form spherical clusters, K -means clusters mix more configurations from different single-linkage clusters, which can be shown by the percentage dropping of dominant end base-pair from 96% to 65% in the smaller cluster. However the structural similarity in both methods suggests that single-linkage clusters are not very far from spherical shapes.

3.5 On Nonlinear Dimensionality Reduction

Although a biomolecular system is typically described by a high dimensional configuration space, it is expected that those configurations often visited in a folding process may concentrate around some low-dimensional manifolds which might be described by a much smaller number of reaction coordinates. Recently Das *et. al.* [8] shows that ISOMAP can be applied to recover such reaction coordinates in simple folding processes with a single pathway. Isomap tries to preserve both the local and global geodesic distance between configurations defined as shortest path distance on a neighborhood graph. However, ISOMAP might not work in complex problems where multiple pathways exist. ISOMAP requires that the data manifold are globally isometric to a convex domain of low dimensional space [2, 5]. The existence of more than two pathways connecting two metastates, may lead to holes in sampled regions which fails the convex domain assumption. Moreover, ISOMAP is too sensitive to the metric in choice. In this paper we use a coarse metric as Hamming distance for contact maps, where the geodesic distance between configurations does not reflect the distance

in folding process. Moreover, The high heterogeneity in distribution is also a hurdle for ISOMAP technique to identify useful intermediates.

The last two issues also challenge other techniques for nonlinear dimensionality reduction, such as LLE [3], Laplacian Eigenmap [4], Hessian Eigenmap [5], and Diffusion map [6], *etc.* These geometric embedding techniques maps the data in high dimensional spaces to a low dimensional space by preserving some local metric relations among neighbors of data points, *e.g.* see [7]. They are thus sensitive to the metric in choice and heterogeneous distribution might distort local metrics. In applications to complex biomolecular systems, successful examples are only found in simple settings such as with a single protein folding pathway [8] or quasi-steady state in dynamics of signal transduction networks [9].

However, as a topological tool Mapper with density filters is shown efficient in dealing with heterogeneous distributions and less sensitive to the metric in choice. In this paper even with such a coarse metric as Hamming distance, it efficiently discloses structural information in pathways which are difficult to other geometric embedding techniques. Thus, one of our ongoing directions is to combine the topological tool Mapper with those geometric embedding techniques, such as applying nonlinear dimensionality reduction separately on components or clusters discovered by Mapper.

4 Conclusions

In this paper we develop Mapper, a topological data analysis tool, in the analysis of simulation data for biomolecular folding pathways. As an application, in the first time we are able to obtain structural evidence from computer simulations in support that RNA hairpin folding has two dominant pathways with multiple intermediate states. It is thus a promising direction to explore with Mapper such structural information in biomolecular folding problems.

We have shown that with proper designs of conditional density filters and clustering schemes, Mapper can address the heterogeneity issue in distribution, deal with multiple pathway data with nontrivial topology, and be less sensitive to the metric in choice. These features can be used to enhance traditional nonlinear dimensionality reduction methods, such as ISOMAP, Laplacian Eigenmap, Diffusion maps, *etc.* One of our ongoing direction is to explore the combinations of the topological tool Mapper with those geometric tools for better characterizations of biomolecular systems.

Only a limited use of data has been pursued for Mapper in this paper, where we ignore the kinetic information in simulation trajectories. Since we did not take into account the kinetic information, the intermediate states we identify are thermodynamically relevant states, but not necessarily to be kinetically relevant states. Our future direction is to incorporate such kinetic information for developments of novel dynamical models.

Acknowledgments

YY would like to thank Wing-Hung Wong, Nancy Zhang, and Qing Zhou for helpful discussions, XH thanks Michael Levitt for his support. YY, JS, GS, LG and GC are funded by DARPA grant HR0011-05-1-0007; XH by NIH Roadmap for Medical Research Grant U54 GM072970; GB by the NSF Graduate Research Fellowship Program; ML by a NDSEG fellowship. YY and GC are also supported by NSF grant DMS 0354543, and LG by NSF grant FRG-0354543 and NIH grant GM-072970. This work was also supported by NIH P01 GM066275. Computing resources were provided by the Folding@home users and NSF award CNS-0619926.

References

- [1] Chodera, J. D., Singhal, N., Pande V. S., Dill, K. A., and Swope W. C. (2007) Automatic discovery of metastable states for the construction of Markov models of macromolecular conformational dynamics. *J. Chem. Phys.*, **126** (15), 155101.
- [2] Tenenbaum, J., de Silva V. and Langford, J. (2000) A global geometric framework for nonlinear dimensionality reduction. *Science* **290**, 2319-2323.
- [3] Roweis, S. T. and Saul, L. K. (2000) Nonlinear dimensionality reduction by locally linear embedding. *Science* **290**, 2323-2326.
- [4] Belkin, M. and Niyogi, P. (2003) Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Comput.* **15**, 1373-1396.
- [5] Donoho, D. and Grimes C. (2003) Hessian eigenmaps: locally linear embedding techniques for high-dimensional data. *Proc. Natl. Acad. Sci. USA* **100**, 5591-5596.
- [6] Coifman, R. R., Lafon, S., Lee, A. B., Maggioni, M., Nadler, B., Warner, F. and Zucker, S. W. (2005) Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. *Proc. Natl. Acad. Sci. USA* **102**, 7426-7431.
- [7] Jones, P. W., Maggioni, M. and Schul, R. (2008) Manifold parametrizations by eigenfunctions of the Laplacian and heat kernels. *Proc. Natl. Acad. Sci. USA* **105**, 1803-1808.
- [8] Das, P., Moll, M., Stamati, H., Kaviraki, L. E. and Clementi C. (2006) Low-dimensional, free-energy landscapes of protein-folding reactions by nonlinear dimensionality reduction. *Proc. Natl. Acad. Sci. USA* **102**, 10141-10146.
- [9] Barbano, P., Spivak, M., Flajolet, M., Nairn, A. C., Greengard, P. and Greengard, L. (2007) A mathematical tool for exploring the dynamics of biological networks. *Proc. Natl. Acad. Sci. USA* **104**, 19168-19174.
- [10] Milnor, J. (1963) *Morse Theory* (Princeton University Press).

- [11] Singh, G., Mémoli, F. and Carlsson, G. (2007) Topological Methods for the Analysis of High Dimensional Data Sets and 3D Object Recognition. In *Eurographics Symposium on Point-Based Graphics*.
- [12] Becker, O. M. and Karplus, M. (1997) The topology of multidimensional potential energy surfaces: theory and application to peptide structure and kinetics. *J. Chem. Phys.*, **106**, 1495-1517.
- [13] Ma, H., Wan, C., Wu, A. and Zewail, A. H. (2007) DNA folding and melting observed in real time redefine the energy landscape. *Proc. Natl. Acad. Sci. USA* **104**, 712-716.
- [14] Bowman, G., Huang, X., Yao, Y., Sun, J., Carlsson, G., Guibas L. and Pande, V. (2008) Structural Insight into RNA Hairpin Folding Intermediates. *J. Am. Chem. Soc.*, preprint.
- [15] Smale, S. (1961) Generalized Poincaré’s conjecture in dimensions greater than four. *Ann. of Math.* **74**, 391-406.
- [16] Reeb, G. (1946) Sur les points singuliers d’une forme de pfaff complètement intégrable ou d’une fonction numérique. *Comptes Rendus Acad. Sci. Paris* **222**, 847-849.
- [17] van Kreveld, M., van Oostrum, R., Bajaj, C., Pascucci, V. and Schikore D. (1997) Contour trees and small seed sets for isosurface traversal. *Proc. 13th Annu. ACM Sympos. Comput. Geom.*, 212-220.
- [18] Hartigan, J. A. (1981) Consistency of single linkage for high-density clusters. *J. Amer. Statist. Assoc.* **76**, 388-394.
- [19] Stuetzle, W. (2003) Estimating the cluster tree of a density by analyzing the minimal spanning tree of a sample. *J. Classification* **20**: 25-47.
- [20] Zhou, Q. and Wong, W.-H. (2008) Reconstructing the energy landscape of a distribution from Monte Carlo samples. *Ann. App. Stat.*, preprint.
- [21] Niyogi, P., Smale, S. and Weinberger, S. (2008) A topological view of unsupervised learning from noisy data. *preprint*.
- [22] Liu, J. (2004), *Monte Carlo Strategies in Scientific Computing* (Springer).
- [23] Silverman, B. W. (1986) *Density Estimation for Statistics and Data Analysis* (CHAPMAN & HALL/CRC).
- [24] Hastie, T., Tibshirani, R. and Friedman, J. (2001) *The Elements of Statistical Learning* (Springer-Verlag).
- [25] Ding, C. and Zha, H. (2007) *Spectral Clustering, Ordering and Ranking – Statistical Learning with Matrix Factorizations* (Springer, ISBN: 0-387-30448-7).

- [26] Huang, X., Bowman, G. R. and Pande, V. D. (2008) Convergence of folding free energy landscapes via application of enhanced sampling methods in a distributed computing environment. *J. Chem. Phys.* preprint.
- [27] Menger, M., Eckstein, F. and Porschke, D. (2000) Dynamics of the RNA hairpin GNRA tetraloop. *Biochemistry*. *Biochemistry* **39**, 4500-4507.

Appendix A: RNA Hairpin Folding Simulations

Our simulations used the AMBER 94 potential. 2,800 SREMD simulations with an aggregate simulation time of $54.6\mu s$ were performed, see [14] for details. Among 2,800 SREMD simulations we obtain 760 trajectories with a complete unfolding event and 550 trajectories with a complete refolding event. Note that an unfolding event defined above only contains one unfolded state as the end point, whose density is thus too low in samples. Therefore, among such trajectories, we randomly choose 149 extended unfolding events and 23 extended refolding events, which includes $m = 10$ more points after the end point of each event. In this way we obtain about 100,000 samples for either class of events.

Note that contact maps are used as a discrete representation of structures, whence different configuration samples might have the same contact map representation. Such repetitions should be kept for density estimation, but can be compressed into unique structures for clustering analysis. In fact, those samples contain 49,332 and 56,118 unique contact maps, for unfolding and refolding extended events, respectively. They are sufficient for the analysis by Mapper.

Appendix B: Parameter Choice in Mapper

Conditional Density Filter. The data generated from SREMD simulations is normally dominated by the folded and unfolded structures. For example, a typical refolding trajectory starts from an unfolded state, undergoing a significant period of stochastic fluctuation around that, then proceeds gradually to the folded state. It is in the neighborhood of folded states that interesting structural information about folding pathways are exhibited. Therefore, in the construction of the conditional density filters, we treat folding and unfolding separately. In the study of folding pathways, we take configurations from refolding events, and then weight heavily a neighborhood around the native states. However in the study of unfolding pathways, we sample from unfolding events, and focus on a neighborhood of the unfolded states.

To be specific, in the study of unfolding, we extract 4,330 configurations around the extended states with no more 6 non-adjacent contacts formed. On the other hand, in the study of refolding, we extract 2,952 configurations of no more that 7-bit Hamming distance away from the native state to avoid the highly populated extended states. This is equivalent to the choice of a weight function $w(x)$ which is a constant in a neighborhood of the

extended states (no more than 6 non-adjacent contacts) or the native state (no more than 7-bit Hamming distance) and zero otherwise.

Since the space of contact maps is discrete, we use Hamming distance and choose $\alpha = \beta = 1$ in kernel density estimation (2), which is equivalent to the Gaussian kernel with the standard Euclidean distance in \mathbb{R}^{55} .

We note that in a range of $1 \leq \alpha \leq 8$ Mapper returns qualitatively similar results. In fact smoothing the density filter without changing the order leads to the same result in Mapper. However decreasing α , even to 0.9, causes the disappearance of the smaller passway. A small choice of α creates a rugged density filter, which alters the results of Mapper. Our experiments show that $\alpha = 1$ is close to this bifurcation point.

Level Sets. We divide the range of the density filters into n overlapped intervals, where each interval contains the same number of samples. In other words, we order the samples according to the filter value, then divide the sample into overlapped bins of equal size. It is also possible to consider division by equal filter value intervals [11], but the former has at least two advantages. First the former method only takes into account the order information about the filters to stratify the data, whence any monotone transformation on $h(x)$, such as $-\log h(x)$, leads to the same result. This makes the result from Mapper relatively more robust to the error in density estimation. Second it is more convenient to control the computational cost where each level has similar running time due to the same sample size, which is suitable for parallel computations.

We have tested the choice of n among 4, 6, 8, 10, 12, 14, so that each level contains around several hundred configurations. Overlap percentage can be chosen from 15% to 75%. All of them give qualitatively similar results though smaller number of levels and larger overlap causes longer computation. The results presented in this paper are generated under the choice of 8 intervals with 25% overlap.

Single-Linkage Clustering.

To determine an appropriate threshold in single-linkage, we build up a histogram based on the edge weights in the MST using $k = 5$ bins. We only focus on those bins containing the largest $p = 20\%$ edges. The threshold value is chosen to be the center of the first empty bin among them, defined as less than $q = 2\%$ samples in the largest bin. The reason we do so is that the empty bin with short edges often appear due to under-sampling which does not tell us information about gaps among components. If there is no such short bin, take the entire level set as one cluster. The results of Mapper will be sensitive to the choice of such k . Generally speaking, increasing k will increase the number of clusters, and vice versa. Considering the fact that the diameter of the data is about 14, we normally choose k an integer between 5 and 10, which all gave qualitatively the same results as $k = 5$. We note that although the threshold found by histogram method is sensitive to the bin number k , the threshold leading to the two clusters on level five in Figure 3 (b) is always 2-bit Hamming distance, which is very robust in different choices of k .

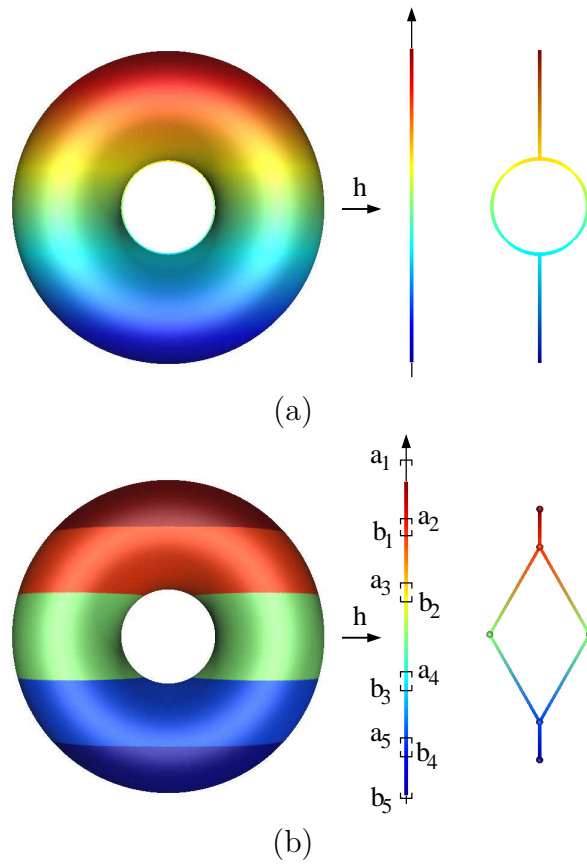
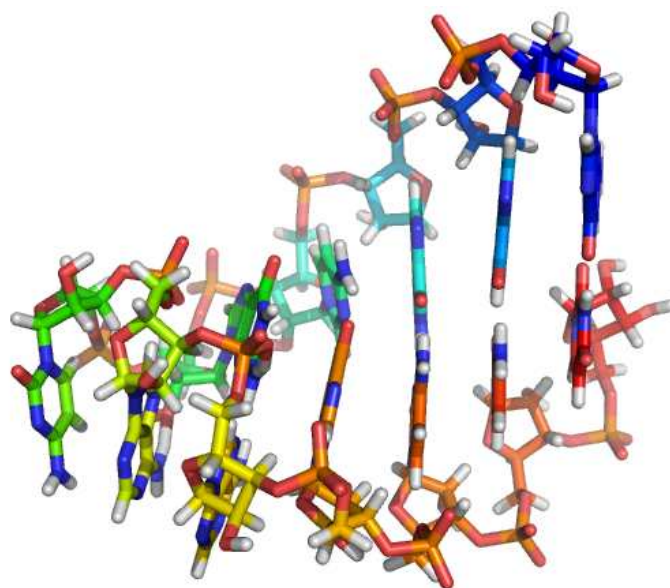
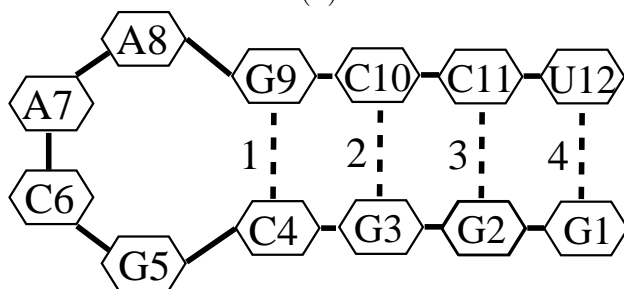


Figure 1: (a) Construction of Reeb graph; (b) Construction of Mapper. h maps each point on torus to its height.

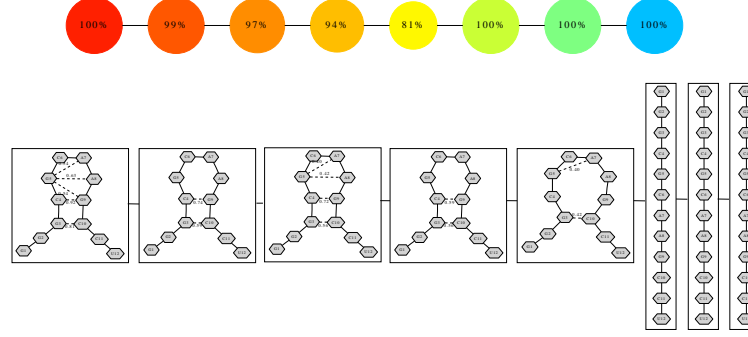


(a)

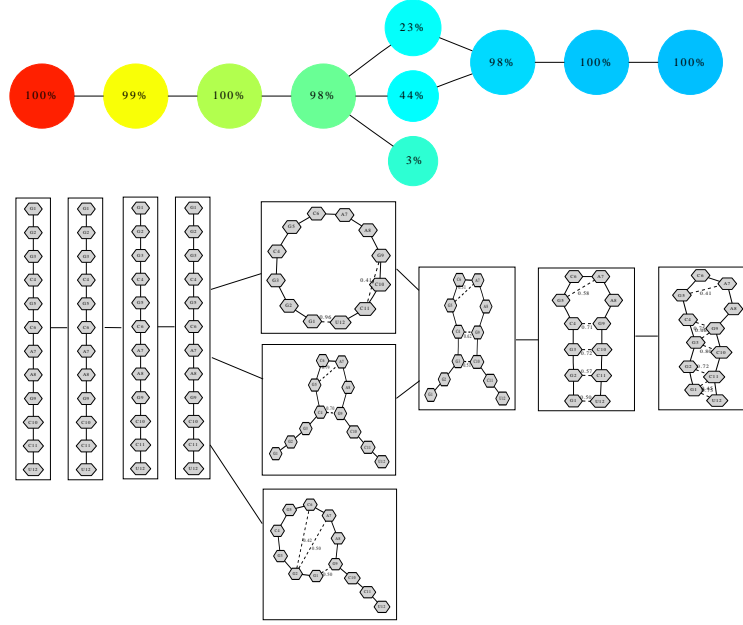


(b)

Figure 2: (a) NMR structure of the GCAA tetraloop. (b) Contact map for the native state. Bases are numbered from 1 to 12 and native basepair contacts (dotted lines) are numbered 1-4.



(a)



(b)

Figure 3: Graphical representation of pathways by Mapper. (a) Unfolding pathway. (b) Folding pathway. In both cases, the top row graphs are the outputs from Mapper, while the bottom row depicts the mean contact maps of the corresponding clusters. For clarity in mean contact maps we drop those mean contacts lower than 0.4. The node colors from red to blue indicate the density from high to low, and the labels (*e.g.* 100%) show the percentage of configurations of the same level included in the cluster corresponding to the node. We dropped all the clusters of size smaller than 3% of the level size. (a) shows that unfolding has a single dominant pathway characterized by unzipping from the end base-pair. (b) shows that folding process has two dominant pathways, passing through either the formation of the closing base-pair or the end base-pair. A noisy cluster consisting 3% of the level size was also shown in (b), which accounts for reptation, *i.e.* sliding of the two strands of the stem.

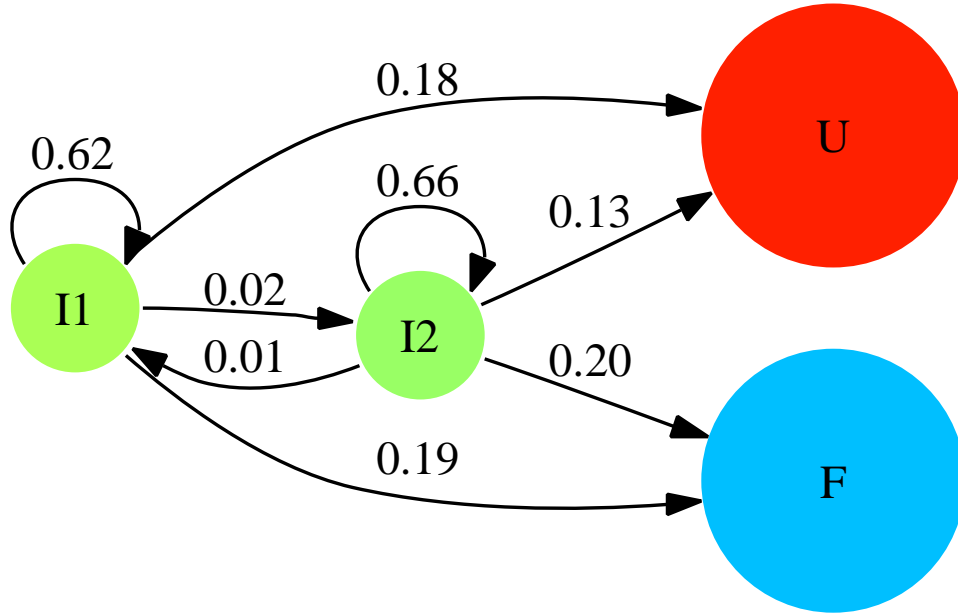


Figure 4: Transition probability from two intermediate states. Lag time is $2ps$. The left four nodes as extended structures (Figure 3 (b)) are merged into node U, and the right three nodes as folded structures are collected in node F. The two intermediate states on pathways are denoted by I1 and I2, respectively. The transition probability from I1 and I2 to other states are noted as numbers on the arrows. One can see that I1 and I2 are kinetically separated.

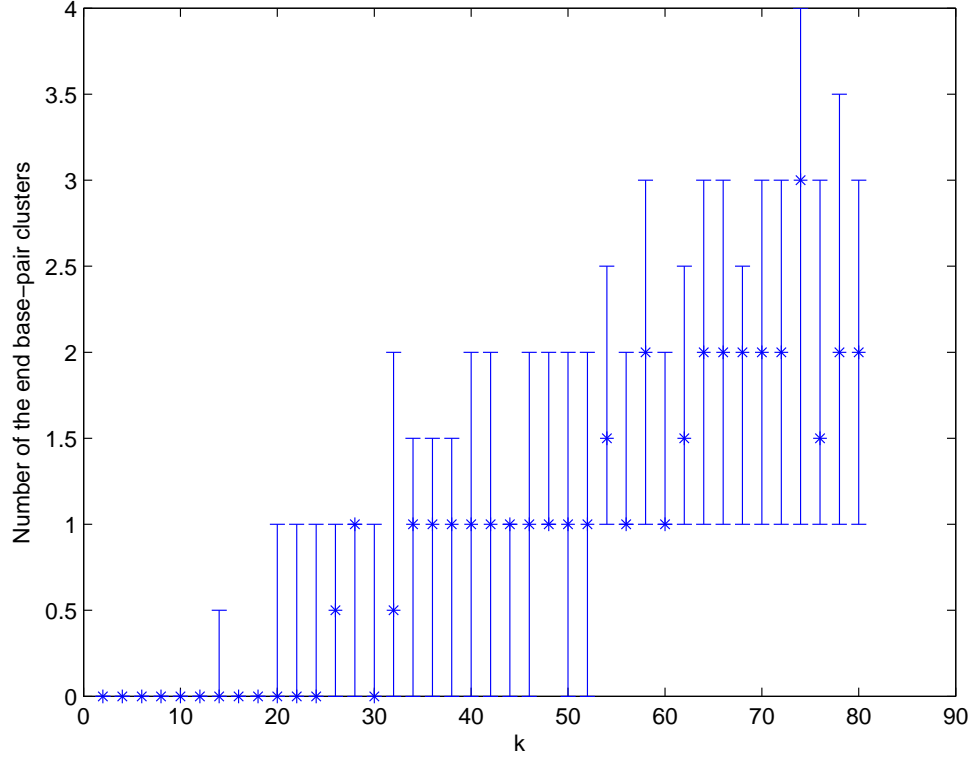
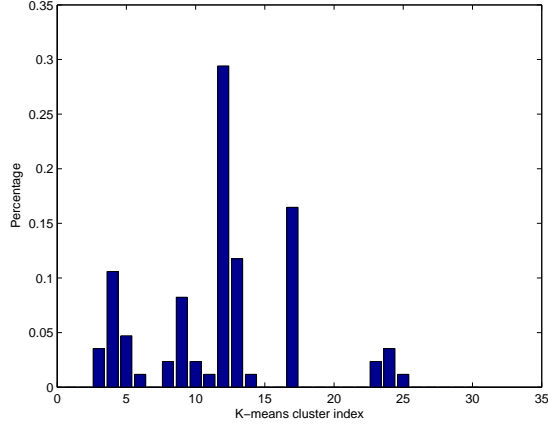
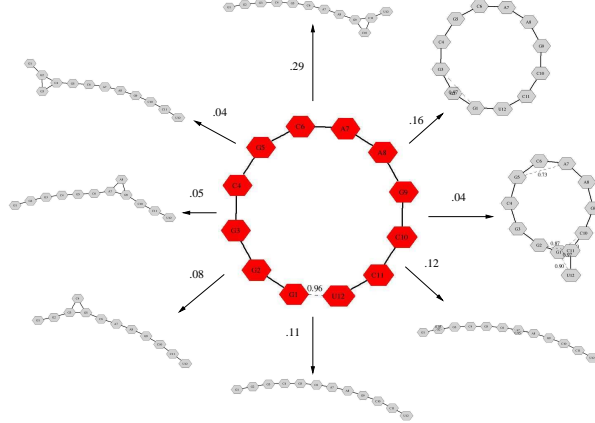


Figure 5: The number of end base-pair clusters found by K -means. Here k ranges from 2 to 80 with step 2. For each k , 20 experiments are repeated with K -means clustering. The number of clusters with end base-pair formed are recorded. The star is the median of such numbers and the bar delimits the distribution range from 10% to 90%. Starting from around $k = 25$, such clusters appear with at least 1/2 probability. Around $k = 55$, such clusters begin to split. The instability of K -means clusters is increasing as k grows, indicated by the expanding ranges.



(a)



(b)

Figure 6: K-means clustering fails to capture the low density intermediate states with one end base-pair formed. The illustration here chooses $k = 30$ for K -means clustering. (a) shows how end base-pair formed structures are distributed in different k-means clusters; (b) illustrates the mean structures of the top eight K-means clusters (gray) which contain base-pair formed structures. K -means splits the Mapper cluster and lumps them with densest clusters.