# Topological Analysis and Prediction of Biomolecular Data

Zixuan Cang[1], Lin Mu[3], Kedi Wu[1], Kristopher Opron[2], Kelin Xia[1], and Guo-wei Wei[1,2]

[1]Department of Mathematics, Michigan State University, MI 48824, USA
[2]Department of Biochemistry and Molecular Biology, Michigan State University, MI 48824, USA
[3]Oak Ridge National Laboratory, TN 37831, USA

## INTRODUCTION

Protein function and dynamics are closely related to its sequence and structure. However, the prediction of protein function and dynamics from its sequence and structure is still a fundamental challenge in molecular biology. Prediction of protein related observables provides advices for experiments and sheds light on how protein functions. Persistent homology is a new branch of algebraic topology that has found its success in the topological data analysis in a variety of disciplines, including molecular biophysics.

### FILTRATION OF SIMPLICIAL COMPLEX

Filtration of Vietoris-Rips complex built on $\alpha$-carbon point cloud of protein (ID:2LJC)

We explore the potential of using persistent homology as an independent tool for protein structure classification and protein-ligand/drug binding affinity prediction. From persistent homology computations, we extract protein topological fingerprints which are generated during a filtration process. We develop topological machine learning in which feature vectors are generated solely based on the output of topological fingerprints.

## PERSISTENT HOMOLOGY

In the past decade, persistent homology has been developed as a new multiscale representation of topological features.

**Simplex** A $k$-simplex denoted by $\sigma^k$ is a convex hull of $k+1$ vertices which is represented by a set of points

$$\sigma^k = \{\lambda_0 u_0 + \lambda_1 u_1 + ... + \lambda_k u_k | \sum \lambda_i = 1, \lambda_i \geqslant 0, i = 0, 1, ..., k\},$$

where $\{u_0, u_1, ..., u_k\} \subset \mathbb{R}^n$ is a set of affinely independent points.

**Simplicial complex** A simplicial complex $\mathcal{K}$ is a finite collection of simplices satisfying two conditions. First, faces of a simplex in $\mathcal{K}$ are also in $\mathcal{K}$; Secondly, intersection of any two simplices in $\mathcal{K}$ is a face of both the simplices. The highest dimension of simplices in $\mathcal{K}$ determines dimension of $\mathcal{K}$.

**Homology** For a simplicial complex $\mathcal{K}$, a $k$-chain is a formal sum of the form $\sum_{i=1}^{N} c_i [\sigma_i^k]$, where $[\sigma_i^k]$ is oriented $k$-simplex from $\mathcal{K}$. A boundary operator $\partial_k$ over a $k$-simplex $\sigma^k$ is defined as,

$$\partial_k \sigma^k = \sum_{i=0}^{k} (-1)^i [u_0, u_1, ..., \widehat{u_i}, ..., u_k],$$

where $[u_0, u_1, ..., \widehat{u_i}, ..., u_k]$ denotes the face obtained by deleting the $i$th vertex in the simplex. The boundary operator induces a boundary homomorphism $\partial_k : C_k(\mathcal{K}) \to C_{k-1}(\mathcal{K})$. The composition operator $\partial_{k-1} \circ \partial_k$ is a zero map,

$$\partial_{k-1}\partial_k(\sigma^k) = \sum_{j<i}(-1)^i(-1)^j[u_0, ..., \widehat{u_i} ... \widehat{u_j}, ... u_k] + \sum_{j>i}(-1)^i(-1)^{j-1}[u_0, ..., \widehat{u_j}, ... \widehat{u_i}, ... u_k]$$
$$= 0$$

A sequence of chain groups connected by boundary operation form a chain complex,

$$\cdots \longrightarrow C_n(\mathcal{K}) \xrightarrow{\partial_n} C_{n-1}(\mathcal{K}) \xrightarrow{\partial_{n-1}} \cdots \xrightarrow{\partial_1} C_0(\mathcal{K}) \xrightarrow{\partial_0} 0.$$

The equation $\partial_k \circ \partial_{k+1} = 0$ is equivalent to the inclusion $\text{Im}\partial_{k+1} \subset \text{Ker } \partial_k$, where Im and Ker denote image and kernel. Elements of $\text{Ker}\partial_k$ are called $k$th cycle group, and denoted as $Z_k = \text{Ker}\partial_k$. Elements of $\text{Im}\partial_{k+1}$ are called $k$th boundary group, and denoted as $B_k = \text{Im}\partial_{k+1}$. A $k$th homology group is defined as the quotient group of $Z_k$ and $B_k$.

$$H_k = Z_k / B_k.$$

The $k$th Betti number of simplicial complex $\mathcal{K}$ is the rank of $H_k$,

$$\beta_k = \text{rank}(H_k) = \text{rank}(Z_k) - \text{rank}(B_k).$$

Betti number $\beta_k$ is finite number, since $\text{rank}(B_p) \leqslant \text{rank}(Z_p) < \infty$. Betti numbers computed from homology group are used to describe the corresponding space.

## FEATURE CONSTRUCTION FOR STRUCTURE CLASSIFICATION TASKS

A number of features describing properties of the samples in different scales are extracted from persistent homology bar codes. Illustrated below are some examples of features used in machine learning.

- The length of the second longest Betti 0 bar.
- The summation of lengths of all Betti 0 bars except for those exceed the max filtration value.
- The average length of Betti 0 bars except for those exceed the max filtration value.
- The onset value of the longest Betti 1 bar.
- The number of Betti 1 bars that locate at [4.5, 5.5], divided by the number of atoms.
- The onset value of the first Betti 2 bar that ends after a given number.

Bar codes in different dimension with different lifespan, birth time, and death time show characteristics of the sample from different scales. With persistent homology, we are able to describe local properties like alpha helices or beta sheets and global properties like size of the cavity of a spherical structure or size of tunnel of a cylindrical structure.

## PERFORMANCE ON CLASSIFICATION TASKS

Influenza A virus drug inhibition: 96% Accuracy

Protein secondary structures: 85% Accuracy

Hemoglobins in their relaxed and taut forms: 80% accuracy

**SCOPe**
55 classification tasks of protein superfamilies over 1357 proteins from *Protein Classification Benchmark Collection*: 82% Accuracy

## PERFORMANCE ON BINDING AFFINITY PREDICTION



The persistent homology based protein-ligand/drug binding affinity predictor named T-Score is tested on the PDBBind v2007 core set with the v2007 refined set as a training set where the testing set has been excluded from the training set. A high Pearson correlation of 0.80 is achieved and our T-Score outperforms all the other eminent methods in computational biophysics.

## REFERENCES

- Zixuan Cang, Lin Mu, Kedi Wu, Kristopher Opron, Kelin Xia and Guo-Wei Wei, "A topological approach to protein classification", Molecular Based Mathematical Biology, 3, 140-62 (2015).
- Kelin Xia and Guo-Wei Wei, "Persistent homology analysis of protein structure, flexibility and folding", International Journal for Numerical Methods in Biomedical Engineering, 30(8):814-844 (2014).
- Kelin Xia, Zhixiong Zhao and Guo-Wei Wei, "Multiresolution persistent homology for excessively large biomolecular datasets", Journal of Chemical Physics, 143, 134103 (2015).

## PROTEIN-LIGAND/DRUG BINDING FREE ENERGY

Structure-based drug design relies on computational methods to identify and optimize potential drugs. Molecule docking is the most widely used approach which predicts the location and orientation (pose) of a ligand bound to a protein to form a stable complex. In the process of search for a pose for the ligand, a scoring function which measures the binding affinity between the two molecules is needed to distinguish the favorable poses from the unfavorable ones. An accurate and efficient protein-ligand/drug binding affinity predictor is therefore the key of molecule docking process. As the dominating forces that regulates protein-ligand binding are mainly weak forces which heavily depends on spacial arrangements, persistent homology becomes a competitive candidate for this job.

## FLOWCHART OF PERSISTENT HOMOLOGY BASED PROTEIN-LIGAND BINDING AFFINITY PREDICTION



## CONCLUSION

We test the performance of persistent homology in various protein structure classification tasks as well as protein-ligand binding/drug affinity prediction. It is found that persistent homology is able to offer a power representation of proteins and capture their intrinsic interactions. Our persistent homology based T-Score outperforms all the other eminent methods in computational biophysics on the blind prediction of protein-ligand/drug binding affinities.

## ACKNOWLEDGMENT