

# Topological modelling and analysis for biomolecular data: 3 TDA-based machine learning

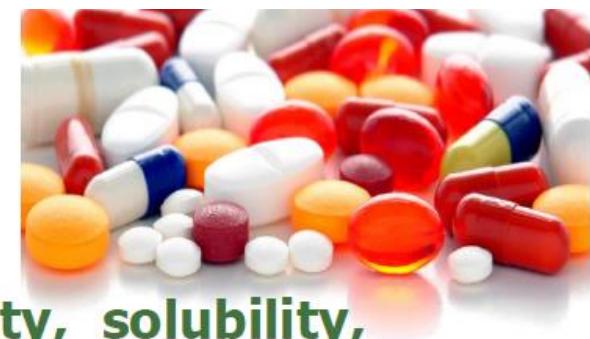
Kelin Xia

*School of Physical and Mathematical Sciences,  
Nanyang Technological University*

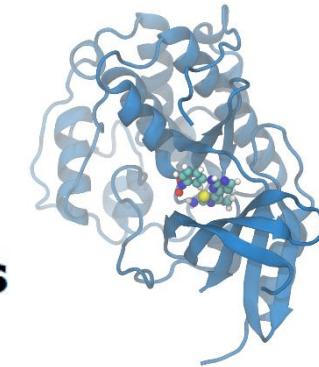
**应用拓扑短期课程与报告会, 12月13-16日, 2020**

**Fund: NTU-JSPS(2019), MOE-Tier 1(2018,2019), MOE-Tier 2(2018,2021), Alibaba-NTU(2020), Merlion(2020)**

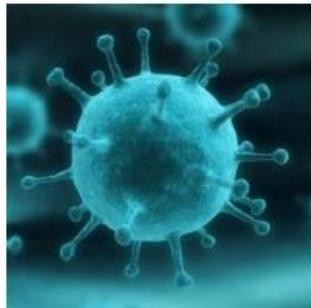
# Drug design and discovery



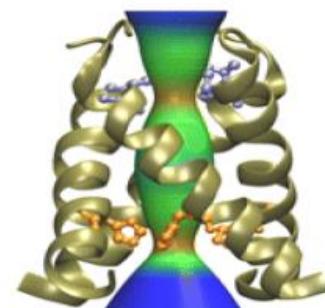
- 1) Disease identification (physiology)
- 2) Target hypothesis (biochem./mole. biol.)
- 3) Virtual screening: drug pose, binding affinity, solubility, partition coefficient, toxicity, and side-effects (biophysics/bioinformatics)
- 4) Drug structural optimization in the target binding site (biochemistry/biophysics/synthetic chem.)
- 5) Preclinical *in vitro* and *in vivo* test
- 6) Clinical trials
- 7) Optimize drug's efficacy, pharmacokinetics, and pharmacodynamics properties (quantitative systems pharmacology)



Influenza -- flu virus



M2 channel



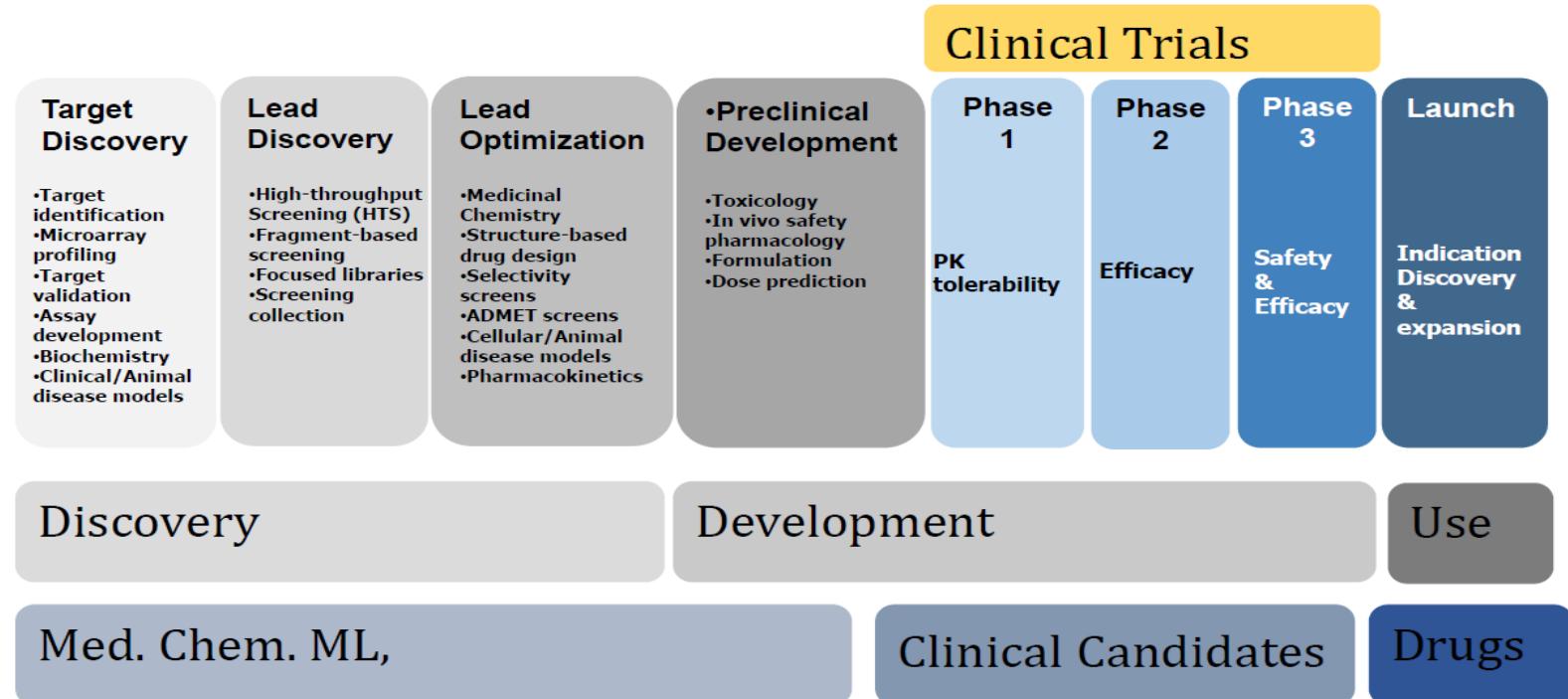
Amantadine



M2-A complex

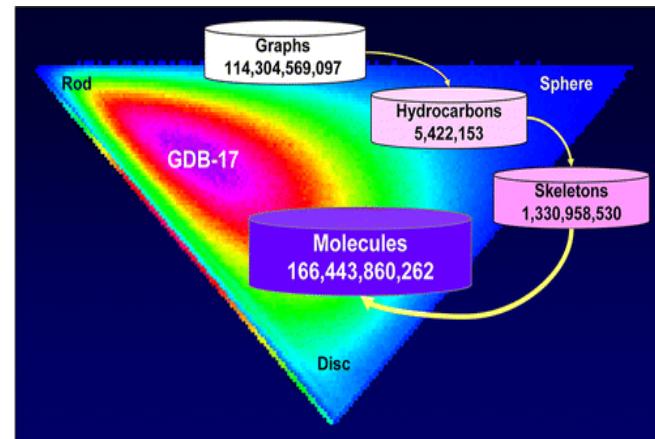


# Drug Discovery Process (simplified)

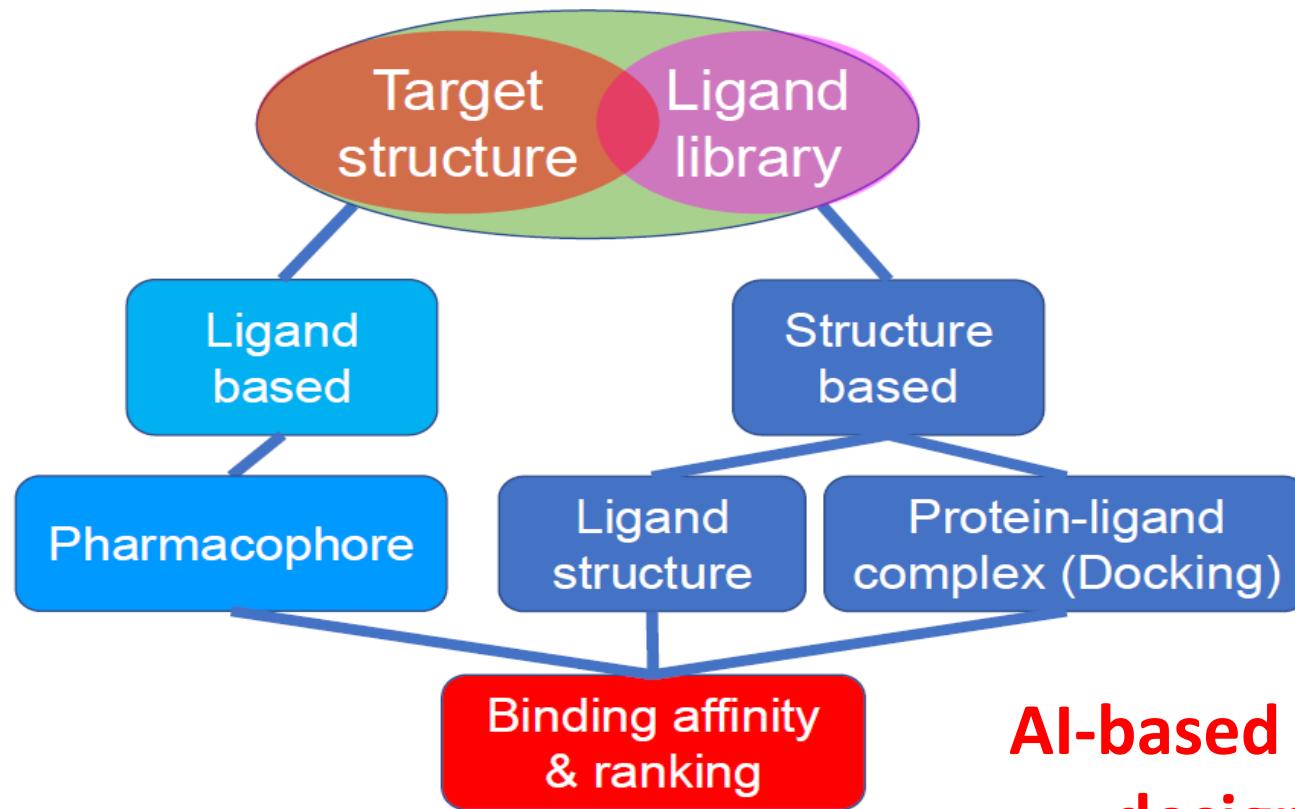


**“...166.4 billion molecules of up to 17 atoms of C, N, O, S, and halogens forming the chemical universe database GDB-17”**

Ruddigkeit, Lars, et al. "Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17." Journal of chemical information and modeling 52(11): 2864-2875, 2012.



# Virtual Screening



Force field

Empirical

Knowledge-based

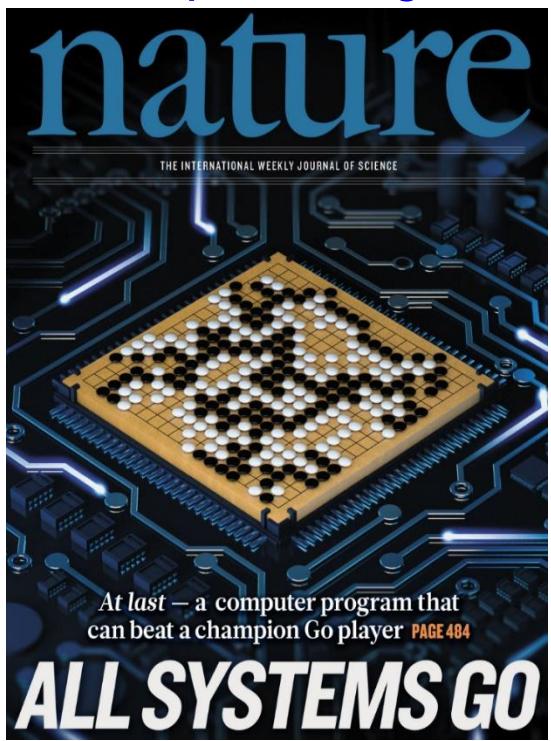
Machine learning

P. J. Ballester and J. B. Mitchell, "A machine learning approach to predicting protein–ligand binding affinity with applications to molecular docking," *Bioinformatics*, 26 (9): 1169–1175, 2010

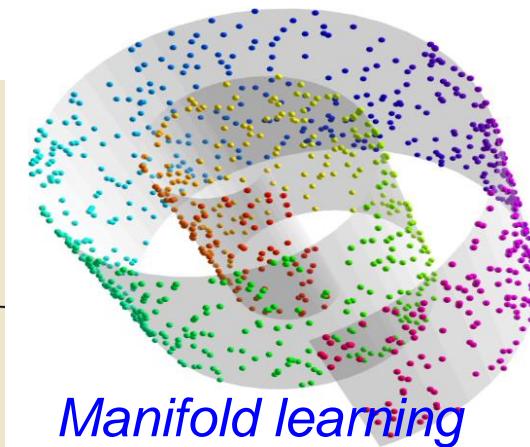
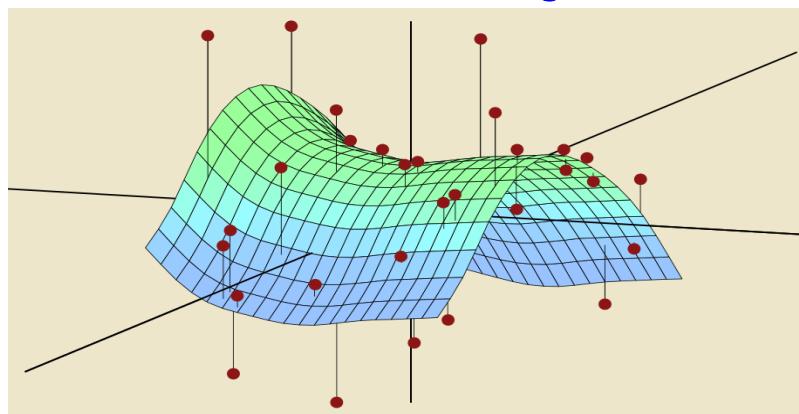
Q. U. Ain, A. Aleksandrova, F. D. Roessler, and P. J. Ballester, "Machine-learning scoring functions to improve structure-based binding affinity prediction and virtual screening," *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 5(6): 405–424, 2015.

Y. C. Lo, S. E. Rensi, W. Torng, and R. B. Altman, "Machine learning in chemoinformatics and drug discovery," *Drug discovery today*, 23 (8), 1538–1546, 2018.

Deep learning



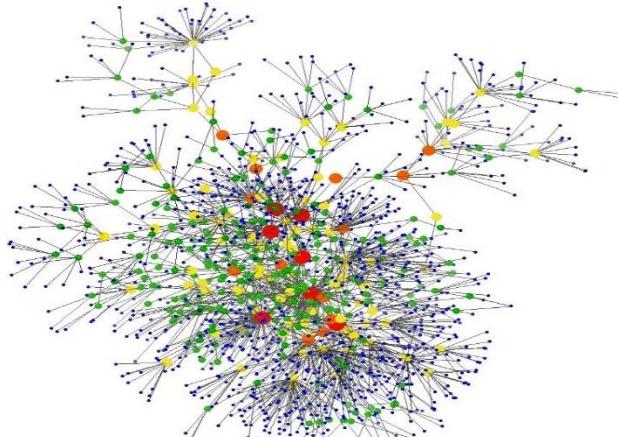
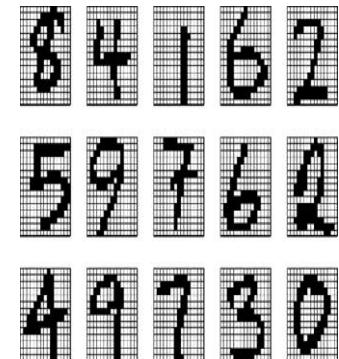
Statistic learning



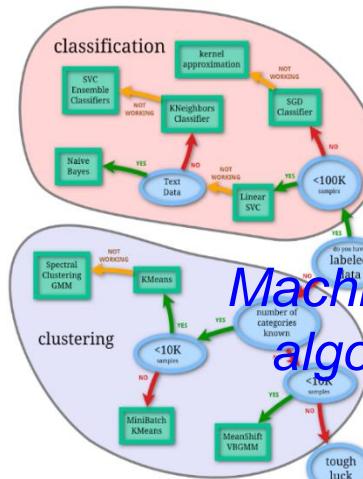
Manifold learning

# How to deal with big data?

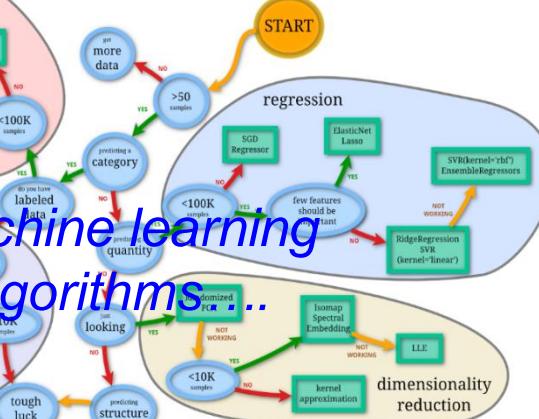
Pattern recognition



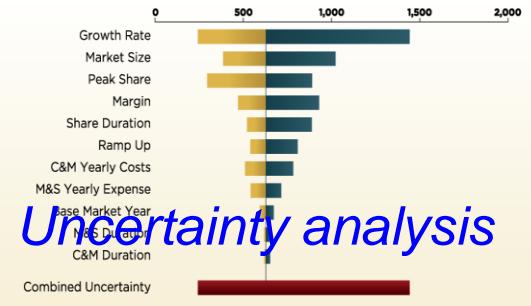
Graph modeling and analysis



Machine learning algorithms...



Impact of Uncertainty



Uncertainty analysis

# Representation and Feature Learning

**“The success of machine learning algorithms generally depends on data representation...”**

*Y. Bengio, etc, “Representation Learning: A Review and New Perspectives”*

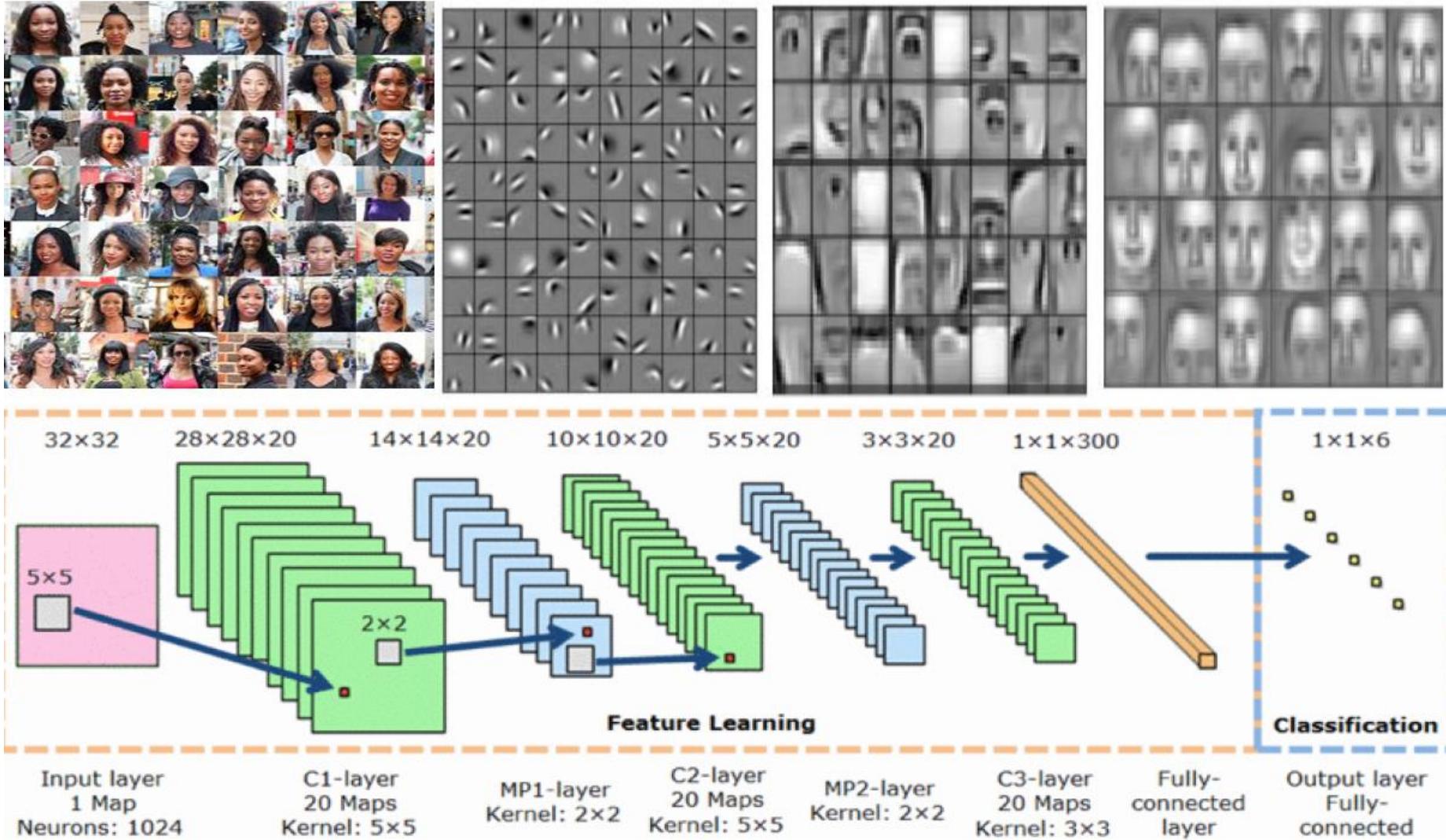


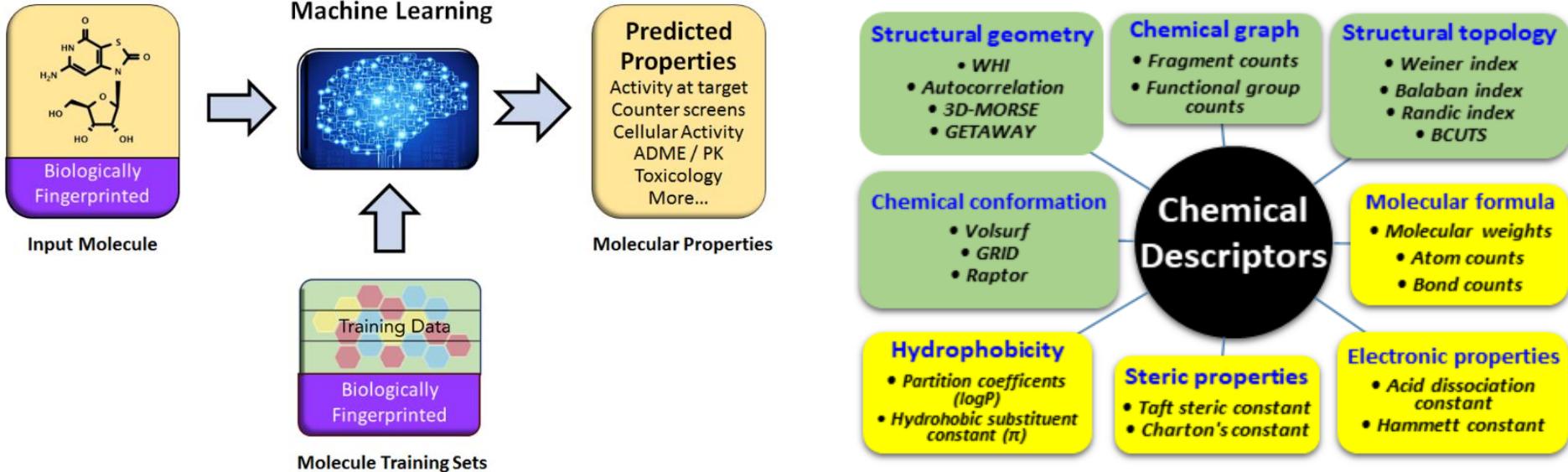
**“The deep learning research aims at discovering learning algorithms that discover multiple levels of distributed representations...”**

*Y. Bengio, “Deep Learning of Representations: Looking Forward”*

# Feature learning is key to data analysis!

Fukushima (1980) – Neo-Cognitron; LeCun (1998) – Convolutional Neural Networks (CNN);...





**Molecular (chemical) descriptors (>5000) directly determine the performance of learning models!**

### Common chemical descriptors for QSAR/QSPR analysis

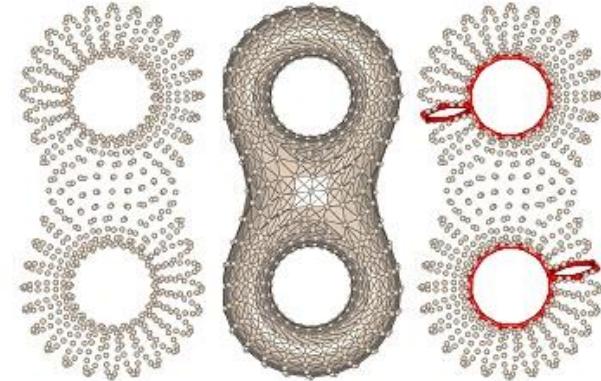
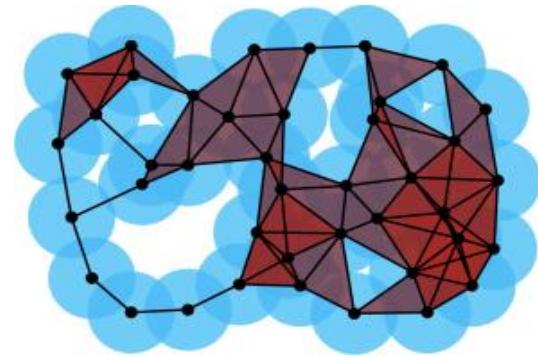
Chemical descriptors	Based on	Examples
Theoretical descriptors		
0D	Molecular formula	Molecular weights, atom counts, bond counts
1D	Chemical graph	Fragment counts, functional group counts
2D	Structural topology	Weiner index, Balaban index, Randic index, BCUTS
3D	Structural geometry	WHIM, autocorrelation, 3D-MORSE, GETAWAY
4D	Chemical conformation	Volsurf, GRID, Raptor
Experimental descriptors		
Hydrophobic parameters	Hydrophobicity	Partition coefficients ( $\log P$ ), hydrophobic substituent constant ( $\pi$ )
Electronic parameters	Electronic properties	Acid dissociation constant, Hammett constant
Steric parameters	Steric properties	Taft steric constant, Charton's constant

Fingerprint	Description	Number of features	Package
FP2	A path-based fingerprint which indexes small molecule fragments based on linear segments of up to 7 atoms <sup>20</sup>	256	RDKit <sup>23</sup>
Daylight	A path-based fingerprint consisting of 2048 bits and encoding all connectivity pathways in a given length through a molecule <sup>21</sup>	2048	
MACCS	A substructure keys-based fingerprint with 166 structural keys based on SMARTS patterns <sup>19</sup>	166	
Estate1	A topological fingerprint based on electro-topological State Indices, which encodes the intrinsic electronic state of the atom as perturbed by the electronic influence of all other atoms in the molecule within the context of the topological character of the molecule. Estate 1 represents the number of times each atom type is hit <sup>22</sup>	79	
Estate2	Similar to estate 1, however it contains the sum of the EState indices for atoms of each type <sup>22</sup>	79	
ECFP4	The de facto standard circular fingerprint based on the Morgan algorithm, <sup>42</sup> which uses an iterative process to assign numeric identifiers to each atom <sup>15</sup>	2048	
Pharm2D	Each bit corresponds to a particular combination of features and interactions needed for a molecule to be active against a given target <sup>23</sup>	990	
ERG	A Pharmacophore fingerprint, which is an extended reduced graph approach using pharmacophore-type node descriptions to encode the relevant molecular properties <sup>24</sup>	315	

# Topological Data Analysis (TDA)

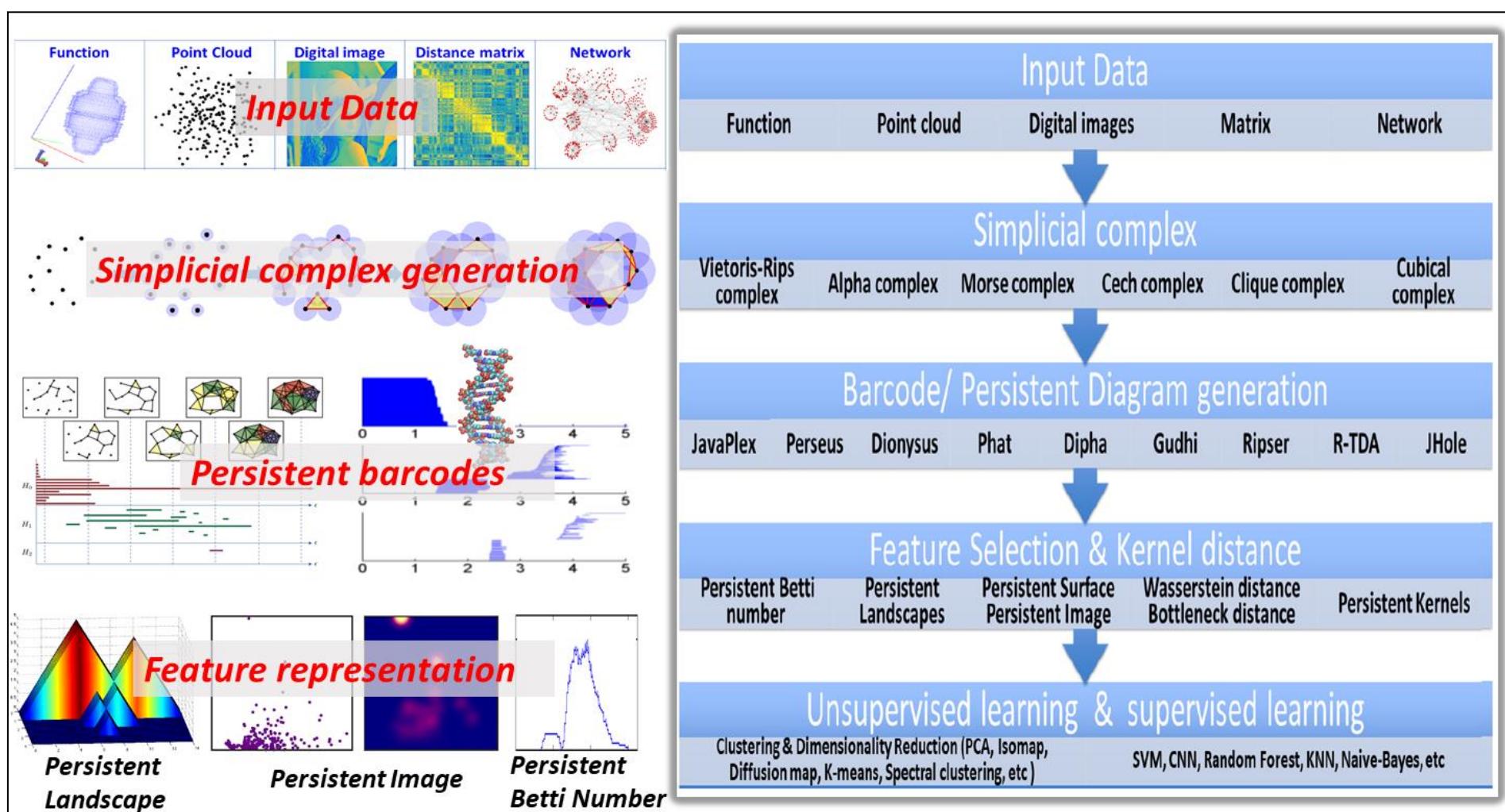
**Topological invariant;  
Homology;  
Homotopy;  
Simplicial complex;  
Morse theory;  
Reeb graph;**

.....



Computational Geometry;  
Computational topology;  
Algebraic topology

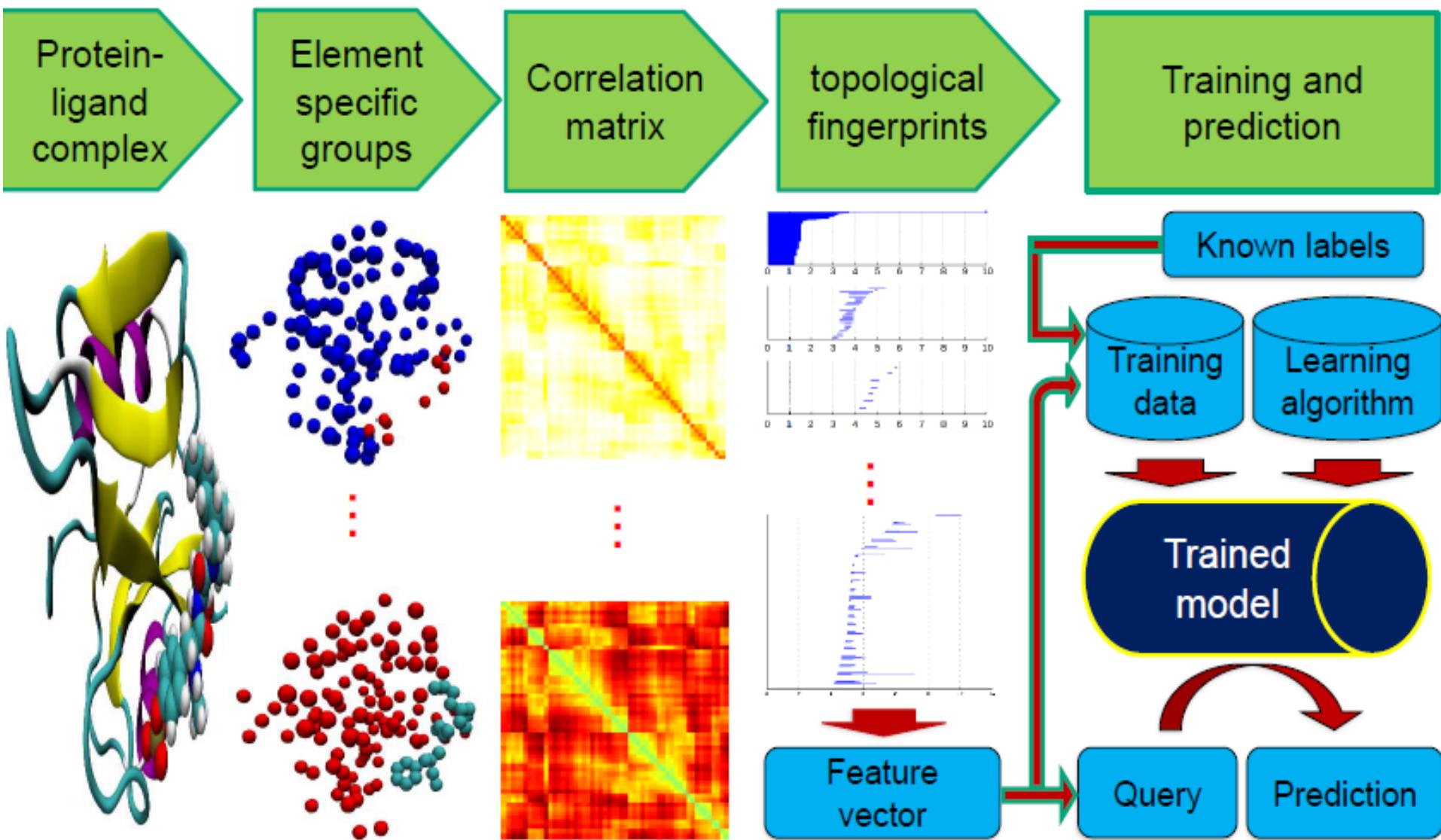
# TDA based machine learning models



(Pun, Xia and Lee, arXiv, 2018)

# Topology based learning architecture

(Cang & Wei, IJNMBE, 2017)



# Recent progress in TDA based drug design

Collaborator  
Guowei Wei  
MSU, USA



MORE AT SIAM | siam news | HOME | HAPPENING NOW | GET INVOLVED | RESEARCH

Get Involved | September 01, 2016

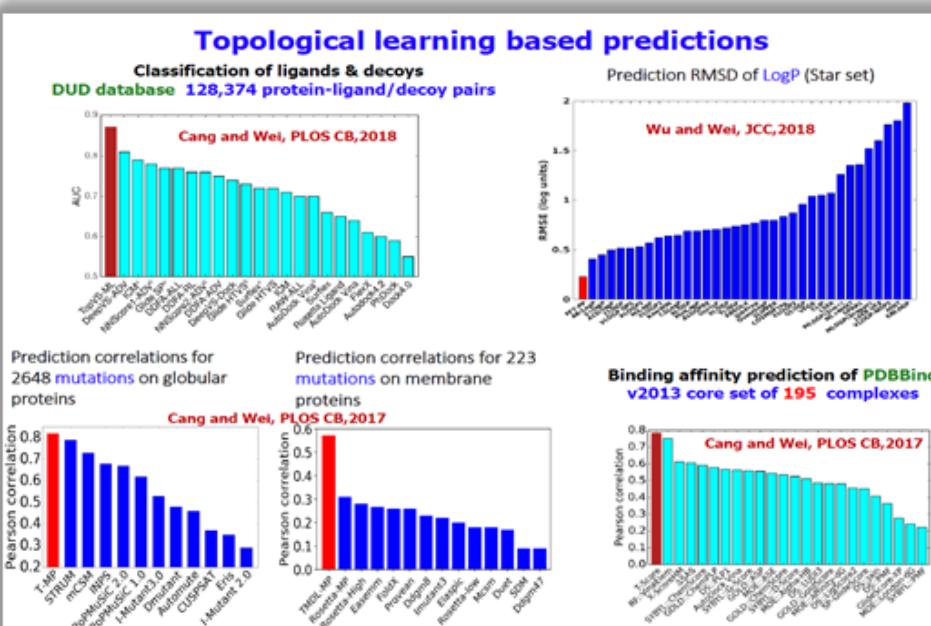
## Mathematical Molecular Bioscience and Biophysics

A Recurring Theme at the SIAM Conference on the Life Sciences

MORE AT SIAM | siam news | HOME | HAPPENING NOW | GET INVOLVED | RESEARCH

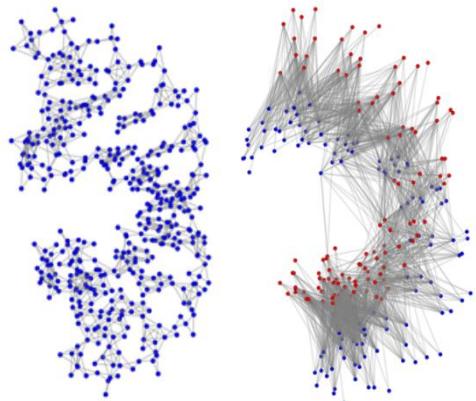
Research | December 01, 2017

## Persistent Homology Analysis of Biomolecular Data

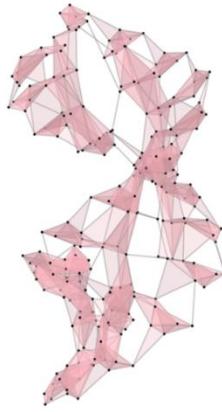


# Mathematical representations for molecules

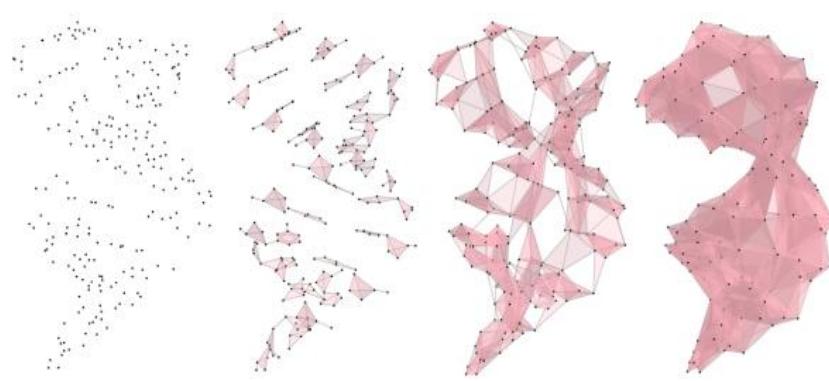
Graph



Simplicial complex



Multiscale simplicial complexes



## Graph models and measurements:

Graph Laplacian; Fiedler Eigenvalue; Fiedler eigenvector; Shortest path; Clique; Cluster coefficient; Closeness; Centrality; Betweenness; Modularity; Cheeger constant; Erdős number; Percolation...

## Advanced mathematical representations and intrinsic invariants

## Simplicial complex models and measurements:

Combinatorial Laplacian; Hodge theory; Betti number; Euler characteristics; Homology; Cohomology; Morse theory; Knot polynomials...

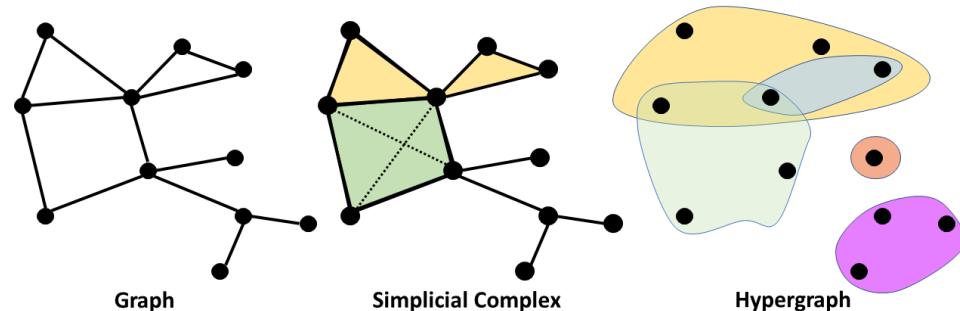
## Multiscale simplicial complex:

Persistent homology; Persistent cohomology...

# Persistent Spectral (PerSpect)

## Spectral models

- **Spectral graph**
- **Spectral simplicial complex**
- **Spectral hypergraph**



## Filtration

- **Nested sequence of Graphs**
- **Nested sequence of Simplicial Complexes**
- **Nested sequence of Hypergraphs**

## Spectral models + Filtration

- ❖ **Persistent spectral (PerSpect) graph**
- ❖ **Persistent spectral (PerSpect) simplicial complexes**
- ❖ **Persistent spectral (PerSpect) hypergraph**

F. Chung, and S. T. Yau. "A Harnack inequality for homogeneous graphs and subgraphs." Comm. Anal. Geom 2.4 (1994): 627-640.

F. Chung, "Spectral graph theory". American Mathematical Society, 1997

D. Spielman, "Spectral graph theory", Combinatorial scientific computing. No. 18. Boca Raton, FL: CRC Press, 2012.

D. Horak, and J. Jost, "Spectra of combinatorial Laplace operators on simplicial complexes". Advances in Mathematics, 244, 303–336, 2013

O. Parzanchevski, and R. Ron "Simplicial complexes: spectrum, homology and random walks." Random Structures & Algorithms 50.2 (2017): 225-261.

M. T. Schaub, et al. "Random walks on simplicial complexes and the normalized Hodge 1-Laplacian." SIAM Review 62.2 (2020): 353-391.

K. Q. Feng, "Spectra of hypergraphs and applications." Journal of number theory 60.1 (1996): 1-22.

J. Cooper, and D. Aaron. "Spectra of uniform hypergraphs." Linear Algebra and its applications 436.9 (2012): 3268-3292.

L. Q. Qi, and Z. Y. Luo. "Tensor analysis: spectral theory and special tensors". Society for Industrial and Applied Mathematics, 2017

.....

# Combinatorial Laplacian

Let  $K$  be a simplicial complex and  $C_k(K)$  be a vector space over some field  $\mathbb{F}$  whose basis is all  $k$ -simplices of  $K$ .

## Definition

The *dual* of  $C_k(K)$ , denoted by  $C^k(K)$ , is the set of all linear functionals on  $C_k(K)$ :

$$C^k(K) = \{\phi : C_k(K) \rightarrow \mathbb{F} : \phi \text{ is linear}\}.$$

**Note:** Both  $C_k(K)$  and  $C^k(K)$  have the same dimension = no. of  $k$ -simplices of  $K$ .

① *Boundary map*  $\partial_k : C_k(K) \rightarrow C_{k-1}(K)$

$$\partial_k([u_0, u_1, \dots, u_k]) = \sum_{i=0}^k (-1)^i [u_0, \dots, u_{i-1}, u_{i+1}, \dots, u_k],$$

where  $[u_0, u_1, \dots, u_k]$  is a basis element of  $C_k(K)$ .

② *Coboundary map*  $\delta_k : C^k(K) \rightarrow C^{k+1}(K)$

$$\delta_k(\phi)(\sigma^{k+1}) = \sum_{i=0}^{k+1} (-1)^i \phi([u_0, \dots, u_{i-1}, u_{i+1}, \dots, u_{k+1}]),$$

where  $\phi \in C^k(K)$  and  $\sigma^{k+1} = [u_0, \dots, u_{k+1}]$  is a basis element of  $C_{k+1}(K)$ .

## Combinatorial Laplacian

Another important map that is crucial in the formulation of Hodge Decomposition Theorem is the combinatorial Laplacian:

## Definition

The  $k$ -dimensional *combinatorial Laplacian* is the linear operator  $\Delta_k : C^k(K) \rightarrow C^k(K)$  is defined as follows:

$$\Delta_k = \begin{cases} \delta_k^* \circ \delta_k + \delta_{k-1} \circ \delta_{k-1}^* & \text{if } k \geq 1, \\ \delta_k^* \circ \delta_k & \text{if } k = 0. \end{cases}$$

$\delta_k^* : C^{k+1}(K) \rightarrow C^k(K)$  is the *adjoint/transpose map* of  $\delta_k$  where

$$\langle \delta_k(f), g \rangle = \langle f, \delta_k^*(g) \rangle$$

for every  $f \in C^k(K)$ ,  $g \in C^{k+1}(K)$  and a suitable inner product  $\langle , \rangle$  for  $C^k(K)$  and  $C^{k+1}(K)$ .

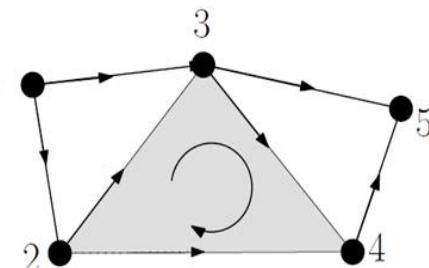
## Boundary matrix

$$B_k(i, j) = \begin{cases} 1, & \text{if } \sigma_i^{k-1} \subset \sigma_j^k \text{ and } \sigma_i^{k-1} \sim \sigma_j^k \\ -1, & \text{if } \sigma_i^{k-1} \subset \sigma_j^k \text{ and } \sigma_i^{k-1} \not\sim \sigma_j^k \\ 0, & \text{if } \sigma_i^{k-1} \not\subset \sigma_j^k. \end{cases}$$

face

same orientation

**For simplicial complexes, if they are upper adjacent, they must be lower adjacent!**



$$\mathbf{B}_1 = \left[ \begin{array}{ccccccc|c} [12] & [13] & [23] & [24] & [34] & [35] & [45] & \\ \hline -1 & -1 & 0 & 0 & 0 & 0 & 0 & [1] \\ 1 & 0 & -1 & -1 & 0 & 0 & 0 & [2] \\ 0 & 1 & 1 & 0 & -1 & -1 & 0 & [3] \\ 0 & 0 & 0 & 1 & 1 & 0 & -1 & [4] \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & [5] \end{array} \right]$$

## Combinatorial Laplacian (Hodge Laplacian)

$$\mathbf{L}_k = \mathbf{B}_k^T \mathbf{B}_k + \mathbf{B}_{k+1} \mathbf{B}_{k+1}^T.$$

upper degree  
(number of cofaces)

lower degree  
(number of faces)

$$L_k(i, j) = \begin{cases} d(\sigma_i^k) + k + 1, & \text{if } i = j \\ 1, & \text{if } i \neq j, \sigma_i^k \not\sim \sigma_j^k, \sigma_i^k \curvearrowleft \sigma_j^k \text{ and } \sigma_i^k \sim \sigma_j^k \\ -1, & \text{if } i \neq j, \sigma_i^k \not\sim \sigma_j^k, \sigma_i^k \curvearrowleft \sigma_j^k \text{ and } \sigma_i^k \not\sim \sigma_j^k \\ 0, & \text{if } i \neq j, \sigma_i^k \curvearrowleft \sigma_j^k \text{ or } \sigma_i^k \not\sim \sigma_j^k. \end{cases}$$

k>0

lower adjacent

upper adjacent

$$\mathbf{B}_2 = \left[ \begin{array}{c|ccccc} [2, 3, 4] & & & & & \\ \hline 0 & [1, 2] & & & & \\ 0 & [1, 3] & & & & \\ 1 & [2, 3] & & & & \\ -1 & [2, 4] & & & & \\ 1 & [3, 4] & & & & \\ 0 & [3, 5] & & & & \\ 0 & [4, 5] & & & & \end{array} \right]$$

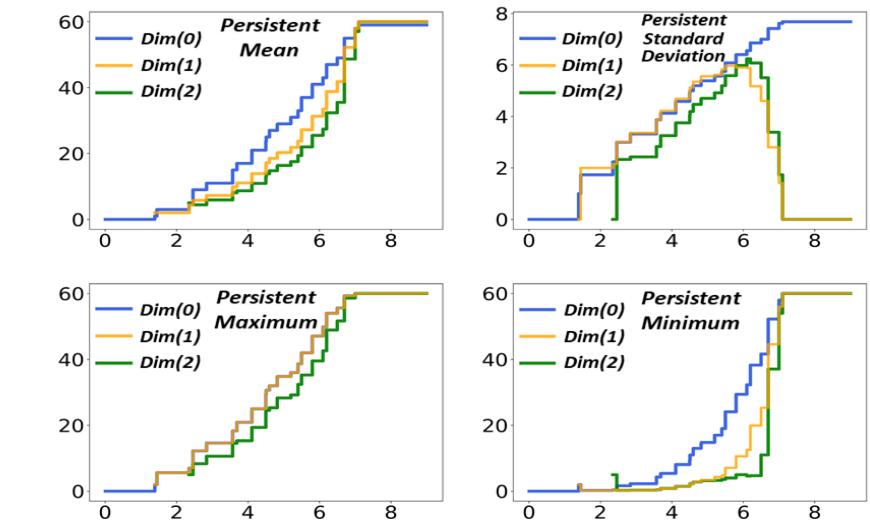
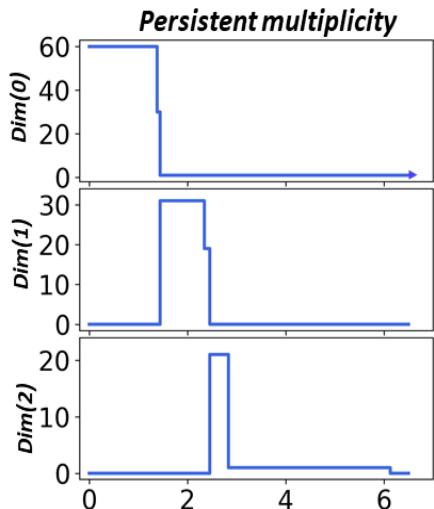
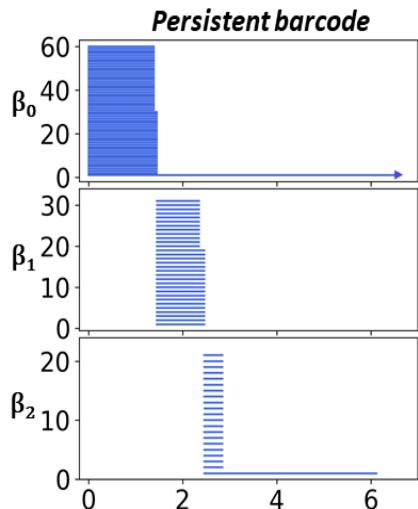
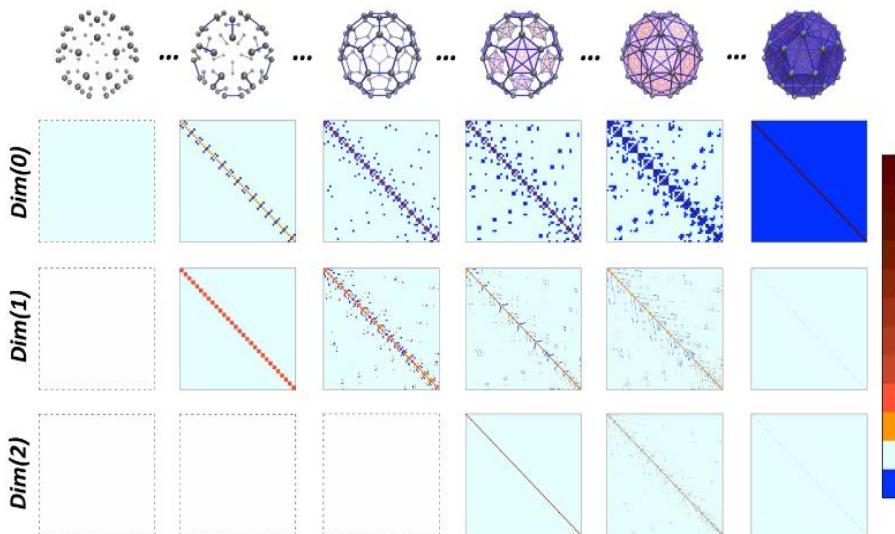
The Laplacians are computed as

$$\mathcal{L}_0 = \begin{pmatrix} 2 & -1 & -1 & 0 & 0 \\ -1 & 3 & -1 & -1 & 0 \\ -1 & -1 & 4 & -1 & -1 \\ 0 & -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & -1 & 2 \end{pmatrix}, \quad \mathcal{L}_1 = \begin{pmatrix} 12 & 13 & 23 & 24 & 34 & 35 & 45 \\ \hline 2 & 1 & -1 & -1 & 0 & 0 & 0 \\ 1 & 2 & 1 & 0 & -1 & -1 & 0 \\ -1 & 1 & 3 & 0 & 0 & -1 & 0 \\ -1 & 0 & 0 & 3 & 0 & 0 & -1 \\ 0 & -1 & 0 & 0 & 3 & 1 & -1 \\ 0 & -1 & -1 & 0 & 1 & 2 & 1 \\ 0 & 0 & 0 & -1 & -1 & 1 & 2 \end{pmatrix}, \quad \mathcal{L}_2 = 3.$$

# Persistent spectral simplicial complex

**Combinatorial Laplacian  
(Hodge Laplacian)**

$$\mathbf{L}_k = \mathbf{B}_k^T \mathbf{B}_k + \mathbf{B}_{k+1} \mathbf{B}_{k+1}^T.$$

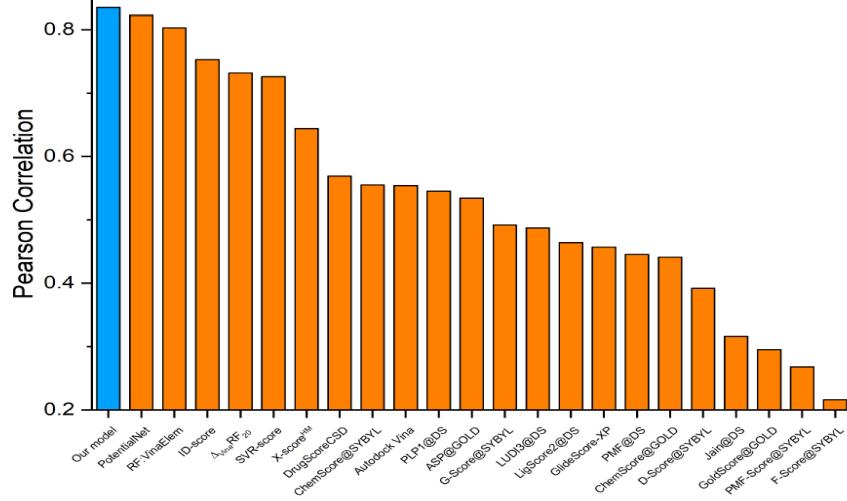


**Multiplicity of zero eigenvalues (Persistent multiplicity) from PerSpect simplicial complex is equivalent to persistent Betti number.**

**PerSpect variables change with filtration parameter and incorporate in them related geometric information.**

# Ours: 0.836

# *Dataset 2007*



**Ours: 0.793**

# *Dataset 2013*

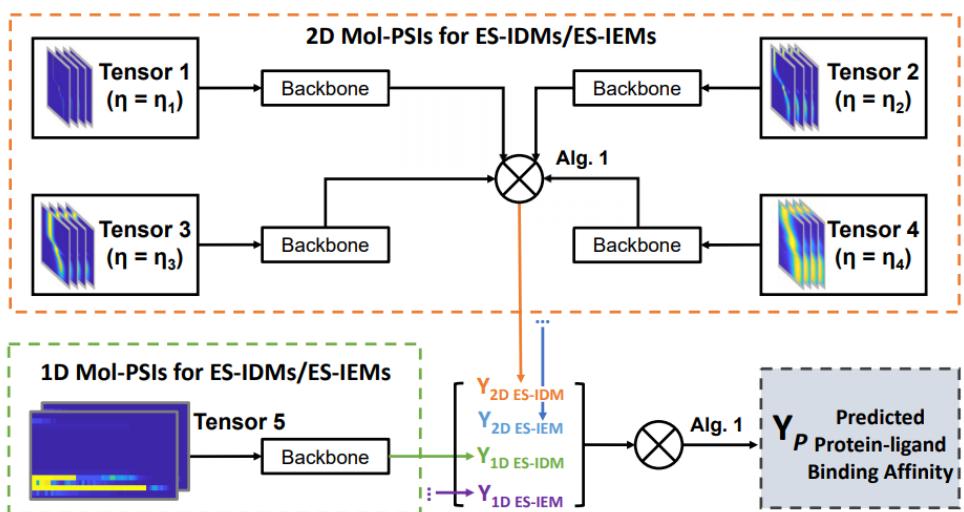
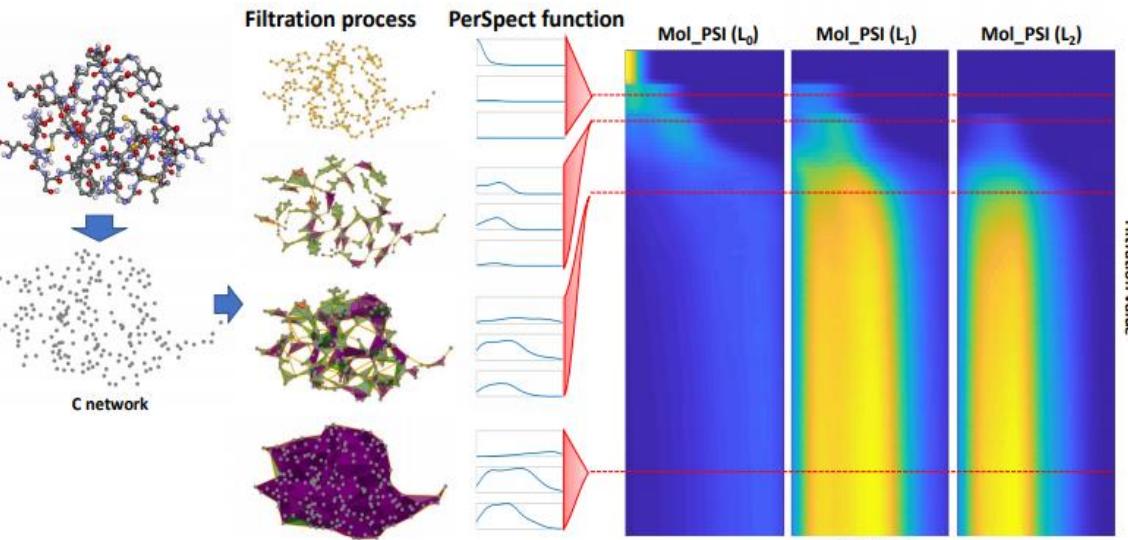
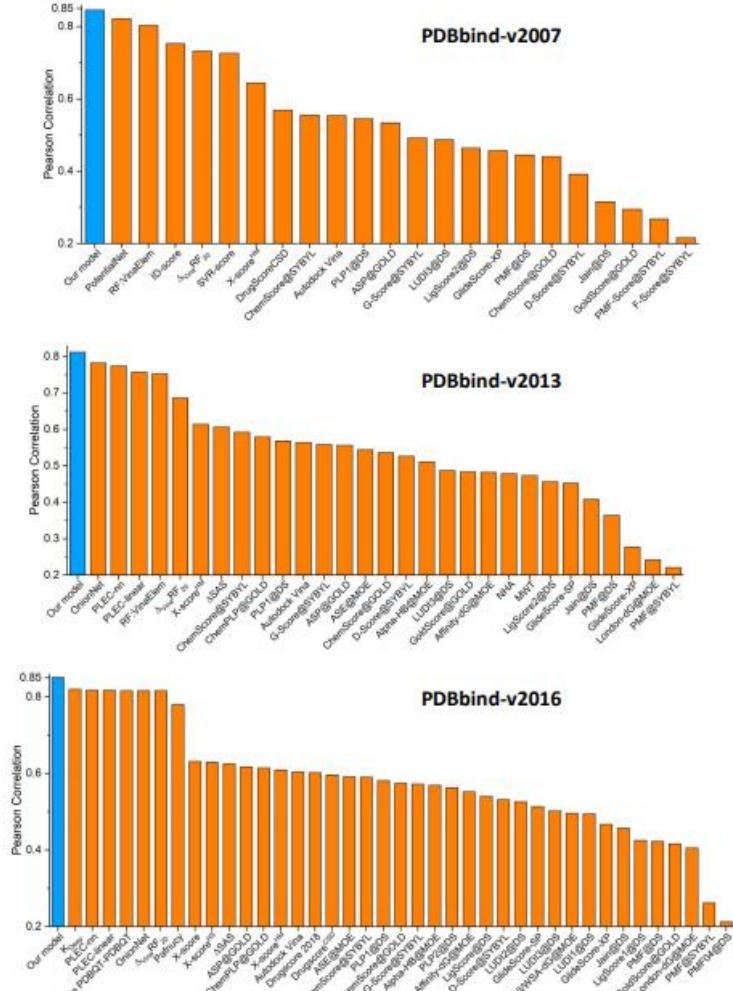
# Benchmark testing with PDBbind datasets

*Model setting:*  
*Spectral vectors*  
+  
*Random forest*

Ours: 0.840

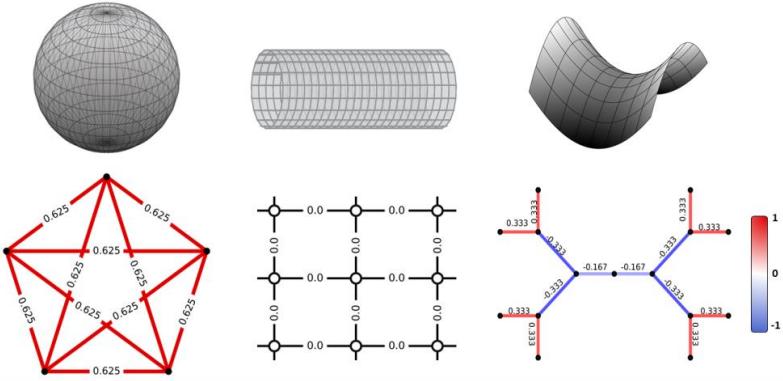
## *Dataset 2016*

# Molecular persistent spectral image (Mol-PSI) based CNNs

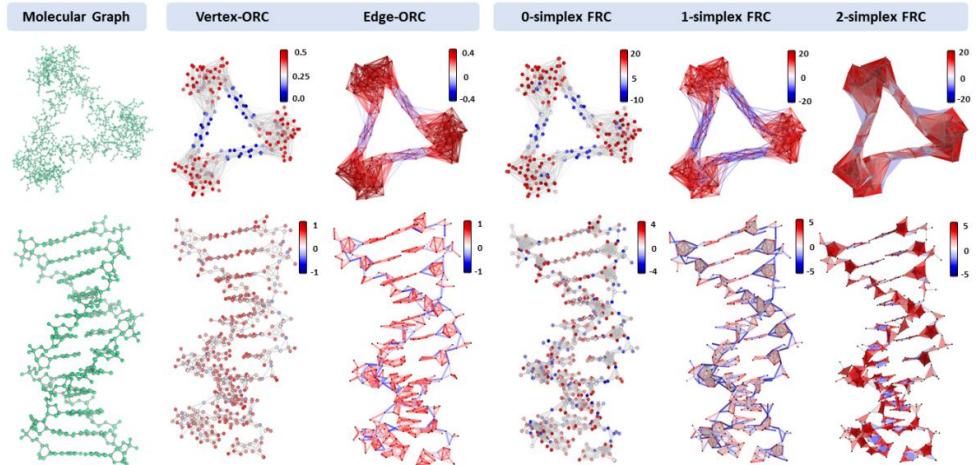


Joint work with Alibaba group (Peiran Jiang, Zhenyu Meng, Ying Chi, Lei Zhang, Xian-Sheng Hua, Kelin Xia, submitted, 2020)

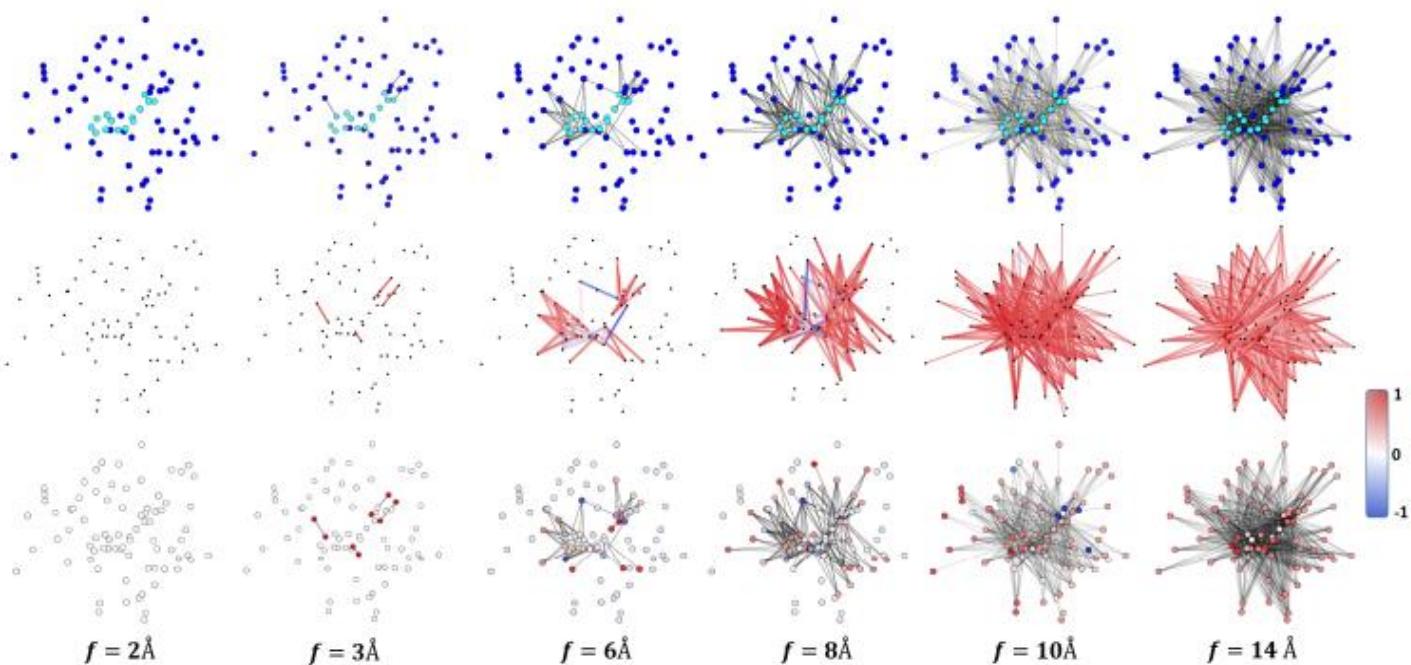
## ORCs for three graphs



## Ollivier Ricci curvature (ORC) VS Forman Ricci curvature (FRC)



## Persistent Ricci curvature



# Olliver Ricci curvature (ORC)

Vertex-x related probability measure:

constant between 0 and 1

$$m_x(x_i) = \begin{cases} \alpha & \text{if } x_i = x. \\ (1 - \alpha)/k_x & \text{if } x_i \in \Gamma_x \\ 0 & \text{otherwise.} \end{cases}$$

set of all the neighboring vertices of x

$L^1$  Wasserstein distance:

distance between two vertices

$$W_1(m_x, m_y) = \inf_{\xi} \sum_{x_i \in V} \sum_{y_j \in V} d(x_i, y_j) \xi(x_i, y_j)$$

$$\sum_{y_j \in V} \xi(x, y_j) = m_x$$

transportation distance

$$\sum_{x_i \in V} \xi(x_i, y) = m_y$$

Ollivier Ricci curvature:

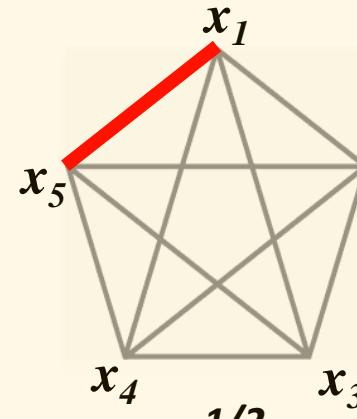
$$c(x, y) := 1 - \frac{W_1(m_x, m_y)}{d(x, y)}$$

An example

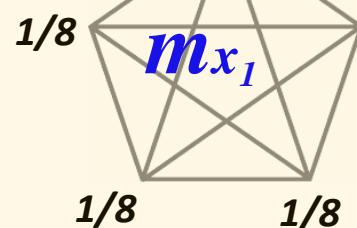
set  $\alpha = 1/2$ , check edge  $(x_1, x_5)$

$$\Gamma_{x_1} = \{x_2, x_3, x_4, x_5\}$$

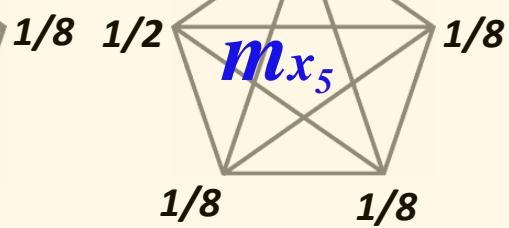
$$\Gamma_{x_5} = \{x_2, x_3, x_4, x_1\}$$



1/2



**$m_{x_1}$**



1/8

**$m_{x_5}$**

$$W_1(m_{x_1}, m_{x_5}) = (\frac{1}{2} - \frac{1}{8}) \times d(x_1, x_5)$$

$$\begin{aligned} c(x_1, x_5) &= 1 - \frac{W_1(m_{x_1}, m_{x_5})}{d(x_1, x_5)} \\ &= 1 - \frac{(\frac{1}{2} - \frac{1}{8}) \times d(x_1, x_5)}{d(x_1, x_5)} \\ &= \frac{5}{8} \end{aligned}$$

F. R. Chung and S.-T. Yau, "Logarithmic harnack inequalities," Mathematical Research Letters, 3(6), 793–812, 1996.

Y. Ollivier, "Ricci curvature of metric spaces," Comptes Rendus Mathématique, 345 (11), 643–646, 2007.

J. Lott and C. Villani, "Ricci curvature for metric-measure spaces via optimal transport," Annals of Mathematics, 903–991, 2009

Y. Lin, L. Lu, and S.-T. Yau, "Ricci curvature of graphs," Tohoku Mathematical Journal, Second Series, 63(4), 605–627, 2011.

Y. Lin and S.-T. Yau, "Ricci curvature and eigenvalue estimate on locally finite graphs," Mathematical research letters, 17(2), 343–356, 2010.

C.-C. Ni, Y.-Y. Lin, J. Gao, X. D. Gu, and E. Saucan, "Ricci curvature of the internet topology," in 2015 IEEE Conference on Computer Communications (INFOCOM), 2758–2766, IEEE, 2015.

W. Zeng, D. Samaras, and X. D. Gu. "Ricci flow for 3D shape analysis." IEEE Transactions on Pattern Analysis and Machine Intelligence 32 (4) 662-677, 2010.

**k-simplex**

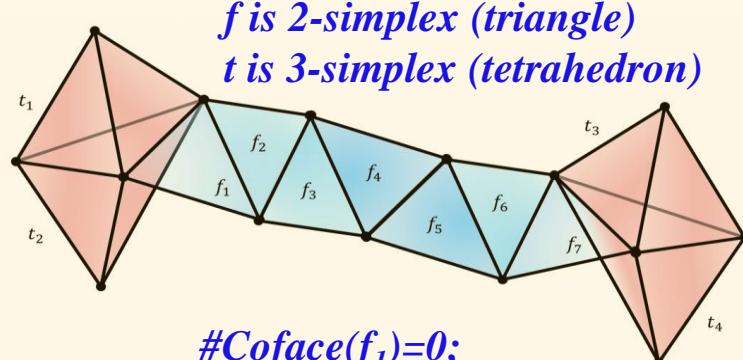
# Forman Ricci curvature (FRC)

$$\mathcal{F}_k^\sharp(\sigma^k) = \#\{\sigma^k \subset \sigma^{(k+1)}\} + \#\{\sigma^{(k-1)} \subset \sigma^k\} - \#\{\text{parallel neighbours of } \sigma^k\}$$

upper degree (number of “cofaces”)

lower degree (number of “faces”)

either upper adjacent or lower adjacent, but not both



$$\mathcal{F}_2^\sharp(f_1) = \mathcal{F}_2^\sharp(f_7) = -1.$$

$$\mathcal{F}_2^\sharp(f_i) = 1 \text{ for } i = 2, 3, \dots, 6.$$

$$\mathcal{F}_3^\sharp(t_i) = 3 \text{ for } i = 1, 2, 3, 4.$$

## ORC VS FRC

- (1) ORC is defined on graph; FRC is defined on cell complex.
- (2) ORC is for network probabilistic properties; FRC is for combinatorial properties.
- (3) ORC is highly correlated with FRC:

*Positive ORCs or FRCs are commonly found in densely-packed clusters or “communities”, while negative ORCs or FRCs usually represent bridges or links between two clusters.*

R. Forman, “Bochner’s method for cell complexes and combinatorial Ricci curvature,” Discrete and Computational Geometry, 29 (3), 323–374, 2003.

Desbrun, Mathieu, Eva Kanso, and Yiying Tong. “Discrete differential forms for computational modeling.” Discrete differential geometry. Birkhäuser Basel, 287-324, 2008.

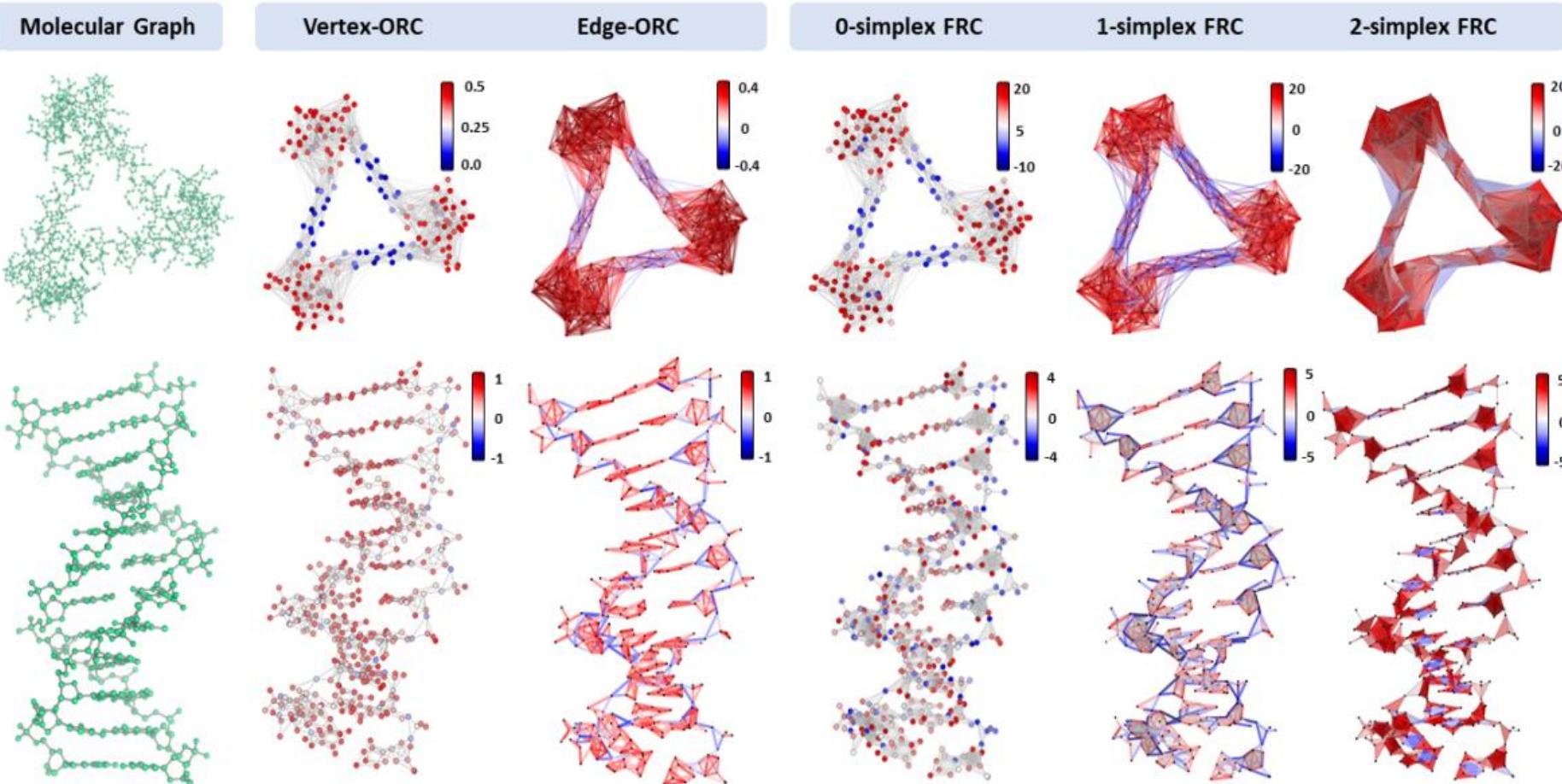
Jost, Jürgen, and Jürgen Jost. “Riemannian geometry and geometric analysis”. Vol. 42005. Berlin: Springer, 2008.

R. Sreejith, K. Mohanraj, J. Jost, E. Saucan, and A. Samal, “Forman curvature for complex networks,” Journal of Statistical Mechanics: Theory and Experiment, 2016(6), 063206, 2016.

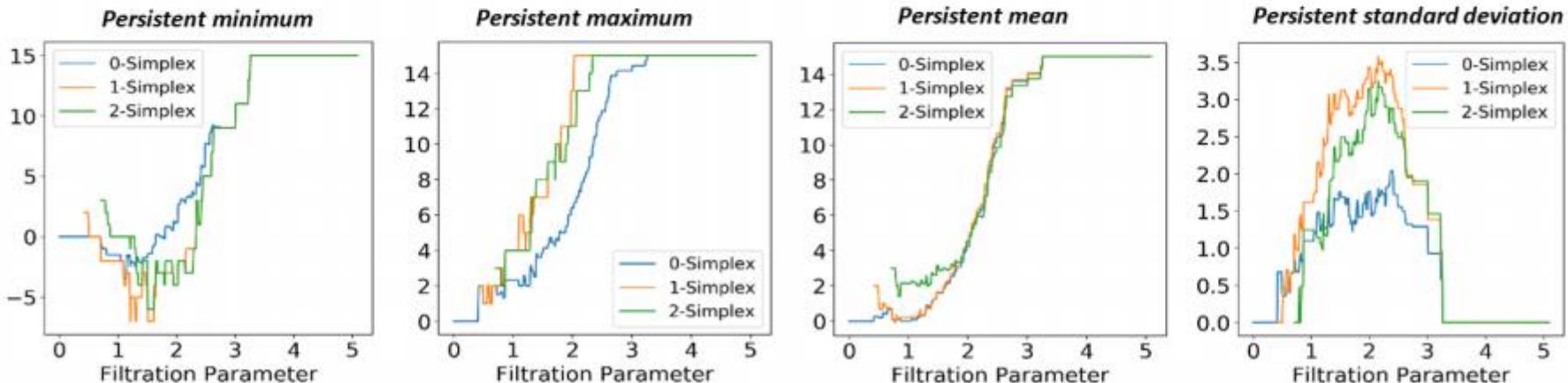
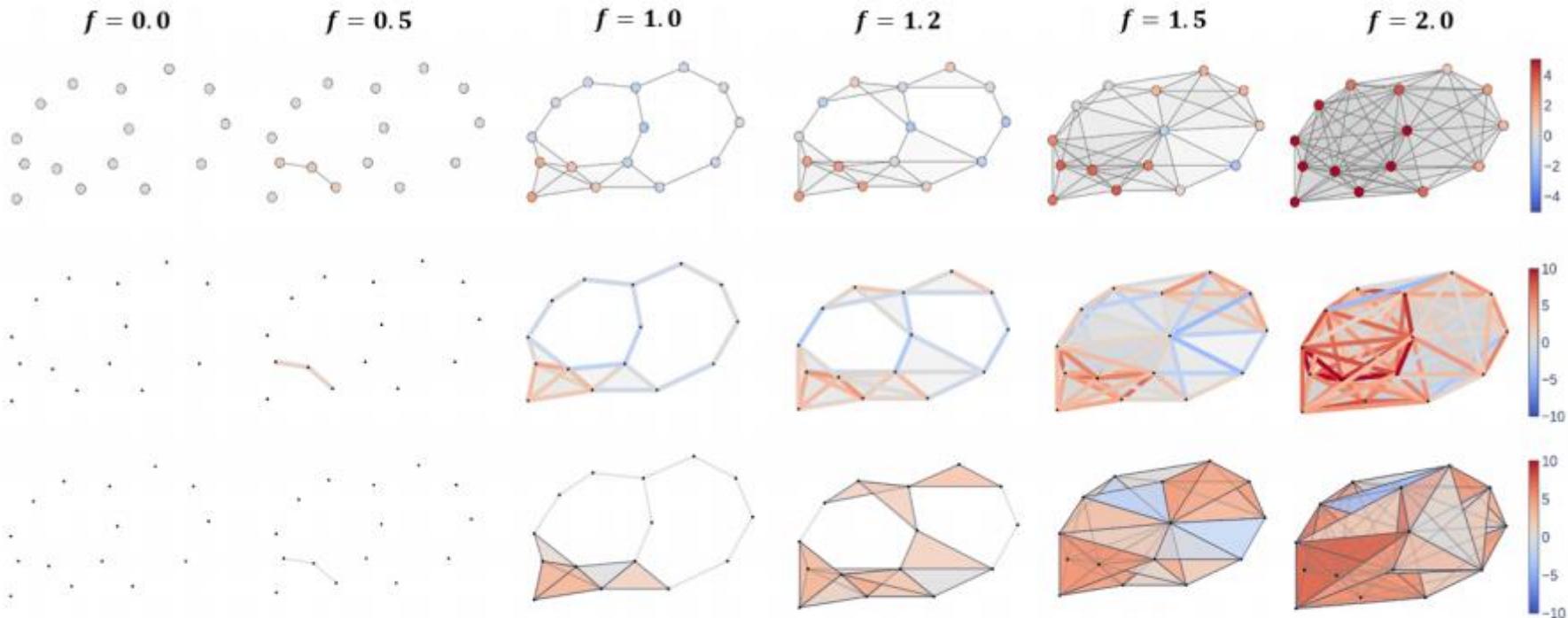
A. Samal, R. Sreejith, J. Gu, S. Liu, E. Saucan, and J. Jost. “Comparative analysis of two discretizations of Ricci curvature for complex networks”. Scientific reports, 8(1):1–16, 2018.

.....

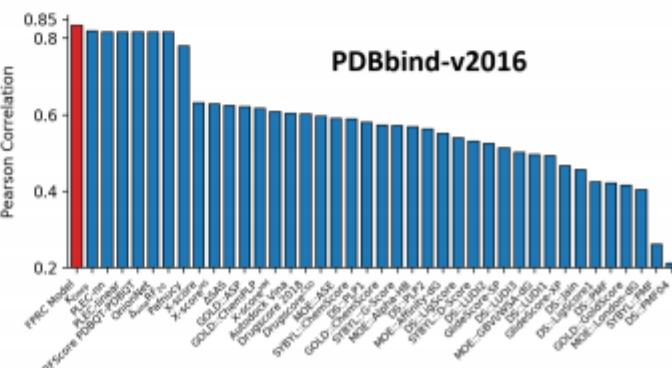
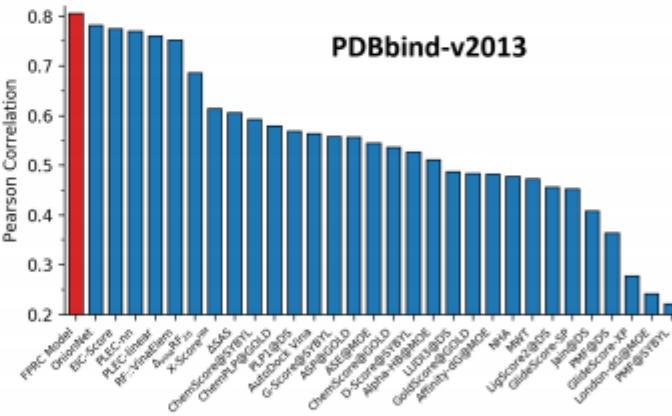
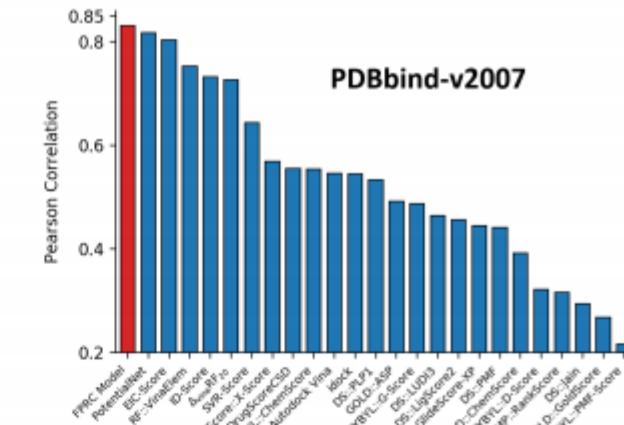
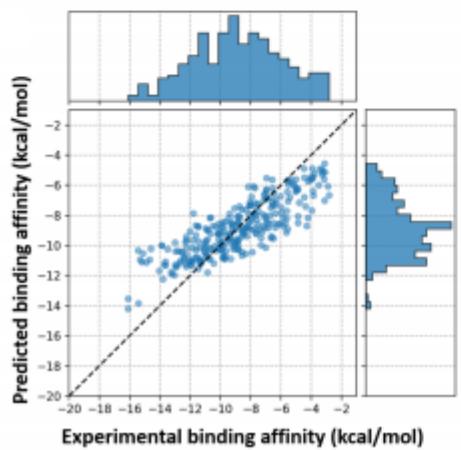
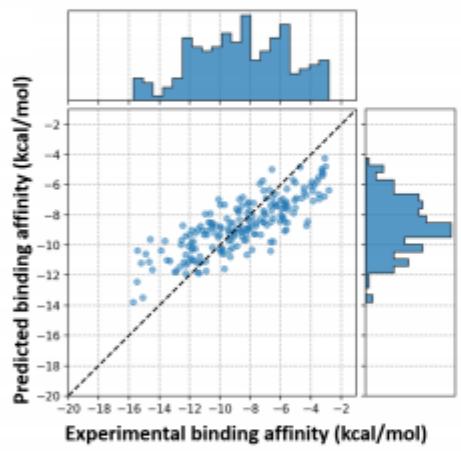
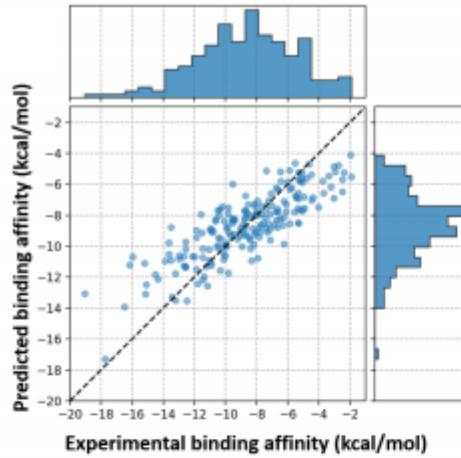
# ORC VS FRC



# (Forman) Persistent Ricci curvature



# Benchmark testing with PDBbind datasets



(Wee, Xia, submitted, 2020)

*Model setting:  
Persistent attributes  
+  
Random forest /  
gradient boosting tree*

# Persistent spectral models

Combinatorial Bochner-Weitzenböck decomposition

Hodge decomposition

## Persistent Ricci curvature

Gauss-Bonnet

## Persistent homology

### Combinatorial Bochner-Weitzenböck decomposition

$$L_k = \Delta_k + \text{Ricc}^{\mathcal{F}}_k$$

#### Combinatorial Bochner Laplacian

$$\Delta_k(i,j) = \begin{cases} \sum_{m \neq i} |L_k(i,m)|, & i = j \\ L_k(i,j), & i \neq j \end{cases}$$

#### Forman Ricci matrix

$$\text{Ricc}_k^{\mathcal{F}}(i,j) = \begin{cases} \mathcal{F}_k^{\sharp}(\sigma_i^k), & i = j \\ 0, & i \neq j \end{cases}$$

#### Forman Ricci curvature

$$\mathcal{F}_k^{\sharp}(\sigma_i^k) = \underbrace{d(\sigma_i^k) + k + 1}_{L_k(i,i)} - \underbrace{\sum_{m \neq i} |(B_{k+1}B_{k+1}^T)(i,m) + (B_k^TB_k)(i,m)|}_{\Delta_k(i,i)}$$

### Hodge decomposition

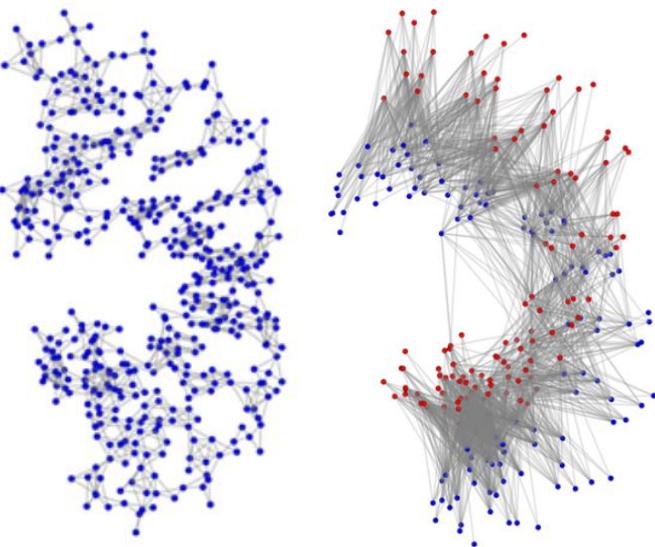
harmonic flow (*topological invariant*)

$$C^1(K) = \text{Im}(\delta_0) \oplus \ker(\Delta_1) \oplus \text{Im}(\delta_1^*)$$

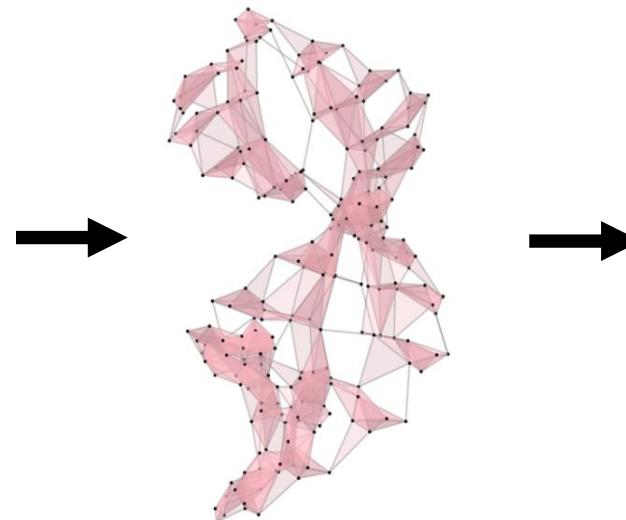
gradient flow

curl flow

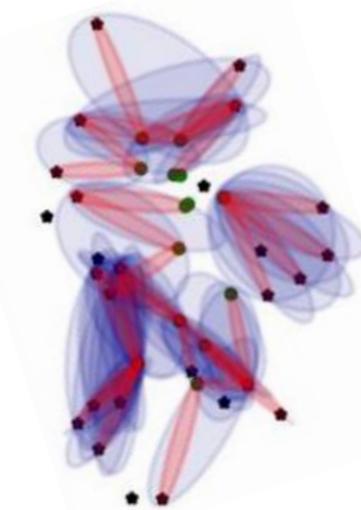
## Graph



## Simplicial complex



## Hypergraph

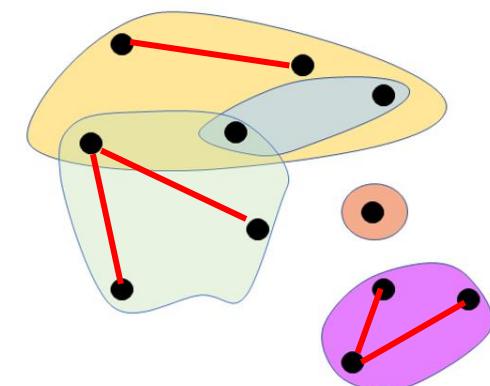
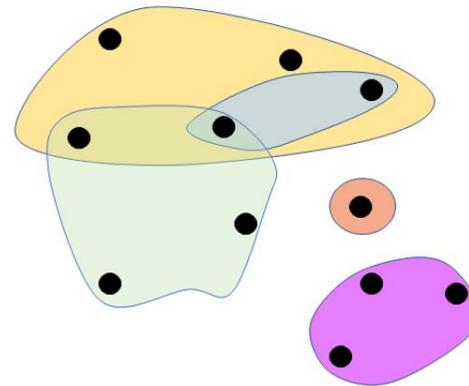


**Collaborator**  
**Jie Wu**  
**Math, HEBNU**



Ren, Shiquan, et al. "Computing the Homology of Hypergraphs." *arXiv preprint arXiv:1705.00151* (2017).  
Ren, Shiquan, Chengyuan Wu, and Jie Wu. "Operators on random hypergraphs and random simplicial complexes." *arXiv preprint arXiv:1712.02045* (2017).  
Ren, Shiquan, and Jie Wu. "Stability of persistent homology for hypergraphs." *arXiv preprint arXiv:2002.02237* (2020).  
Ren, Shiquan, et al. "A Discrete Morse Theory for Hypergraphs." *arXiv preprint arXiv:1804.07132* (2018).

## Hypergraph VS super hypergraph

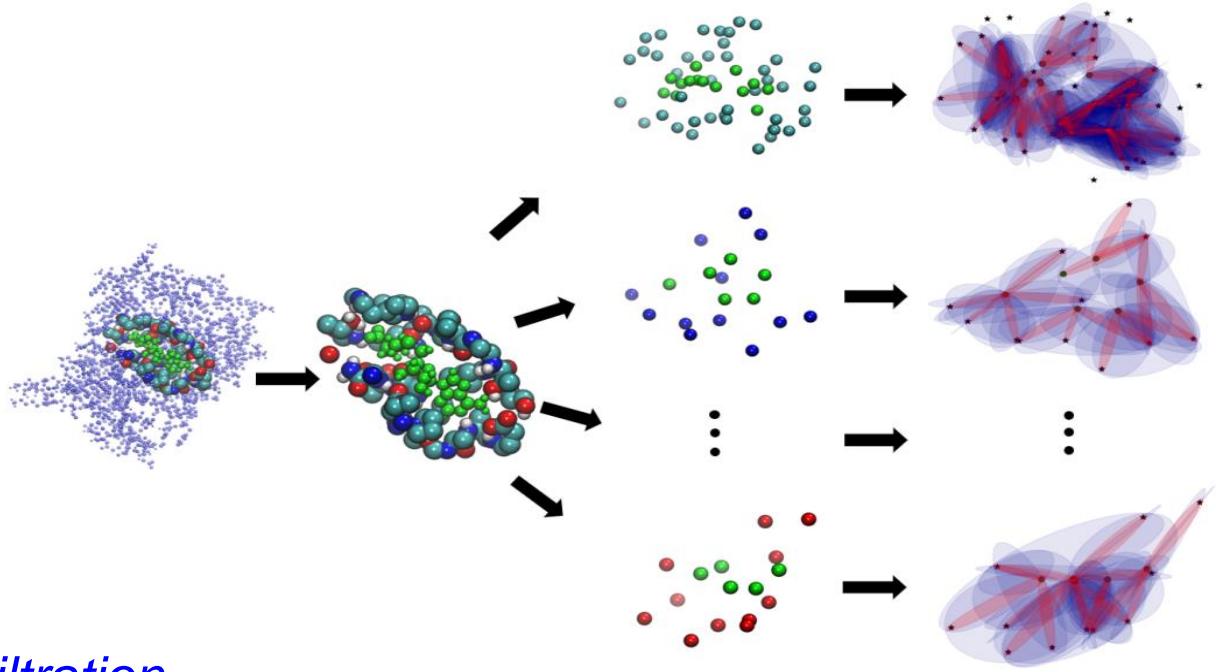


**hyperedge is a  
set of vertices**

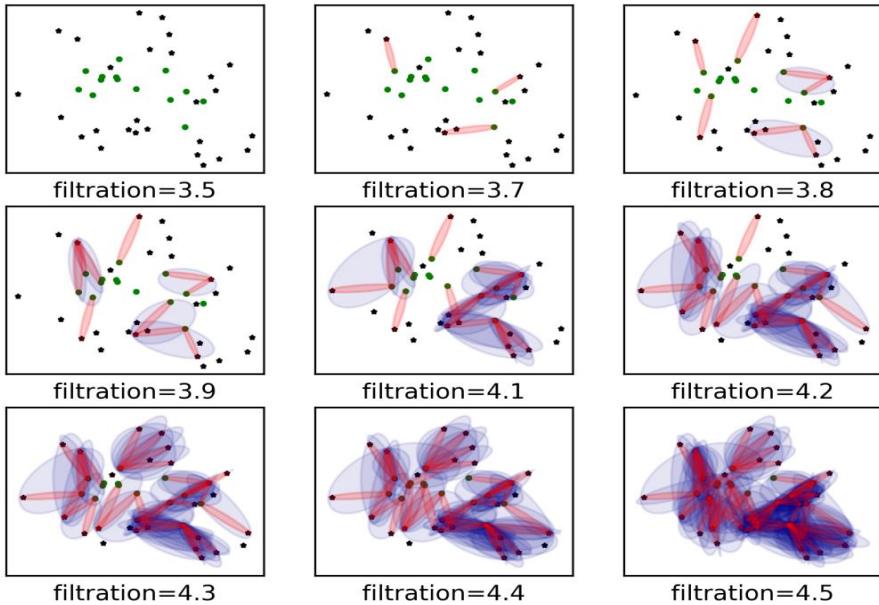
**super hyperedge  
is a subgraph**

# Protein-ligand interaction modeled as hypergraph

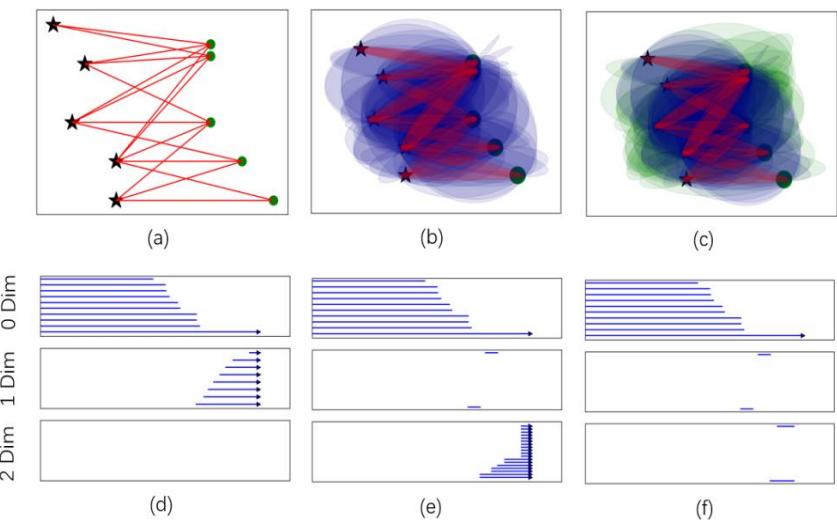
## Hypergraph-based models



## Hypergraph-based filtration



## Bipartite graph VS Hypergraph



# Hypergraph for protein-ligand interactions

Hypergraph:  $(V_{\mathcal{H}}, \mathcal{H})$

set of vertices

set of protein atoms

set of hyperedges

$$V_{\mathcal{H}} = V_P \cup V_L$$

set of ligand atoms

Hyperedge:

At least one atom from protein and another atom from ligand

n-hyperedge

$$\sigma_n = \begin{cases} \{\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_n\}; \mathbf{v}_k \in V_{\mathcal{H}} (0 \leq k \leq n), \exists i, j \in [0, n], \mathbf{v}_i \in V_P, \mathbf{v}_j \in V_L, & n > 0 \\ \{\mathbf{v}_0\}; \mathbf{v}_0 \in V_{\mathcal{H}}, & n = 0. \end{cases}$$

Filtration value for hyperedge:

$$f(\sigma_n) = \begin{cases} \max_{0 \leq i < j \leq n} d(\mathbf{v}_i, \mathbf{v}_j), & n > 0 \\ 0, & n = 0. \end{cases}$$

Interactive distance:

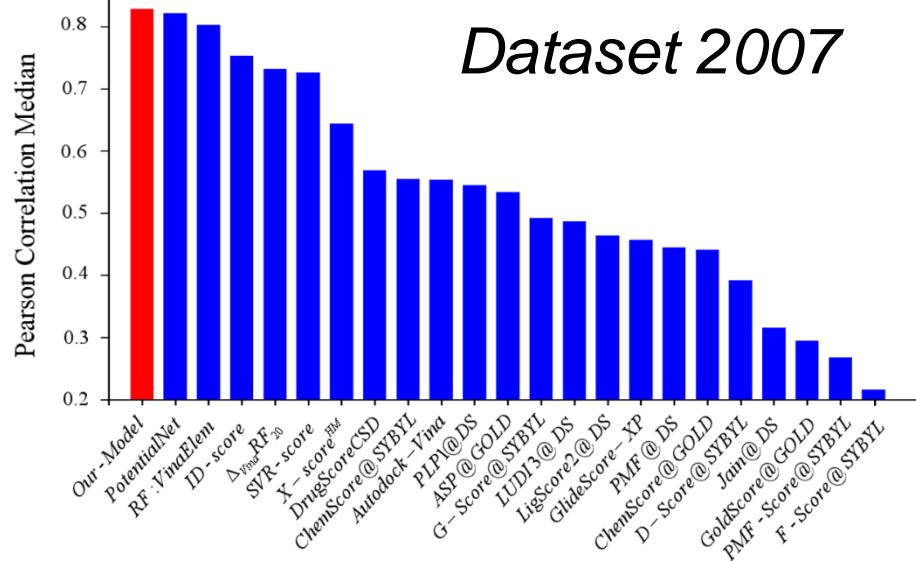
$$d(\mathbf{v}_i, \mathbf{v}_j) = \begin{cases} \|\mathbf{v}_i - \mathbf{v}_j\|, & \text{if } \mathbf{v}_i \in V_P, \mathbf{v}_j \in V_L \text{ or } \mathbf{v}_i \in V_L, \mathbf{v}_j \in V_P \\ g(\mathbf{v}_i, \mathbf{v}_j), & \text{otherwise.} \end{cases}$$

Distance between two atoms from the same molecule is defined as the maximal distance of the two atoms to the other molecule

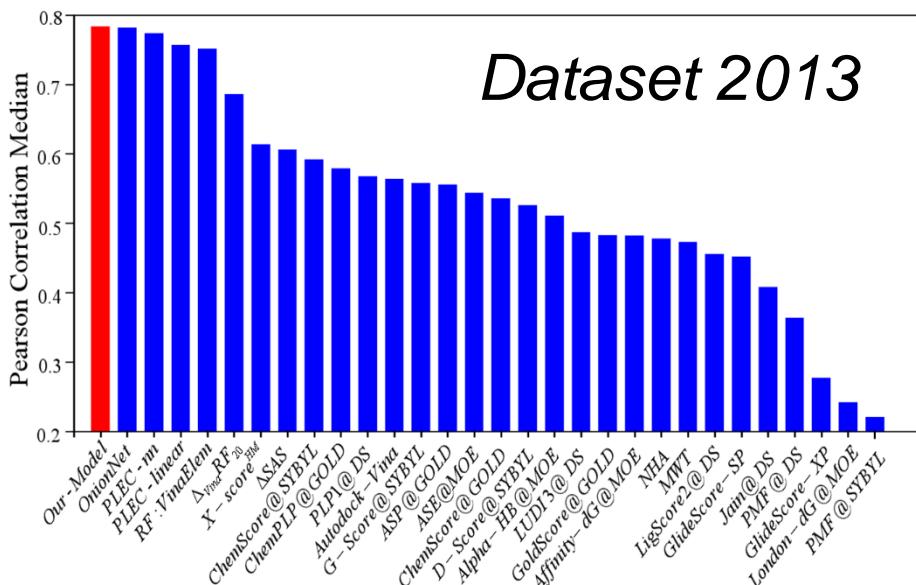
$$g(\mathbf{v}_i, \mathbf{v}_j) = \begin{cases} \max_{\mathbf{v}_k \in V_P, \{\mathbf{v}_i, \mathbf{v}_j, \mathbf{v}_k\} \in \mathcal{H}} \{\|\mathbf{v}_i, \mathbf{v}_k\|, \|\mathbf{v}_j, \mathbf{v}_k\|\} + d_0, & \text{if } \mathbf{v}_i, \mathbf{v}_j \in V_L \\ \max_{\mathbf{v}_k \in V_L, \{\mathbf{v}_i, \mathbf{v}_j, \mathbf{v}_k\} \in \mathcal{H}} \{\|\mathbf{v}_i, \mathbf{v}_k\|, \|\mathbf{v}_j, \mathbf{v}_k\|\} + d_0, & \text{if } \mathbf{v}_i, \mathbf{v}_j \in V_P \end{cases}$$

# Benchmark testing with PDBbind datasets

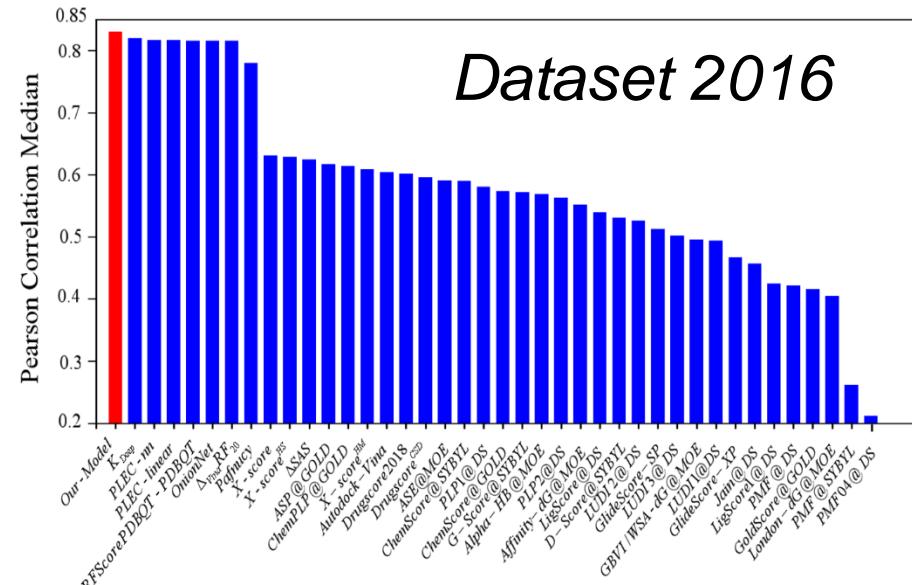
*Model setting:  
homology vectors*  
+  
*Random forest*



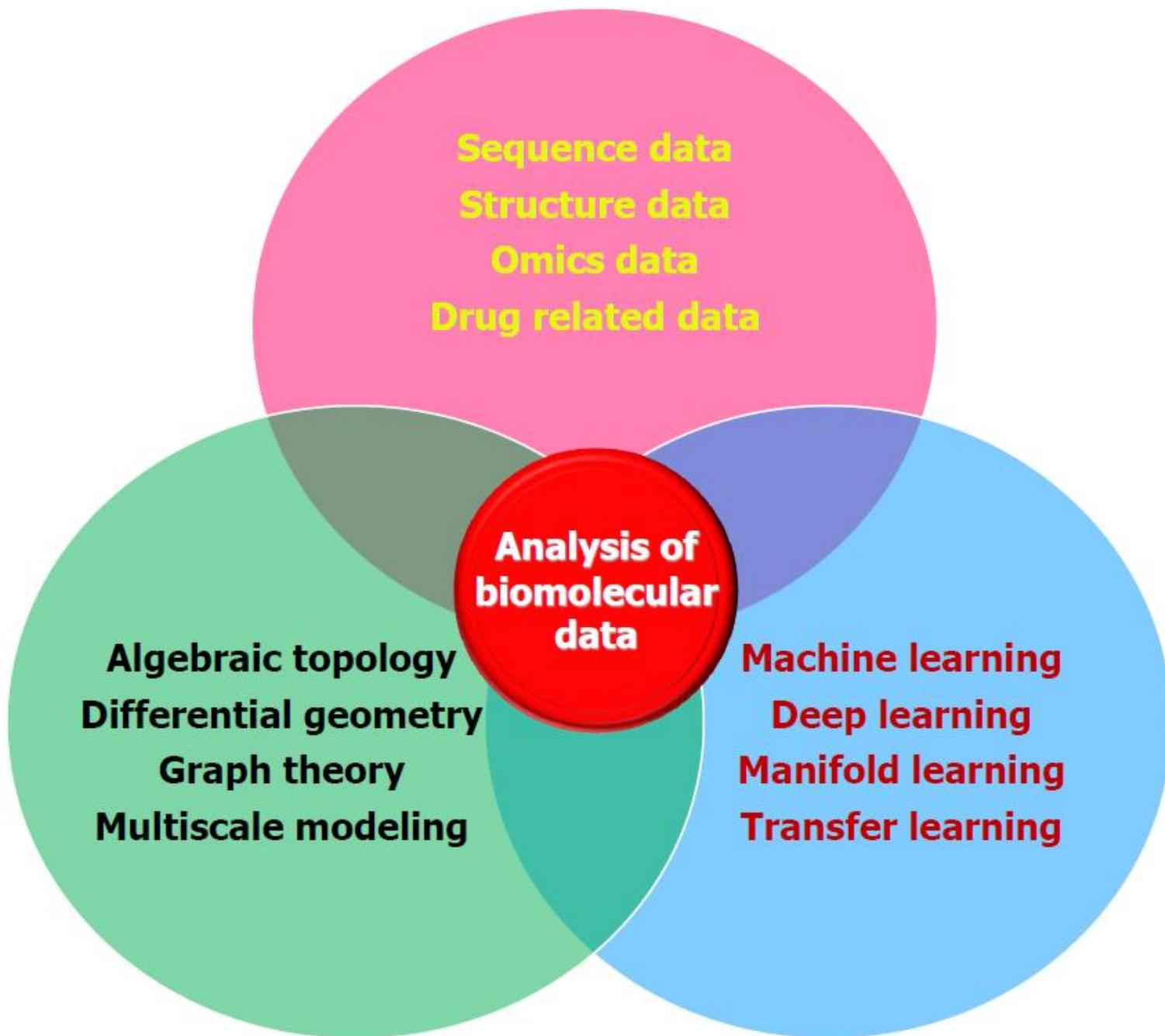
*Dataset 2007*



**(Liu, Wu, Wang, Xia, briefings in bioinformatics, accepted, 2020)**



*Dataset 2016*



# Group members

*Bin Liu*



*Xiang Liu*



*Zhenyu Meng*



*Joon Wei KOH*



*JunJie WEE*



*Vijai Anand*

**WE'RE  
HIRING!**

***Two postdocs (2021)***

***Geometric/topological modelling,  
Machine learning, Mathematical  
biology, Bioinformatics...***

**阿里巴巴达摩院AI医疗团队实  
习生 (机器学习+拓扑几何)**

# Grant support

NTU-JSPS (2019-2022)

Alibaba-NTU (2020-2021)

Merlion (2020-2022)

MOE-Tier 1 (2018-2021, 2019-2022)

MOE-Tier 2 (2018-2021, 2021-2024)

*Email:  
xiakelin@ntu.edu.sg*



## Conference Program

# The Third Conference on Computational and Mathematical Bioinformatics and Biophysics

Tsinghua Sanya International Mathematics Forum

December 20 – 24, 2020, (Beijing Time)

December 19 – 23, 2020, (US Central Time)

Zoom Meeting ID: 981 2055 6545

Password: CM<sup>B</sup>B

Zoom Link:

<https://uasystem.zoom.us/j/98120556545?pwd=OXjhZENlTDhp aUsrNDQ5a2NvQVNudz09>

Organizing Committee:

Stephen Shing-Toung Yau, Tsinghua University, China

Guowei Wei, Michigan State University, USA

Changchuan Yin, University of Illinois at Chicago, USA

Shan Zhao, University of Alabama, USA