# Can AI discover the drugs of the future?

**Guo-Wei Wei**

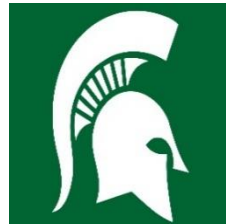**Mathematics**

**Michigan State University**

**http://www.math.msu.edu/~wei**

**The 8th International Congress of Chinese Mathematicians of 2019**

**Beijing, June 9 - 15, 2019**

# The Biggest Crisis of the Contemporary Science

The number of researchers in the world who know both graduate-level mathematics and molecular-level biology is smaller than the number of fields medalists!
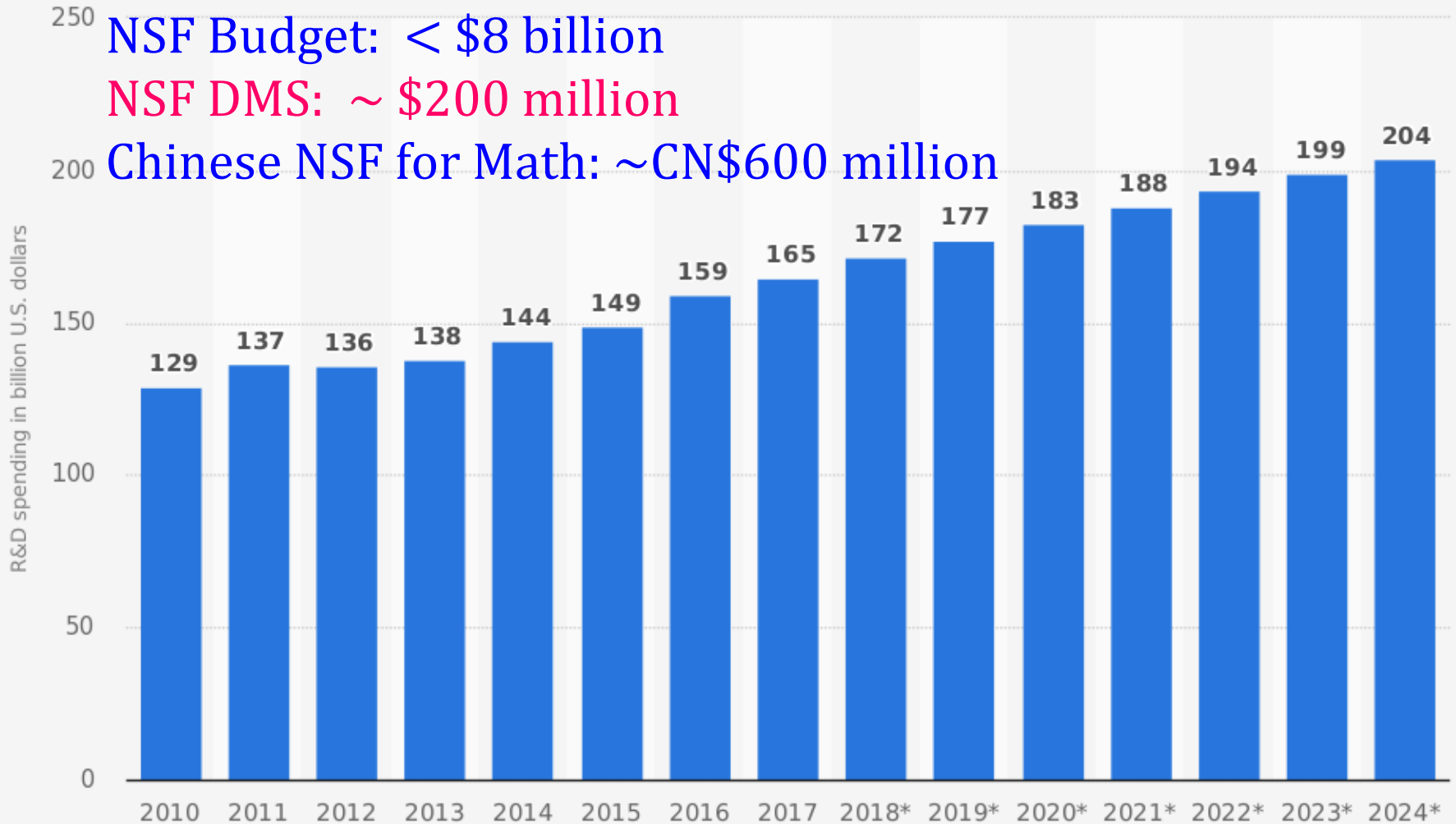
The rule of life has been there for billions of years but very little is known about it!

None knows the existence and uniqueness of mathematical foundation for life!

# A Brief Summary of Modern Biological Science

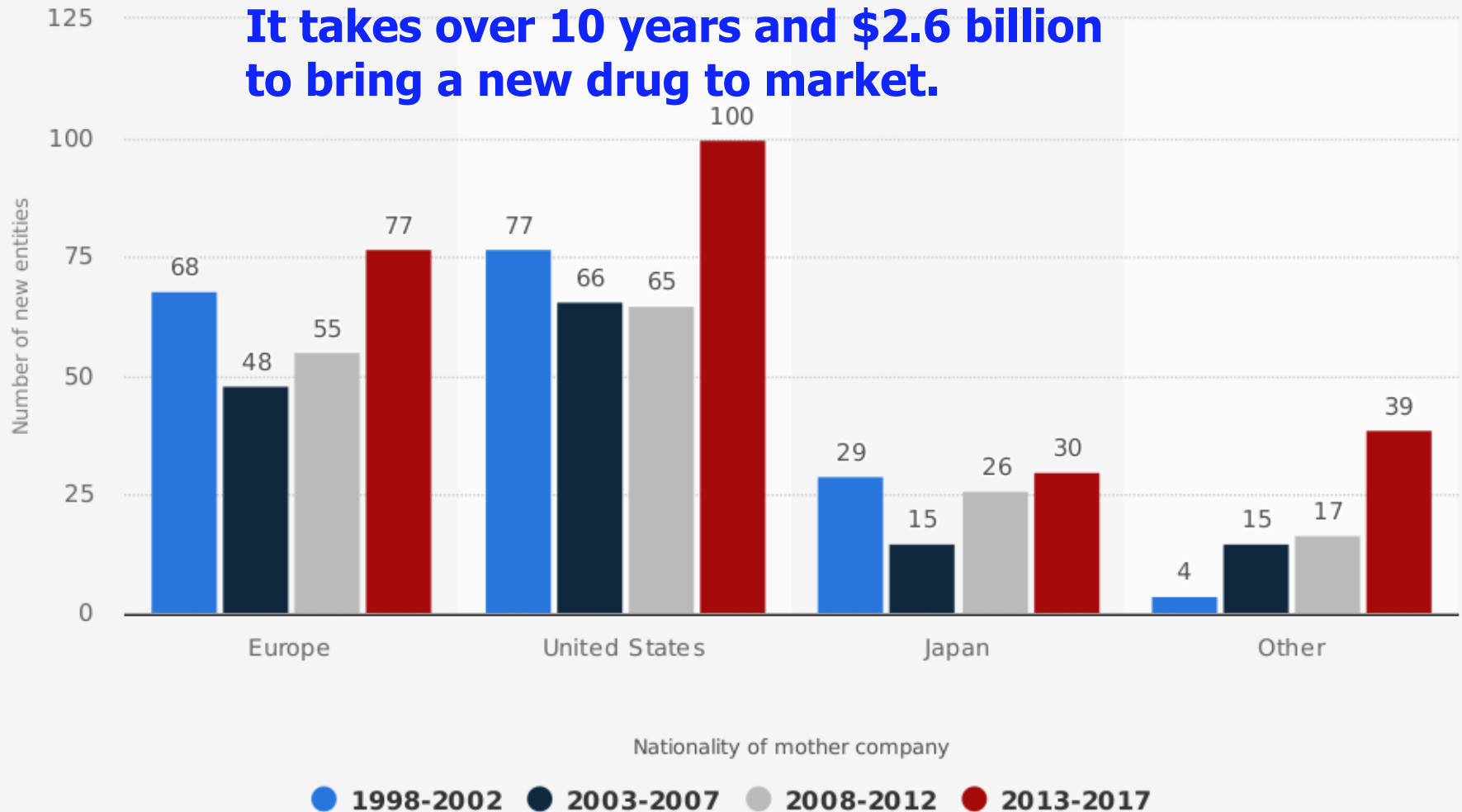| 1960 | 2000 | 2019 |
|---|---|---|
| Organismal biology (i.e., nonliving organisms, living organisms, developmental biology, morphology, anatomy, physiology, and medicine) | Molecular organismal biology, organomics, connectomics, foodomics, physiomics, pharmacogenomics, ... | |
| Ecology | Molecular ecology | |
| Evolution (i.e., life, and evolutionary biology) | Molecular evolution | |
| Molecular and cellular biology (i.e., cell biology, biochemistry, molecular biology, and genetics) | Omics (e.g., genomics, proteomics, metabolomics, metagenomics, lipidomics, glycomics, transcriptomics, epigenomics, ...) | |

Macroscopic    Mesoscopic    Microscopic

# Total global spending on pharmaceutical research and development from 2010 to 2024 (in billion U.S. dollars)



NSF Budget: $< \$8$ billion

NSF DMS: $\sim \$200$ million

Chinese NSF for Math: $\sim$CN$600 million

Funding for mathematical research is $\sim \varepsilon^3$

# Number of new chemical or biological entities developed between 1992 and 2017, by region of origin

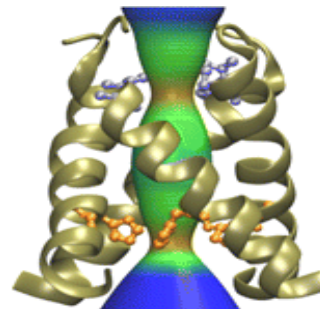**It takes over 10 years and $2.6 billion to bring a new drug to market.**



Number of new entities (y-axis), Nationality of mother company (x-axis)

| Region | 1998-2002 | 2003-2007 | 2008-2012 | 2013-2017 |
|---|---|---|---|---|
| Europe | 68 | 48 | 55 | 77 |
| United States | 77 | 66 | 65 | 100 |
| Japan | 29 | 15 | 26 | 30 |
| Other | 4 | 15 | 17 | 39 |

Source
EFPIA
© Statista 2018

Additional Information:
Worldwide

# 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**

NH₂

# Biological data



National Library of Medicine Twenty Four Years of Growth: NCBI Data and User Services

**GenBank**

Legend: GenBank Base Pairs; Users (Average)

Labels: BLAST, Entrez, Genomes, GenBank at NCBI, OMIM, PubMed, Human Genome, PubMed Central, NIH Public Access, PubChem, Genome-Wide Association Studies, Genome Reference Consortium, 1000 Genomes, ClinVar, Genetic Testing Registry, MedGen PubReader

**PDB**

# Artificial Intelligence & Deep learning

Bryson and Ho (Backpropagation 1969); Fukushima (Neo-Cognitron 1980); LeCun (CNN 1998); Hopfield (RNN 1982); Hochreiter and Schmidhuber (LSTM 1997); Goodfellow et al (GAN 2014); Autoencoder; Image translation, ...

won 25 of 43 contests and was ranked 1$^{St}$ among 98 competitors in CASP 13.

Google DeepMind

Protein Sequence

S Q E T R K K C T E M K K K F K N C E V R C D E S N H C V E V R C S D T K Y T L C

Neural Network

Databases

Distance Predictions

Angle Predictions

Score

Gradient Descent

Structure

T0954 / 6CVZ    T0965 / 6D2V    T0955 / 5W9F

Ground truth

Average predicted distance

# How to do deep learning for 3D biomolecular data?

**Obstacles for deep learning of 3D biomolecules:**

- Geometric dimensionality: $\mathbb{R}^{3N}$ , where $N$ ~5000 for a protein.
- Machine learning dimensionality:  > $1024^3 m$, where $m$ is the number of atom types in a protein.
- Molecules have different sizes --- non-scalable.
- Complexity: intermolecular & intramolecular interactions

**Solution:**

-  Geometric simplification, dimension reduction & scale unification

# Two schools of thinking

Given a protein with **N** atom and an average of **n** electrons in each atom

**Fundamentalism**

**Quantum Mechanics** $\mathbb{R}^{3Nn+3N}$

**QM/MM** $\mathbb{R}^K$
$3N < K < 3N(n+1)$

**Molecular Mechanics** $\mathbb{R}^{3N}$

**Multiscale Coarse-grain** $\mathbb{R}^M$ ($3 < M < 3N$)

**Poisson-Boltzmann, PNP, etc.** $\mathbb{R}^3$

**Differentiable Manifold** $\mathbb{R}^2$

**Algebraic Topology** $\mathbb{R}^1$

**Graph Theory** $\mathbb{R}^0$

**Geo-Top Indices** $\mathbb{R}^0$

**Reductionism**

**Basic hypothesis:**
**Intrinsic physics lies on low-dimensional manifolds in a high dimensional space**

# Classical Topology

## Möbius Strips (1858)



## Klein Bottle (1882)





**Leonhard Paul Euler**
(Swiss Mathematician,

April 15, 1707 – Sept 18 1783)

## Torus

## Double Torus



### Seven Bridges of Konigsberg





2016 NOBEL PRIZE IN PHYSICS

"For the greatest benefit to mankind"

The Royal Swedish Academy of Sciences has decided to award the

David J. Thouless
F. Duncan M. Haldane
J. Michael Kosterlitz

"for theoretical discoveries of topological phase transitions and topological phases of matter"

Nobelprize.org

Augustin-Louis Cauchy,
Ludwig Schläfli,
Johann Benedict Listing,
Bernhard Riemann, and
Enrico Betti

Leonhard Euler (1735)

# Topological invariants: Betti numbers

$\beta_0$ **is the number of connected components.**
$\beta_1$ **is the number of tunnels or circles.**
$\beta_2$ **is the number of cavities or voids.**

| **Point** | **Circle** | **Sphere** | **Torus** |
|:---:|:---:|:---:|:---:|
|  |  |  |  |
| $\beta_0 = 1$ | $\beta_0 = 1$ | $\beta_0 = 1$ | $\beta_0 = 1$ |
| $\beta_1 = 0$ | $\beta_1 = 1$ | $\beta_1 = 0$ | $\beta_1 = 2$ |
| $\beta_2 = 0$ | $\beta_2 = 0$ | $\beta_2 = 1$ | $\beta_2 = 1$ |

# Vietoris-Rips complexes of planar point sets

**Simplexes:**



*0*-simplex         *1*-simplex         *2*-simplex                 *3*-simplex
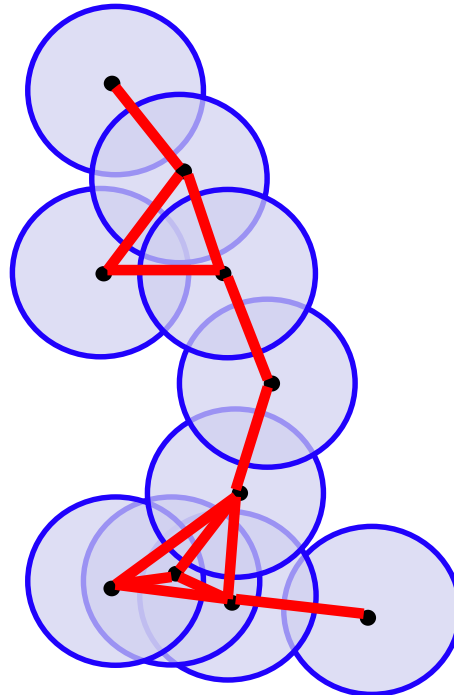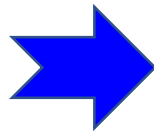
**Simplicial complexes of ten points:**

# Persistent homology

**Simplexes:**



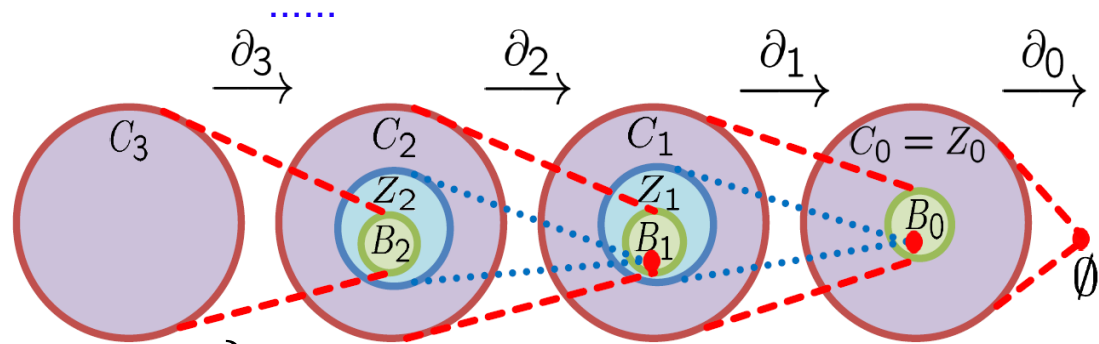*0*-simplex   *1*-simplex   *2*-simplex   *3*-simplex

Frosini and Nandi (1999), Robins (1999), Edelsbrunner, Letscher and Zomorodian (2002),  Zomorodian and Carlsson (2005),  Edelsbrunner and Harer, (2007) Kaczynski, Mischaikow and Mrozek (2004), Ghrist (2008),

......

**k-chain:** $\quad K = \left\{ \sum_j c_j \sigma_j^k \right\}$

**Chain group:** $\quad C_k(K, \mathbb{Z}_2)$

**Boundary operator:**

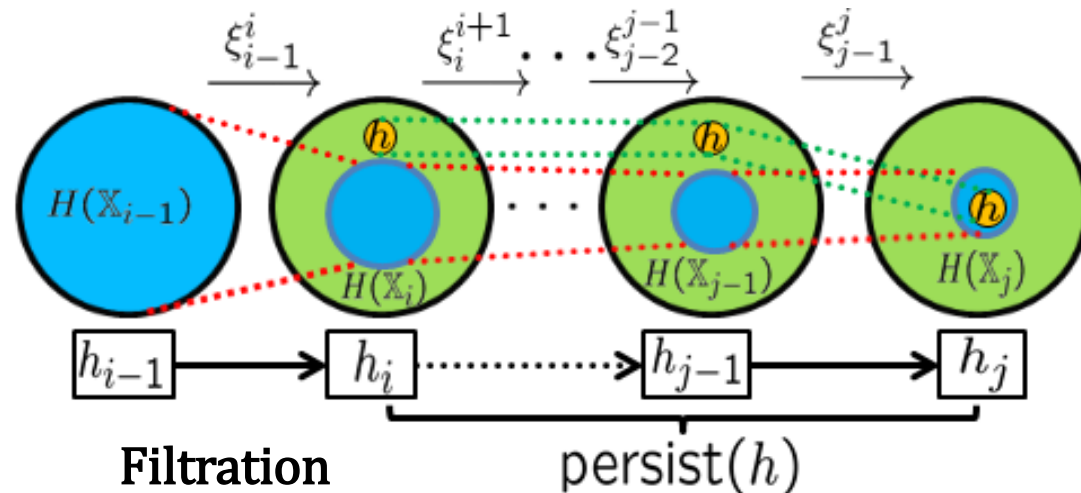$$\partial_k \sigma^k = \sum_{j=0}^{k} (-1)^j \left\{ v_0, v_1, \ldots, \widehat{v_j}, \ldots, v_k \right\}$$



**Cycle group:** $\quad Z_k = \mathrm{Ker}\, \partial_k$

**Boundary group:** $\quad B_k = \mathrm{Im}\, \partial_{k+1}$

**Homology group:** $\quad H_k = \dfrac{Z_k}{B_k}$

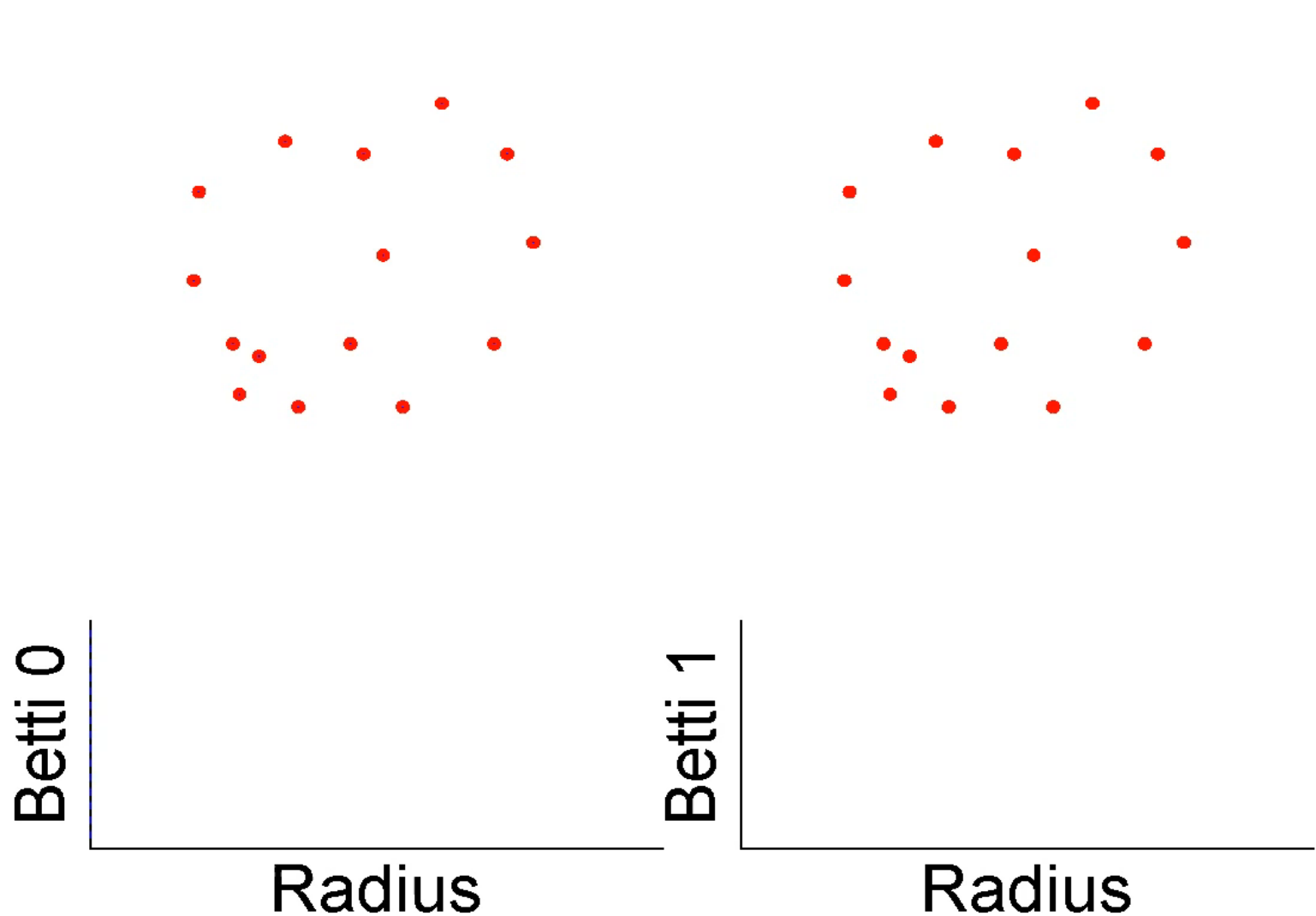**Betti number:** $\quad \beta_k = \mathrm{Rank}(H_k)$



Filtration     persist(h)

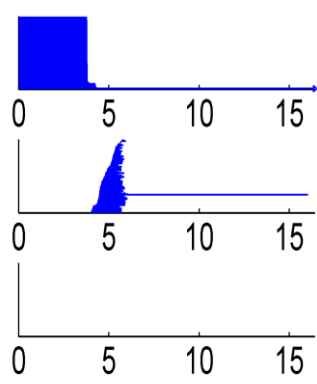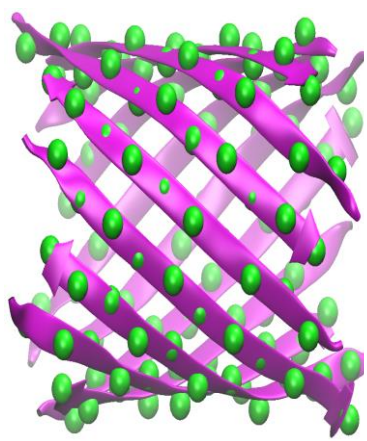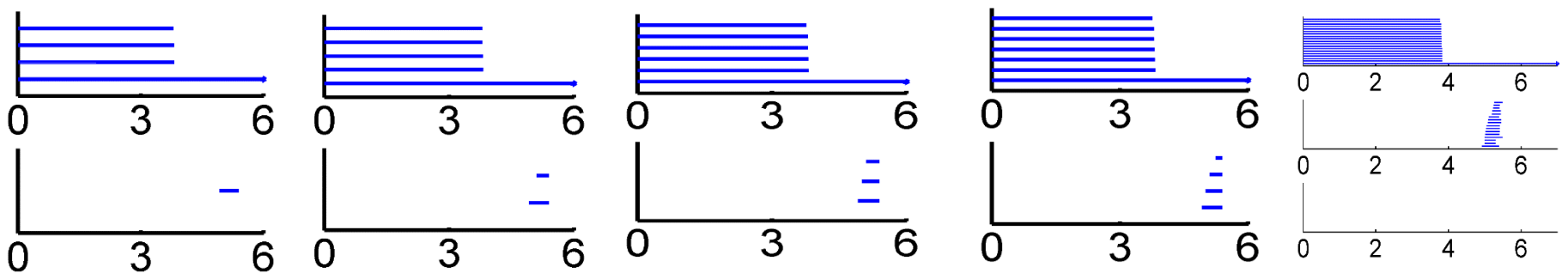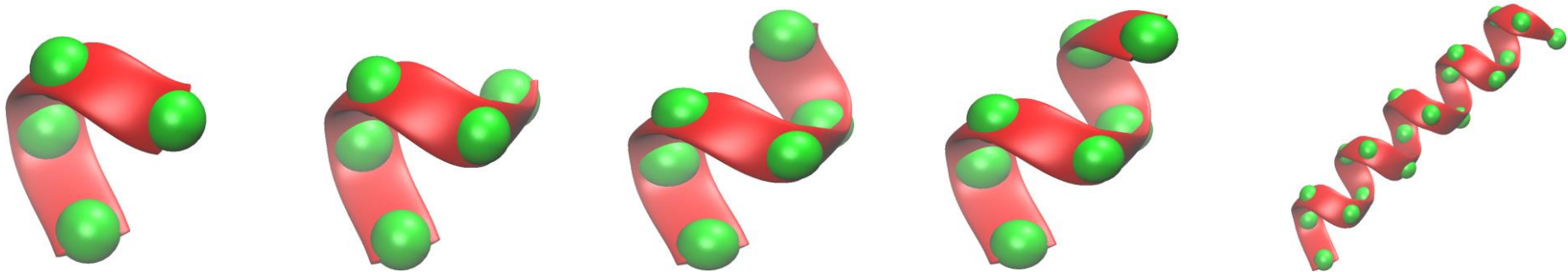Xia, Wei, IJNMBE, 2014;
Xia, Feng, Tong, Wei, JCC, 2015

# Algebraic Topology

Vietoris-Rips complexes, persistent homology and topological fingerprint

(Xia, Wei, 2014)

# Topological fingerprints of an alpha helix
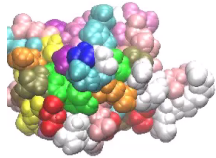
# Algebraic Topology
## 2D persistent homology of protein 1UBQ unfolding



$\beta_0$

$\beta_1$

$\beta_2$

Radius

Time

Kelin Xia

# Topological convolutional deep Learning architecture



Convolution (128x200)

Pooling (128x100)

Flattening (1xN)

Prediction

Multichannel images (54x200)

Original protein-ligand Complex

Classify atoms into element specific groups

Generate topological, DG & graph fingerprints

Convolutional deep learning neural network

(Cang & Wei, PLOS CB, 2017)

**Helicoid**

Viral morphology

**Leonhard P. Euler**
(Swiss Mathematician, April 15, 1707 – Sept 18 1783

Joseph L. Lagrange (Italian Mathematician, January 25 1736 – April 10, 1813)

**Minimal Surfaces**
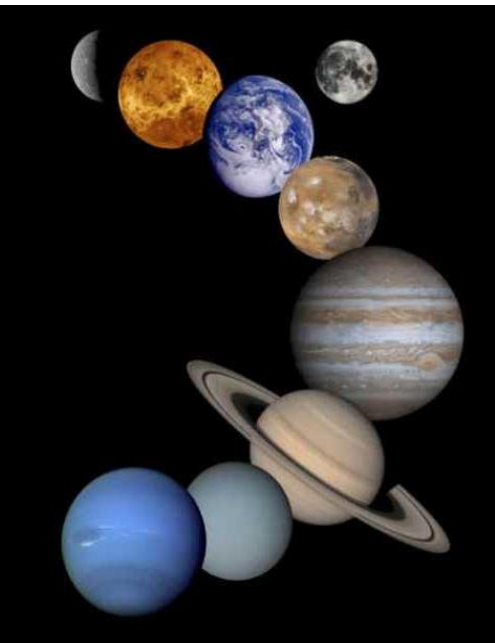**A way to minimize energy and maximize stability**

Man-made life, Mycoplasma mycoides
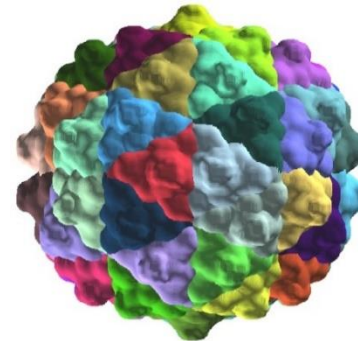
# Differential geometry based minimal surface model

$$G = \int \gamma[\text{area}]\, d\boldsymbol{r} \qquad \text{area} = |\nabla S|$$

where $G$ is the surface energy, gamma $(\gamma)$ is the surface tension, and $S$ is a surface characteristic function:

Generalized Laplace-Beltrami flow:

$$\frac{\partial S}{\partial t} = |\nabla S|\left[\nabla \cdot \frac{\gamma \nabla S}{|\nabla S|}\right]$$

Mean curvature



S=1

S=0





Shan Zhao

(Bates, Wei, Zhao, 2006; JCC,2008; Zhao, Cang, Tong & Wei, Bioinformatics 2018 )

# Differential Geometry


Pit, Valley, Saddle Valley, Flat, Minimal Surface, Saddle Ridge, Ridge, Peak


Kelin Xia

**Gauss**



**Mean**



## Mean curvatures of subcellular structures



(Feng, Xia, Tong and Wei, JCP, IJNMBI, 2012)

**Minimum**



**Maximum**



**CRISPR**



NaturalNews.com

# Protein binding site prediction by the product of curvature and electrostatics

# de Rham-Hodge theory and discrete exterior calculus

**Hodge decomposition:**

(Zhao, Wang, Tong & Wei, 2018)



**A vector field = Harmonic + curl-free + divergent-free**

**Cryo-EM data:**



Input          Normal Gradient          Tangential Curl          Tangential Harmonic          Central Harmonic
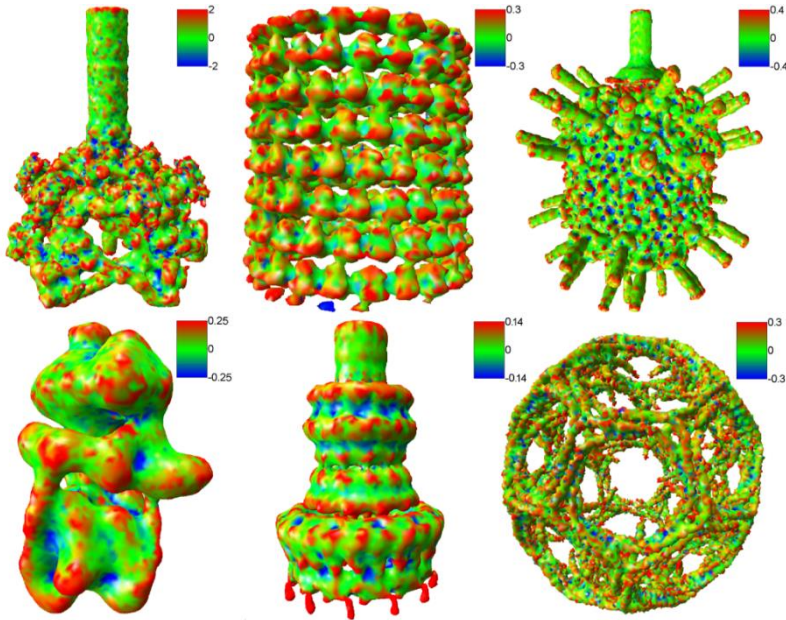


0-Form $\xrightarrow{\;d_0\;}$ 1-Form $\xrightarrow{\;d_1\;}$ 2-Form $\xrightarrow{\;d_2\;}$ 3-Form

$\star_0$  $\star_0^{-1}$       $\star_1$  $\star_1^{-1}$       $\star_2$  $\star_2^{-1}$       $\star_3$  $\star_3^{-1}$

Dual 3-Form $\xleftarrow{\;d_0^T\;}$ Dual 2-Form $\xleftarrow{\;d_1^T\;}$ Dual 1-Form $\xleftarrow{\;d_2^T\;}$ Dual 0-Form

point-based scalar field $\xrightarrow{\;\nabla\;}$ edge-based vector field $\xrightarrow{\;\nabla\times\;}$ face-based vector field $\xrightarrow{\;\nabla\cdot\;}$ cell-based scalar field

$\updownarrow \star$          $\updownarrow \star$          $\updownarrow \star$          $\updownarrow \star$

cell-based scalar field $\xleftarrow{\;\nabla\cdot\;}$ face-based vector field $\xleftarrow{\;\nabla\times\;}$ edge-based vector field $\xleftarrow{\;\nabla\;}$ point-based scalar field

# Persistent cohomology



Wasserstein curves

Zixuan Cang

(Cang & Wei, 2018)

# Algebraic Graph Theory for Biomolecules

## Molecular graph G(V,E)



## Adjacency matrix of $G(V_{ON},E)$

$$\begin{pmatrix} 0 & \Phi_{12} & \Phi_{13} & 0 \\ \Phi_{12} & 0 & 0 & \Phi_{24} \\ \Phi_{13} & 0 & 0 & \Phi_{34} \\ 0 & \Phi_{24} & \Phi_{34} & 0 \end{pmatrix}$$

Eigenvalues: $\lambda_1^A, \lambda_2^A, \ldots$

## Laplacian matrix of $G(V_{ON},E)$

$$\begin{pmatrix} \Phi_{12} + \Phi_{13} & -\Phi_{12} & -\Phi_{13} & 0 \\ -\Phi_{12} & \Phi_{12} + \Phi_{24} & 0 & -\Phi_{24} \\ -\Phi_{13} & 0 & \Phi_{13} + \Phi_{34} & -\Phi_{34} \\ 0 & -\Phi_{24} & -\Phi_{34} & \Phi_{24} + \Phi_{34} \end{pmatrix}$$

Eigenvalues: $\lambda_1^L, \lambda_2^L, \ldots$

## Can one hear the shape of a drum?

(Nguyen and Wei, 2018)

# Algebraic graph theory for biomolecules

**Protein**

**Hypergraph representation**

**Laplacian matrices and adjacency matrices**



Eigenvalue multiplicities in Laplacian and adjacency matrices are associated with structural self-similarity, stability, flexibility and activity and hotspots, etc.

Corresponding eigenvalues

$\lambda_1^L, \lambda_2^L, \ldots$

$\lambda_1^A, \lambda_2^A, \ldots$

Mark Kac: Can one hear the shape of a drum?
Can one hear the interaction of molecules?

# Geometric Graph Theory

- **Multiscale weighted colored graphs (MWCG)**
- MWCG is about 40% more accurate than Gaussian network model (GNM) in B-factor prediction, based on 364 proteins.
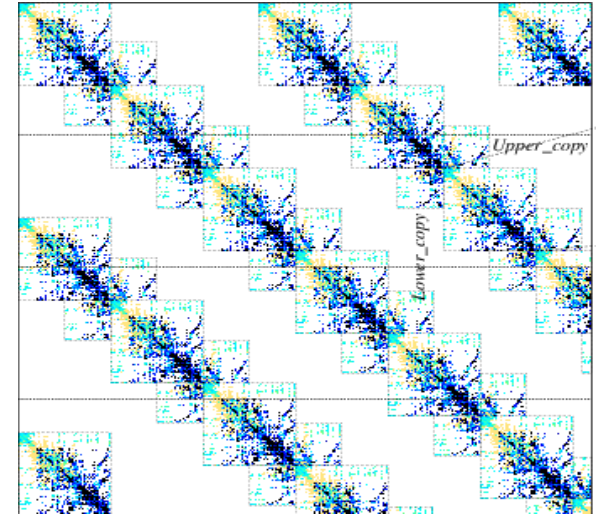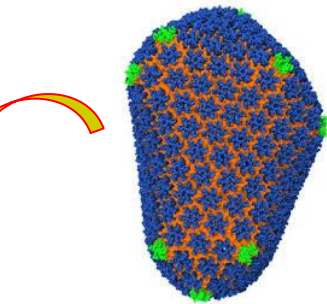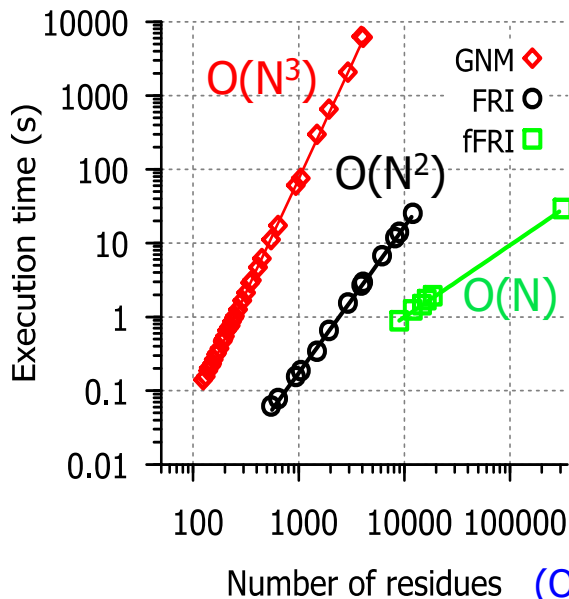
K. Opron





HIV capsid (313,236 residues) would takes GNM 120 years to compute!

$$\Gamma_{ij}(\Phi) = \begin{cases} -\Phi(r_{ij}, \eta), & i \neq j, \\ -\sum_{j, j \neq i}^{N} \Gamma_{ij}, & i = j \end{cases}$$

$$\Phi(r_{ij}, \eta) = 1, \quad r_{ij} \to 0$$

$$\Phi(r_{ij}, \eta) = 0 \quad r_{ij} \to \infty$$

$$\Phi(r_{ij}, \eta) = e^{-(r_{ij}/\eta)^{\kappa}}$$

$$B_i^{FRI} = a(\Gamma_{ii}(\Phi))^{-1}$$

(Opron, Xia and Wei, JCP, 2013; JCP 2014; JCP, 2015; Nguyen, et al, JCIM, 2017, Bramer and Wei, JCP, 2018. Nguyen and Wei, 2018)

# Multiscale: The Poisson-Boltzmann equation

- Discontinuous dielectric constant at the interface
- Non-smooth interface (geometric singularity)
- Singular charges (delta functions)

Chern et al, 2003; Geng, Yu, Wei, JCP, 2007; Geng, Zhao, JCP 2017



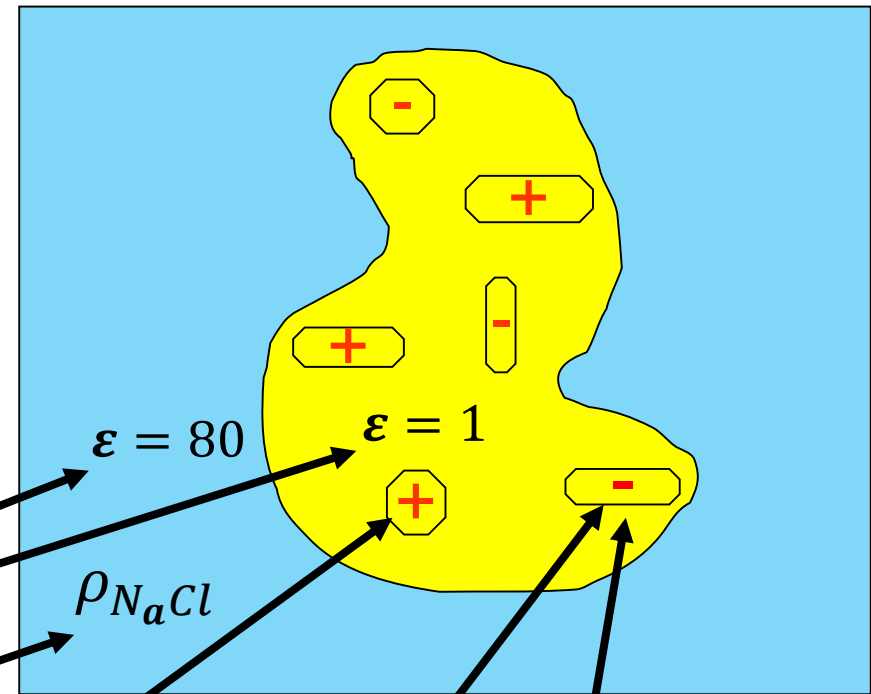$\varepsilon = 80$    $\varepsilon = 1$

$\rho_{N_aCl}$

$$-\boldsymbol{\nabla} \cdot (\boldsymbol{\varepsilon}(\boldsymbol{r})\boldsymbol{\nabla}\boldsymbol{\phi}) = \sum_i q_i c_i e^{-\frac{q_i \phi}{kT}}$$

$$+ \sum_i \left( Q_i \delta(\boldsymbol{r} - \boldsymbol{r_i}) - \boldsymbol{d_i} \cdot \boldsymbol{\nabla}\delta(\boldsymbol{r} - \boldsymbol{r_i}) + \boldsymbol{\Theta_i} : \boldsymbol{\nabla}\boldsymbol{\nabla}\delta(\boldsymbol{r} - \boldsymbol{r_i}) \right)$$

Point charge          Charge polarization (Amoeba)

# MIBPB for solving the Poisson equation with protein interface



ESES

(Yu, Geng, Wei, JCP 2007)

(Liu, Wang, Zhao, Tong, Wei, JCC 2017)

**Relative solvation energy deviations over grid refinement for 947 proteins in the Amber test set**

(Wang, Wei, 2015)

**Electrostatic binding energies of 14 RNA-protein complexes over grid refinement**
(Nguyen, Wang, Wei, JCC, 2015)

# MIBPB
## Wei @ MICHIGAN STATE

MIBPB is a software package for obtaining electrostatic potential and solvation free energy via solving the Poisson-Boltzmann (PB) equation. It makes use of the second order convergent MIB t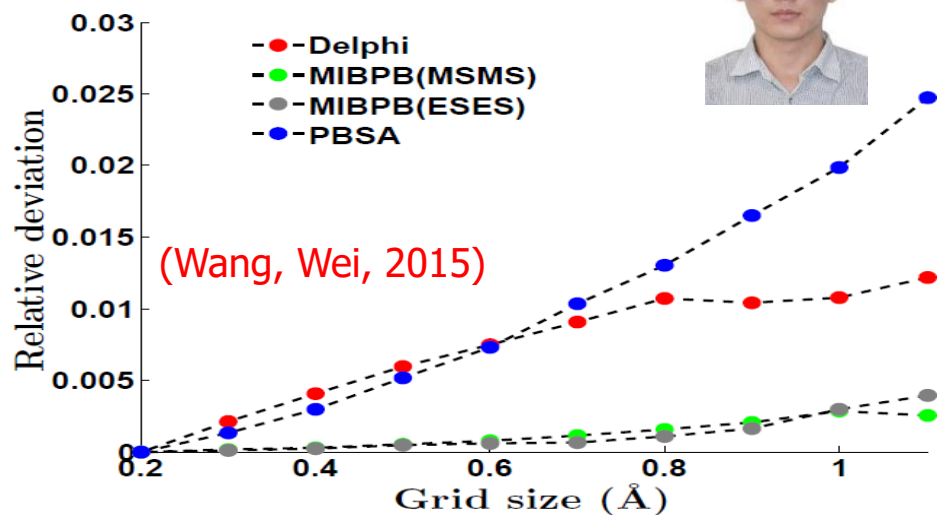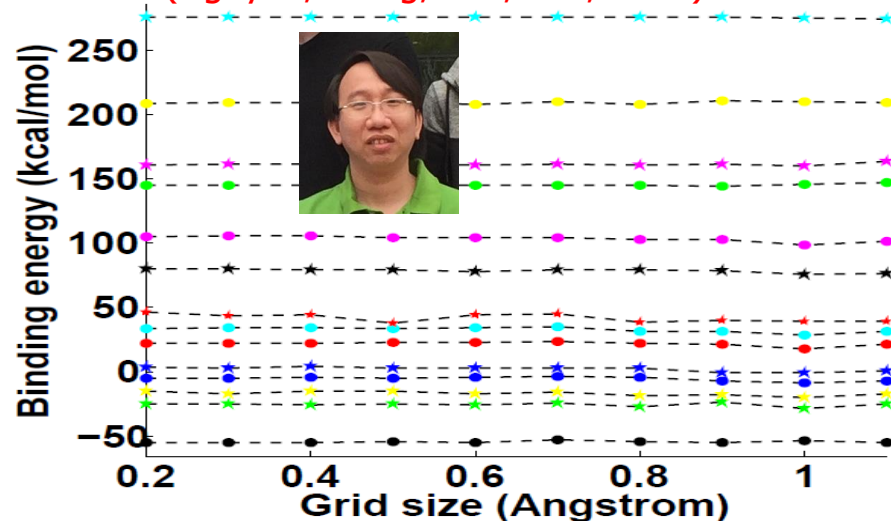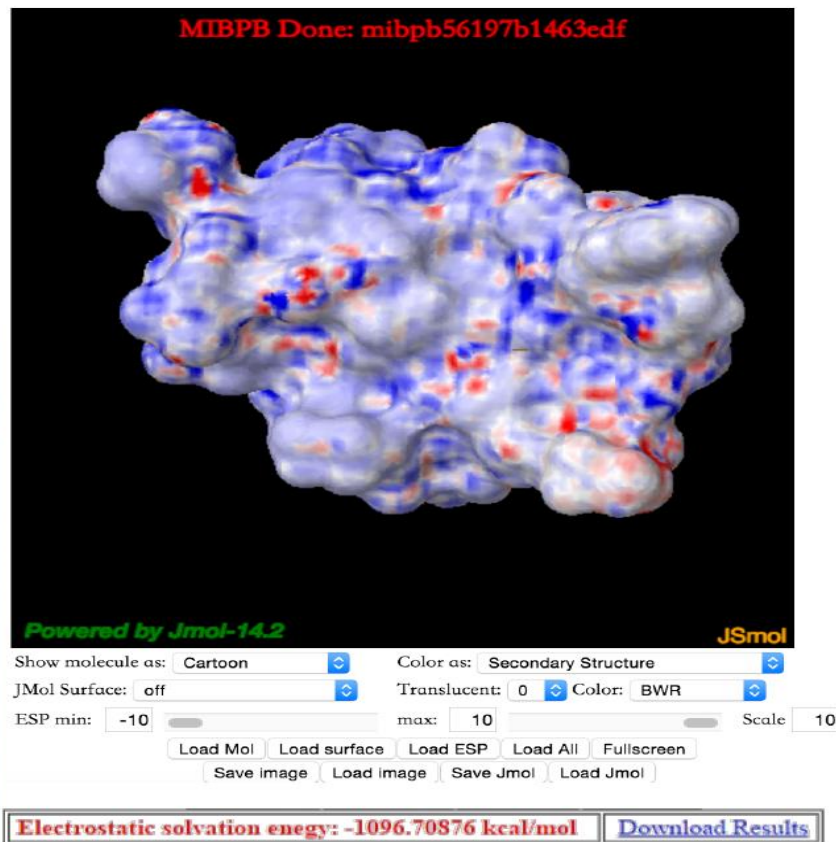echnique and is essentially grid independent. Its mean relative error is less than 0.5% for about 1000 test proteins when the grid size is refined from 1.1 to 0.2 Angstrom. [User Manual]

**MIBPB Done: mibpb56197b1463edf**

| Input_File: | ⦿ PDB ID: ____ Chains: * |
| | ◯ User File: 选择文件 未选择任何文件    * |

| MIBPB Options | 1.0 | Interior Dielectric |
| | 80 | Outerior Dielectric |
| | 0.8 | Grid Resolution |
| | 0.0 | Ion Strength |
| | 1.4 | Surface Probe Radius for MIBPB |
| | Linearized PB: ☑ Yes |
| | Simplified Solver: ☐ Yes |

| Resulting Surface Options ☑ Yes | 1.4 | Probe Radius |
| | 0.8 | Grid Resolution |
| | 2.0 | Grid Extension |

| PDB2PQR Options ⦿ Yes | Force Field: AMBER |
| | Protonation: at pH: 7.0 by: PROPKA |
| | Remove Water: ☑ Yes |
| | Remove Hydrogen: ☐ Yes |
| | Only assign charges and radius: ☐ Yes |

| Small molecule to PQR Options ◯ Yes | Charge Type: AM1-BCC |
| | Radius Type: mbondi |

| pKa Calculation ☐ Yes | Residue Type: ASP |
| | Residue ID: ____ Analyze online PDB |

| Job Title: | _____ |
| User Email: | _____ |

Default   Submit   Clear Job

Show molecule as: Cartoon    Color as: Secondary Structure
JMol Surface: off    Translucent: 0   Color: BWR
ESP min: -10 ____    max: 10 ____   Scale 10

Load Mol   Load surface   Load ESP   Load All   Fullscreen
Save image   Load image   Save Jmol   Load Jmol

Powered by Jmol-14.2    JSmol

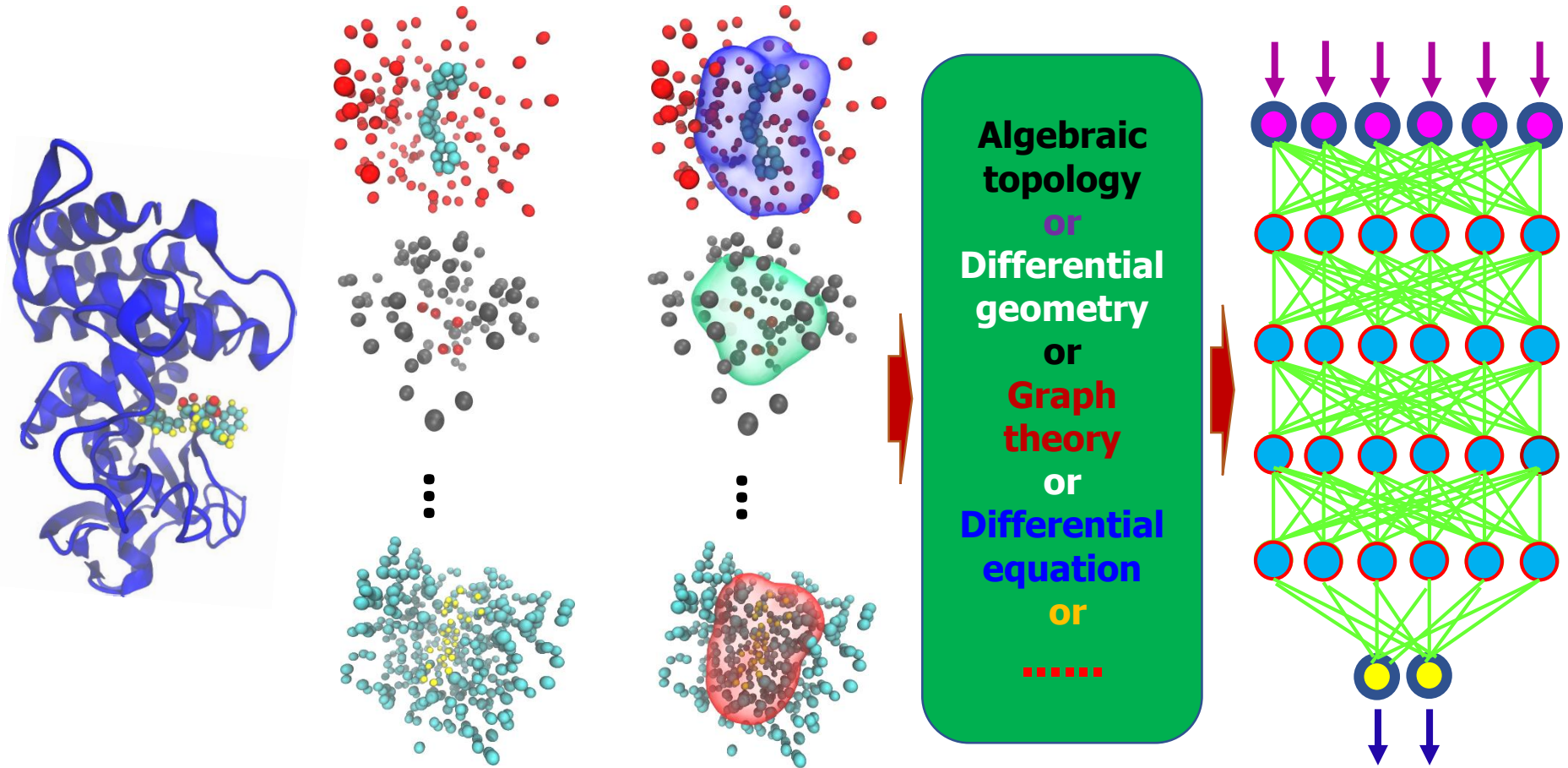**Electrostatic solvation enegy: -1096.70876 kcal/mol**   **Download Results**

## 🖥 DOWNLOAD

- For academic/governmental users, you may download and use MIBPB for free under a license agreement. Please follow the instructions below to register with us and download MIBPB.
- For industrial/commercial users, a moderate license fee may apply. Please contact us directly at wei@math.msu.edu.
- If you have any questions or bugs to report, please feel free to contact: wangbaonj@gmail.com

(Wang, Zhao, Wei, 2015)

# Mathematical deep learning
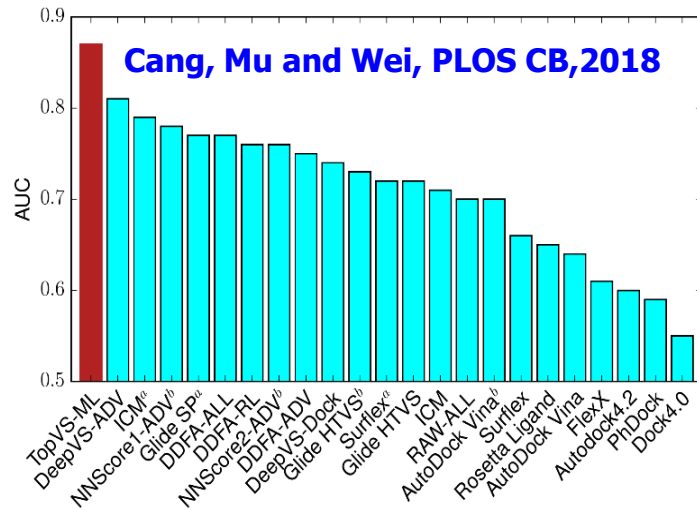


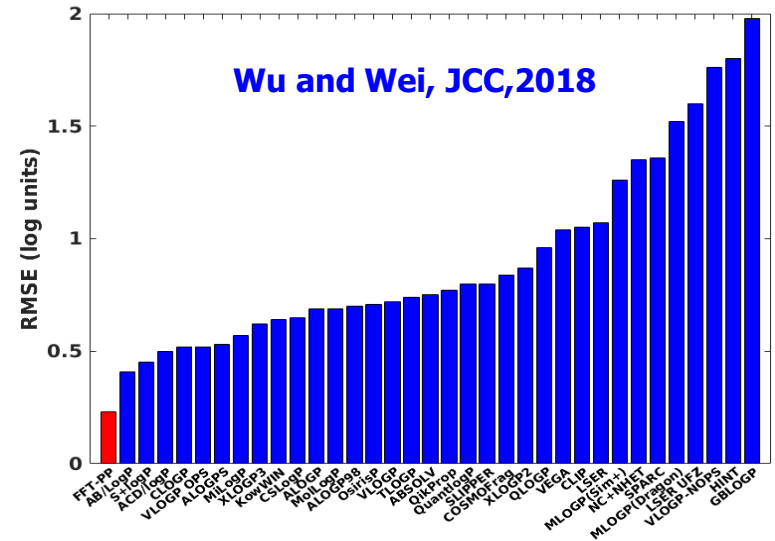Protein-ligand complex → Element specific groups → Element interactive manifolds → Various Mathematical features → Machine learning prediction

Algebraic topology
or
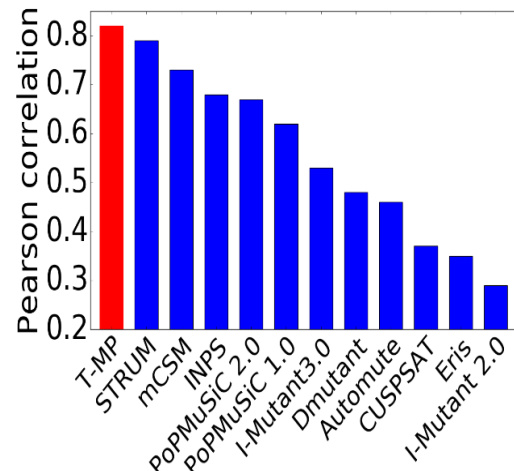Differential geometry
or
Graph theory
or
Differential equation
or
......

# Topological learning based predictions

### Classification of ligands & decoys
### DUD database  128,374 protein-ligand/decoy pairs

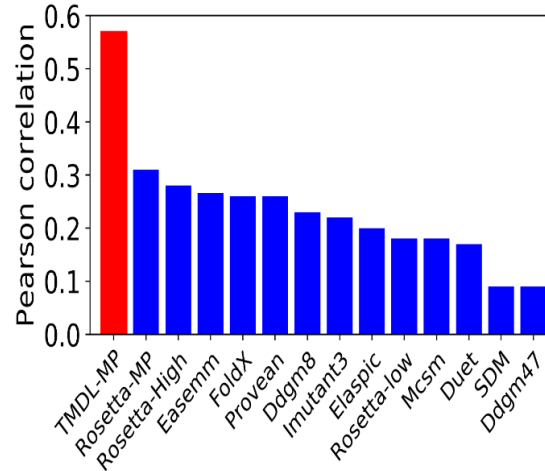Cang, Mu and Wei, PLOS CB,2018

### Prediction RMSD of LogP (Star set)

Wu and Wei, JCC,2018

### Predicting mutations on 2648 globular proteins

### Predicting mutations on 223 membrane proteins

(Cang and Wei, Bioinformatics, 2017)

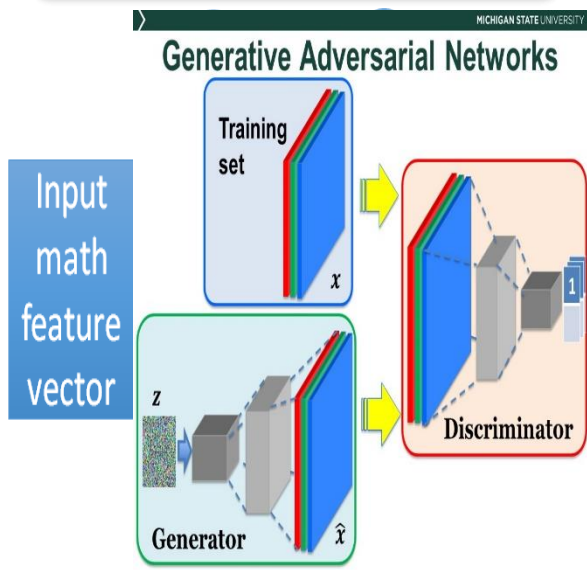### Binding affinity prediction of PDBBind v2013 core set of 195 protein-ligand complexes

Cang and Wei, PLOS CB,2017

# Drug Design Data Resource (D3R) Grand Challenge

**Given data**



Primary structure
amino acid sequence

**Math based GAN**



MICHIGAN STATE UNIVERSITY

Generative Adversarial Networks

Input math feature vector
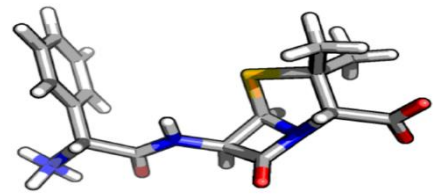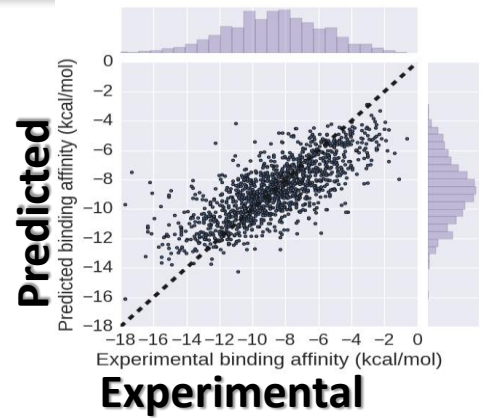
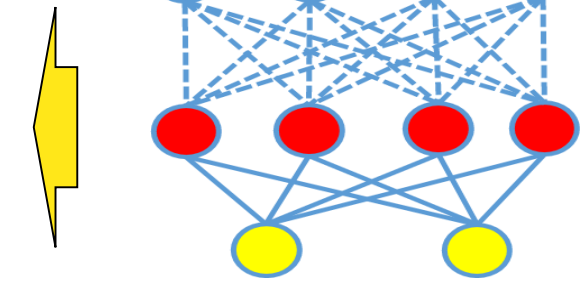Training set

$x$

$z$

Generator

$\hat{x}$

Discriminator

**Predicted complex**



**Final predictions to be compared with experiments**



Predicted binding affinity (kcal/mol)

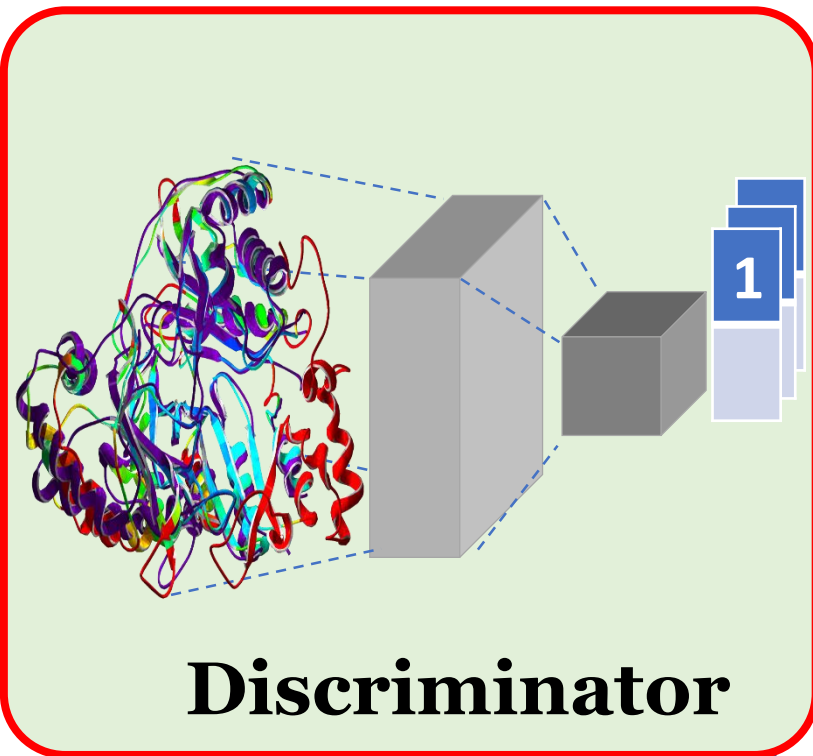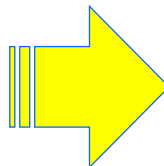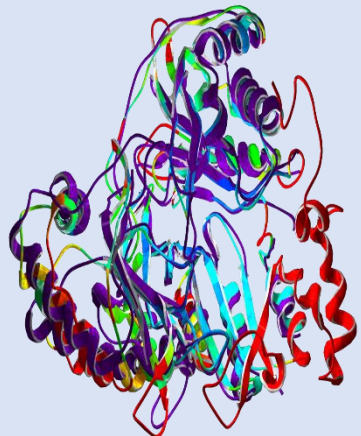Experimental binding affinity (kcal/mol)

**Experimental**

**Drug pose**



(Nguyen et al, JCAMD, 2018)

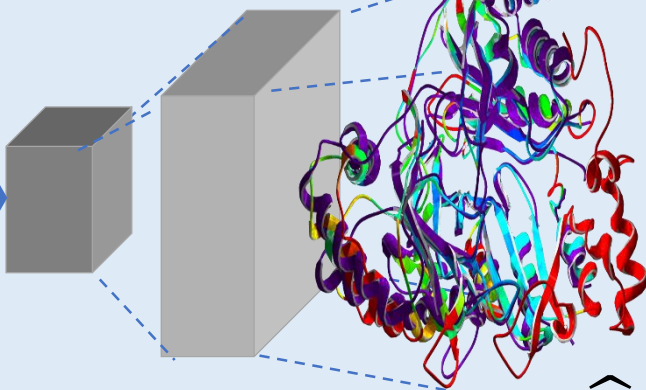# Generative Adversarial Networks for Drug Design

# D3R Grand Challenge 2     (2016-2017)

**Given:** Farnesoid X receptor (FXR) and 102 ligands
**Tasks:** Dock 102 ligands to FXR, and predict their poses, binding free energies and energy ranking

## Stage 1

Pose Predictions (partials)
Scoring (partials)
Free Energy Set 1 (partials)
Free Energy Set 2 (partials)

## Stage 2

Scoring (partials)
Free Energy Set 1 (partials)
Free Energy Set 2 (partials)

(Nguyen et al, JCAMD, 2018)

Dr D Nguyen



**Grand Challenge 2**

Free Energy Set 1 (Stage 2) - Kendall's Tau

Our prediction

Filled circle indicates an incomplete set of predictions
Green circle indicates your predictions (requires login)



**Grand Challenge 2**

Free Energy Set 1 (Stage 1) - RMSD

Our prediction

Green circle indicates your predictions (requires login)

# D3R Grand Challenge 3 (2017-2018)

(Nguyen et al, JCAMD, 2018)

## Pose Prediction

**Cathepsin Stage 1A**
Pose Predictions (partials)

**Cathepsin Stage 1B**
Pose Prediction

## Affinity Rankings excluding Kds > 10 µM

**Cathepsin Stage 1**
Scoring (partials)
Free Energy Set

**Cathepsin Stage 2**
Scoring (partials)
Free Energy Set

**VEGFR2**
Scoring (partials)

**JAK2 SC2**
Scoring (partials)

**p38-α**
Scoring

**JAK2 SC3**
Scoring
Free Energy Set 🥇

**TIE2**
Scoring 🥇
Free Energy Set 2 🥇

**ABL1**
Scoring (partials) 🥇

## Active / Inactive Classification

**VEGFR2**
Scoring (partials)

**JAK2 SC2**
Scoring (partials)

**p38-α**
Scoring (partials)

**JAK2 SC3**
Scoring
Free Energy Set 🥇

**TIE2**
Scoring (partials) 🥇
Free Energy Set 1 🥇

**ABL1**
Scoring (partials)

## Affinity Rankings for Cocrystalized Ligands

**Cathepsin Stage 1**
Scoring (partials)
Free Energy Set 🥇

**Cathepsin Stage 2**
Scoring (partials) 🥇
Free Energy Set

Cathepsin S    Kinase: p38-α

Zixuan Cang    Dr D Nguyen

# D3R Grand Challenge 4 (2018-2019)

## Pose Predictions

**BACE Stage 1A**
Pose Predictions (Partials) 🥇 2/3 🥈 2/3

**BACE Stage 1B**
Pose Prediction (Partials) 🥈 2/2 🥉 1/2

## Affinity Predictions

**Cathepsin Stage 1**
Combined Ligand and Structure Based Scoring 🥇 2/5 🥈 2/3 🥉 2/4

Ligand Based Scoring (No participation)

Structure Based Scoring 🥇 2/4 🥈 3/3 🥉 3/3

Free Energy Set 🥇 1/7 🥈 1/7 🥉 2/5

Dr. Kaifu Gao  Dr. D Nguyen

**BACE Stage 1**
Combined Ligand and Structure (No participation)

Ligand Based Scoring (Partials) (No participation)

Structure Based Scoring (Partials) (No participation)

Free Energy Set (No participation)

**BACE Stage 2**
Combined Ligand and Structure

Ligand Based Scoring (No participation)

Structure Based Scoring (Partials)

Free Energy Set 🥈 3/4 🥉 1/4

| | Algebraic topology | Differential topology | Geometric topology | |
|---|---|---|---|---|
| Differential equation | | | | Number theory |
| Algebraic graph | | | | Algebraic geometry |
| Structural graph | | Biology became microscopic (i.e., molecular) in 1960s and added an omics dimension around the dawn of the millennium. | | Differential geometry |
| Exyreme graph | | **Driving by mathematics, biology is transforming from qualitative, phenomenological and descriptive to quantitative, predictive and analytical.** | | Distance geometry |
| Geometric algebra | | The last frontier of science is biology, while the last frontier of biology is mathematics. | | Lie algebra |
| | Complex analysis | Real analysis | Stochastic analysis | |