

Lecture 9: Mathematics for Biomolecular Data

Guowei Wei

Mathematics

Michigan State University

<https://users.math.msu.edu/users/wei/>

NSF-CBMS Conference on Mathematical Molecular Bioscience and
Biophysics

University of Alabama

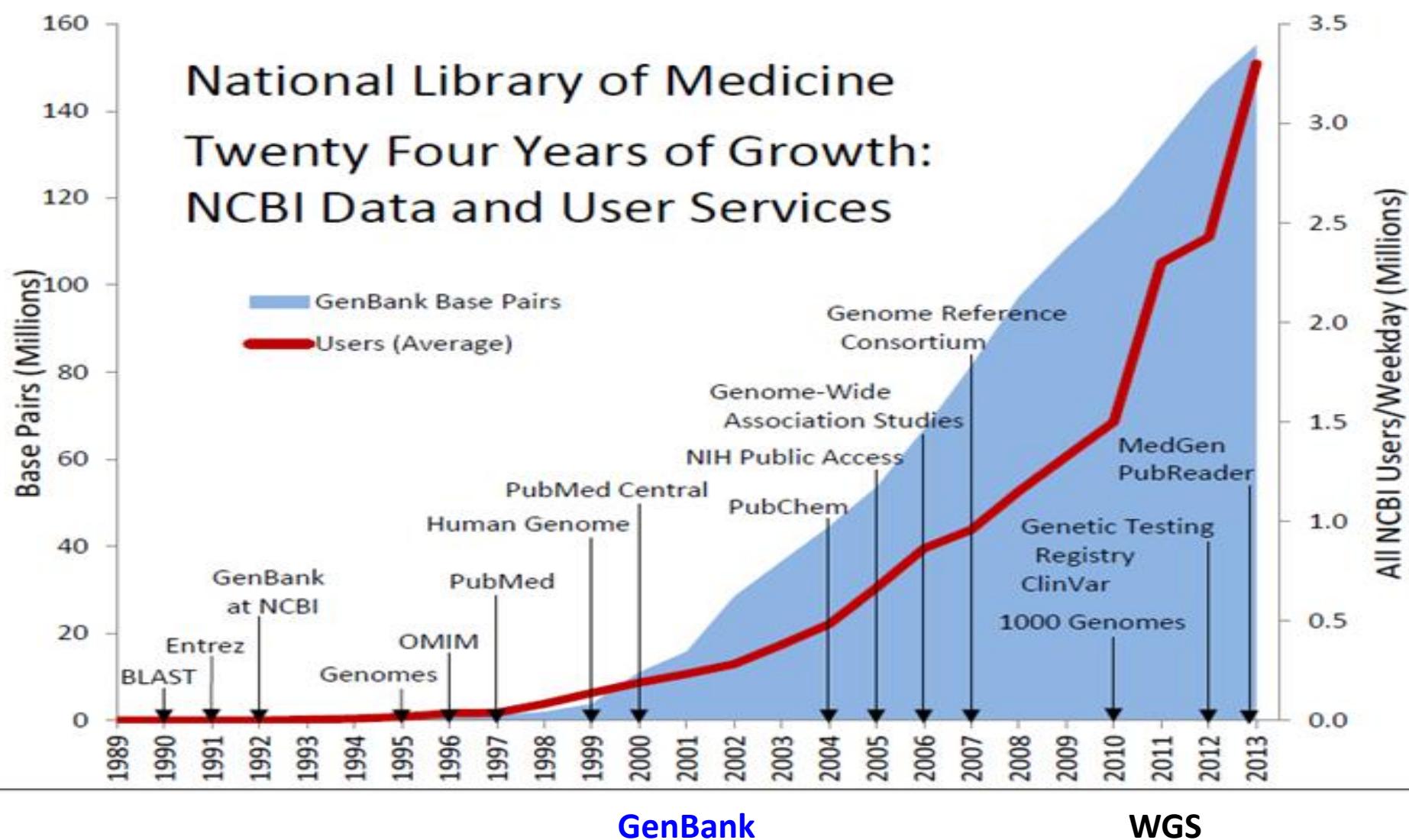
Tuscaloosa, May, 13-17, 2019

Grant support: NSF, NIH, MSU, BMS, and Pfizer



GenBank (GB)

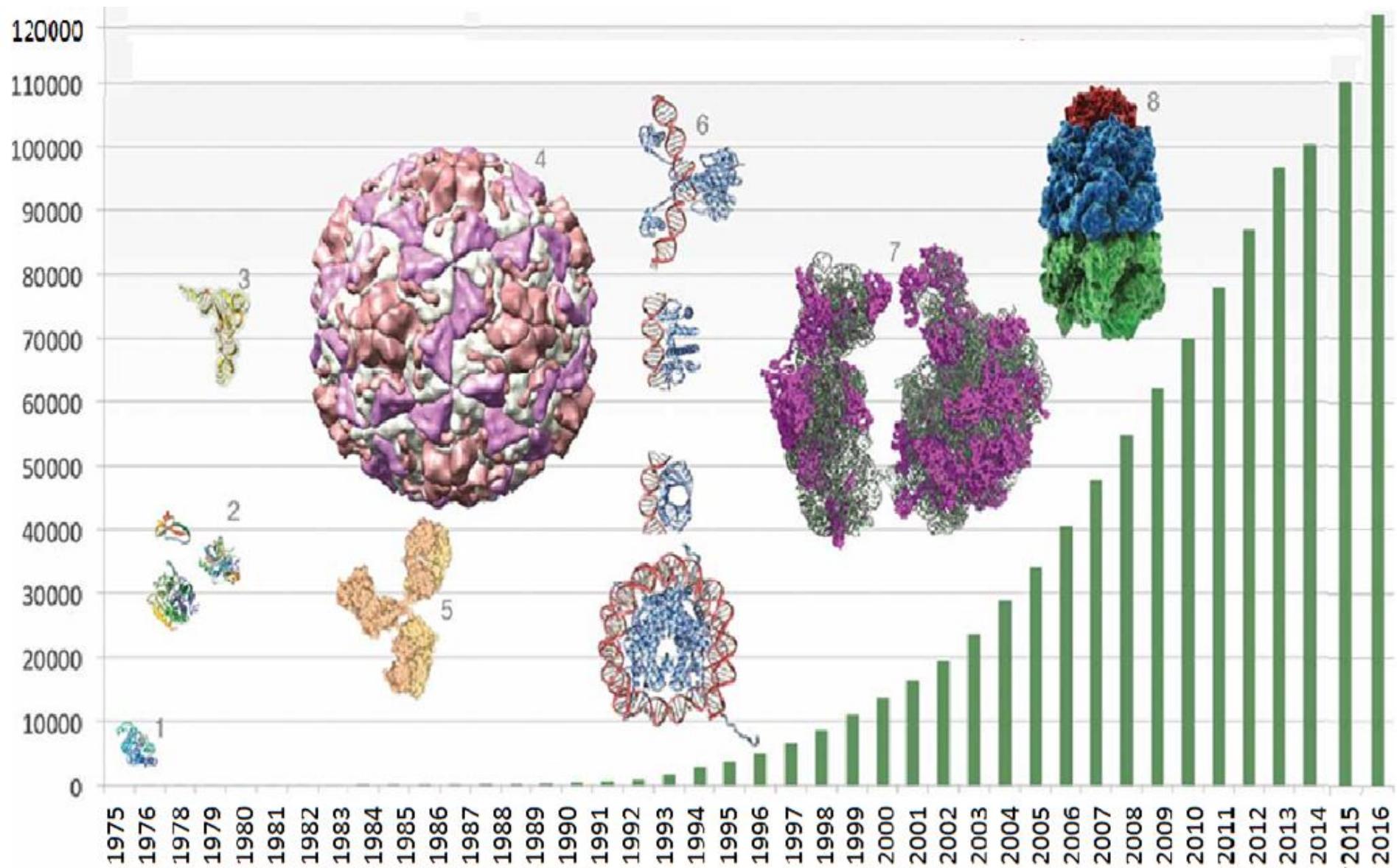
National Library of Medicine
Twenty Four Years of Growth:
NCBI Data and User Services



GenBank

Release	Date	Bases	Sequences	Bases	Sequences
231	Apr 2019	321680566570	212775414	4421986382065	993732214

Protein Data Bank (PDB)



PDB: 151,579 Biological macromolecular structures as 04/30/2019

Protein-ligand (-protein) binding

Protein (P) and ligand (L) form a protein-ligand complex (PL):



The association and dissociation constants are

$$K_a = \frac{[PL]}{[P][L]}, \quad K_d = \frac{[P][L]}{[PL]}.$$

Binding affinity: $\Delta G = RT \ln K_d$

Database:

ChEMBL (as 02/28/2019):

Targets: 12,091

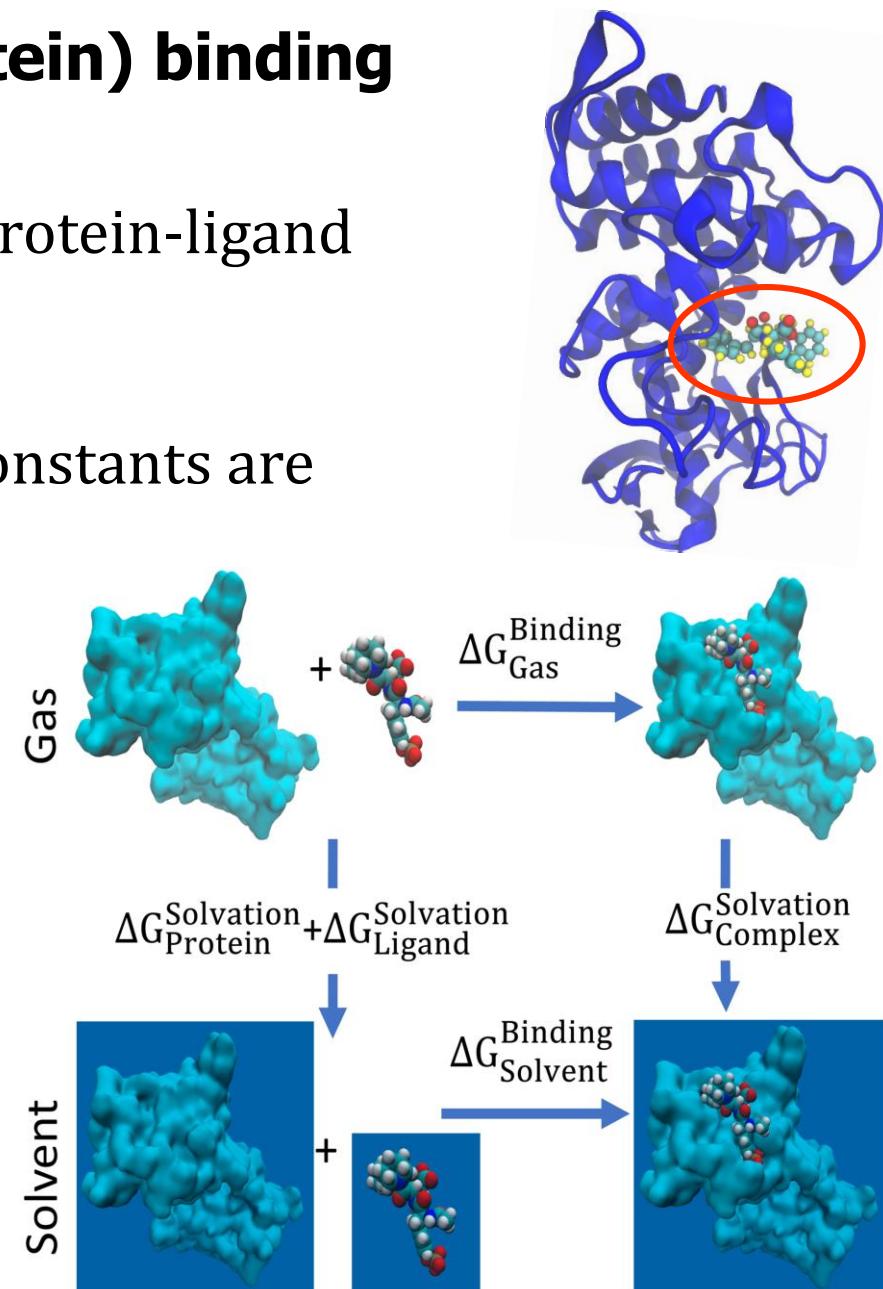
Compound records: 2,275,906

Distinct compounds: 1,828,820

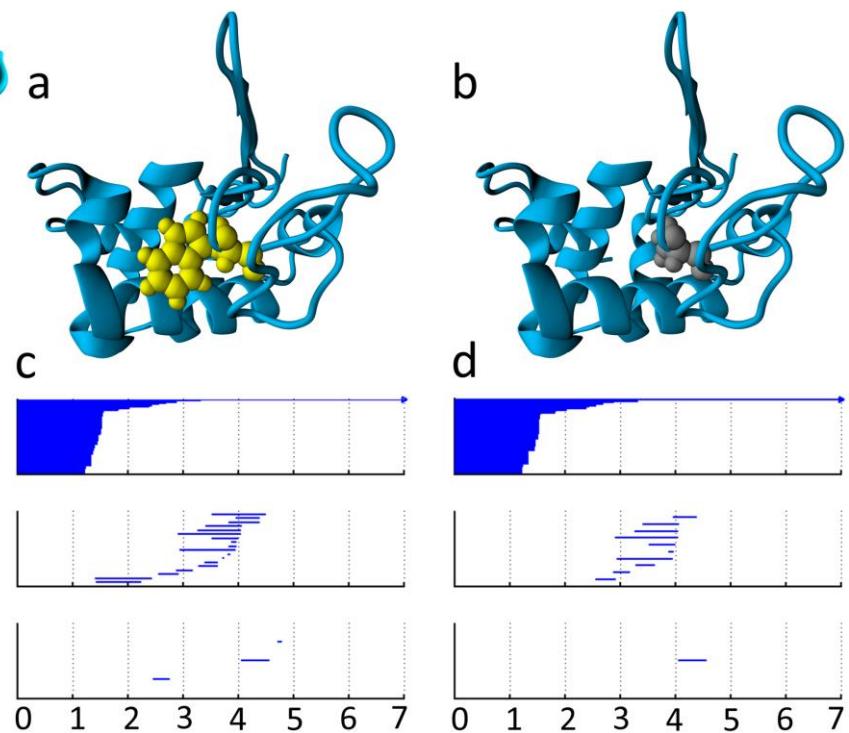
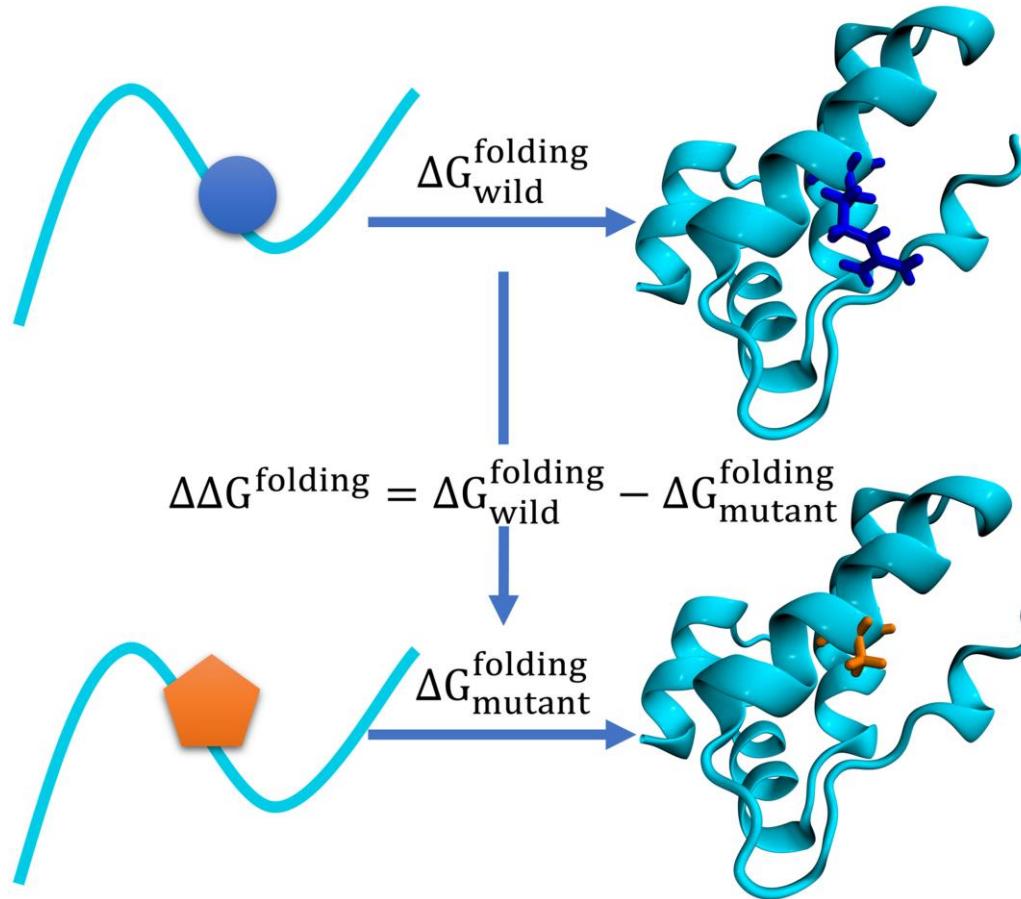
Activities: 15,207,914

Publications: 69,861

Binding DB, PubChem, PDBbind, K_i Database.



Mutation induced free energy change



Protein-Protein interaction

Structures via
X-ray
crystallography,
NMR ,
Cryo-EM, ...

Properties via
Isoth. Titr. Calor.
Surf. Plasm. Res.
Fluorescence
Mutation, ...

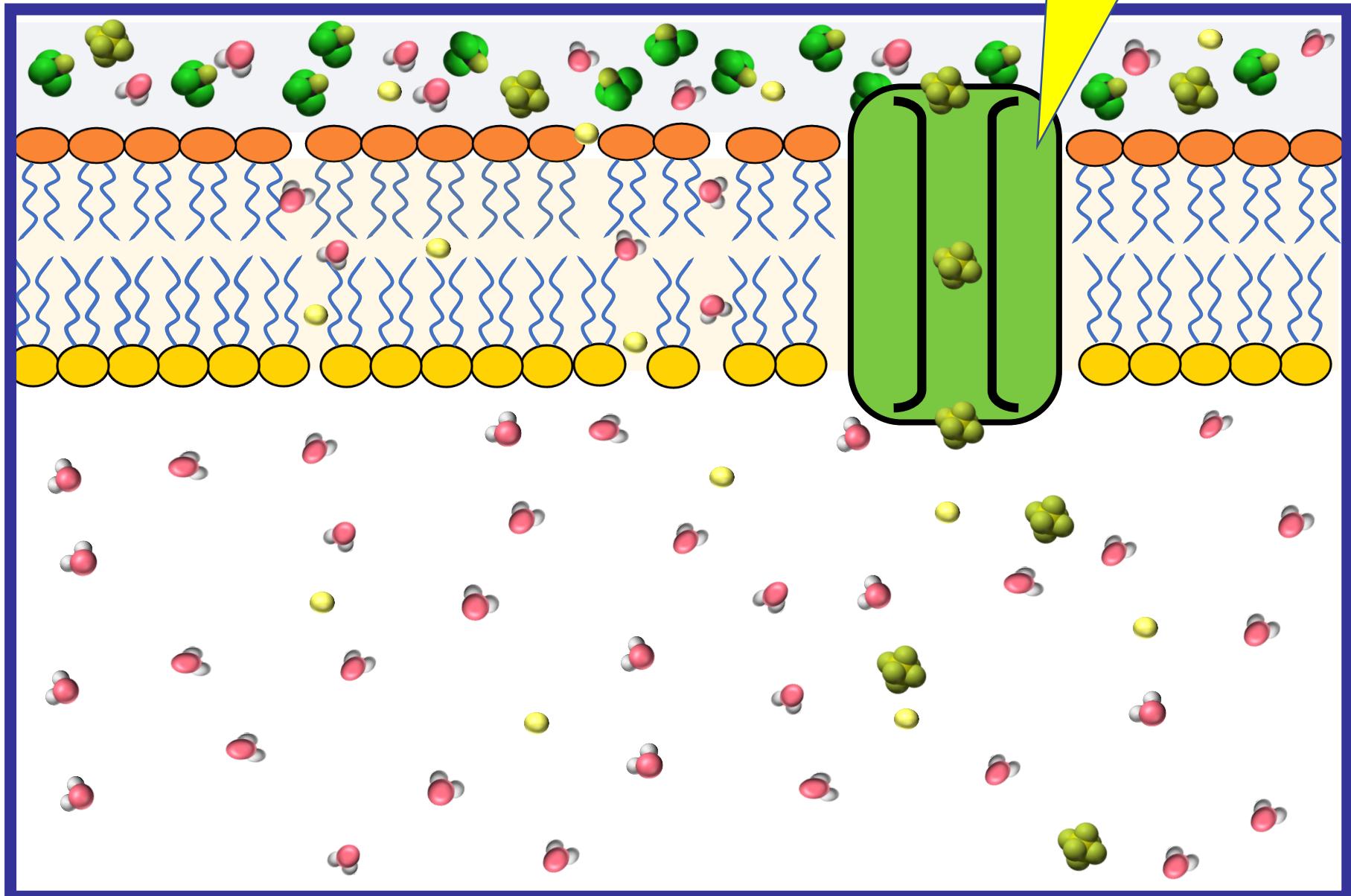
Protein-protein
Interaction

Models via
Docking,
MD, Statistics,
Thermodynamics,
...
...

Predictions via
Scoring,
Empirical,
Machine learning,
...
...

Permeability

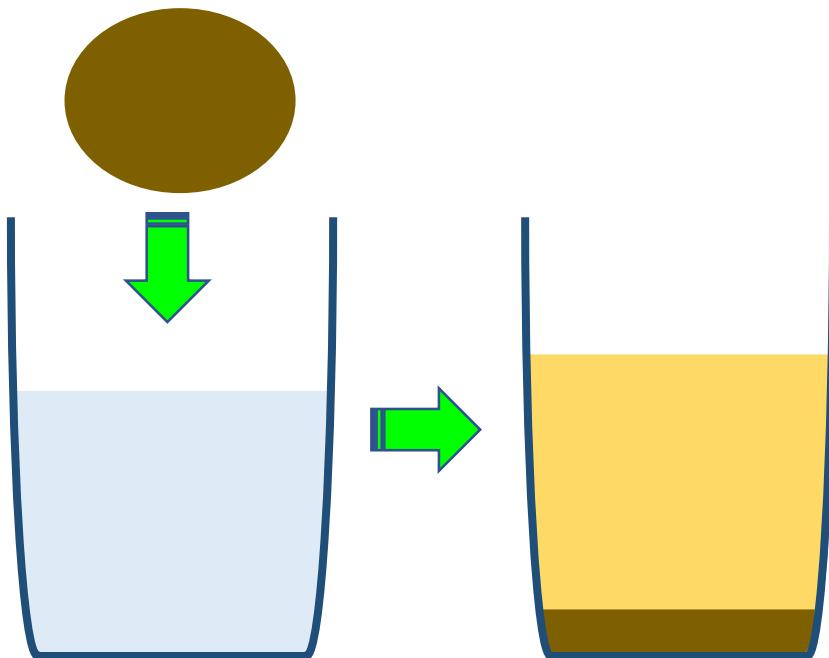
Transporter



Solubility and partition coefficient

Solubility is commonly expressed as a concentration; for example, as g of solute per kg of solvent.

Partition coefficient is defined as a particular ratio of the concentrations of a solute between the two solvents



Toxicity

Toxicity: The degree to which a substance (a toxin or poison) can harm humans or animals. **Drug toxicity** occurs when a person has accumulated too much of a drug in his bloodstream, leading to adverse effects on the body.

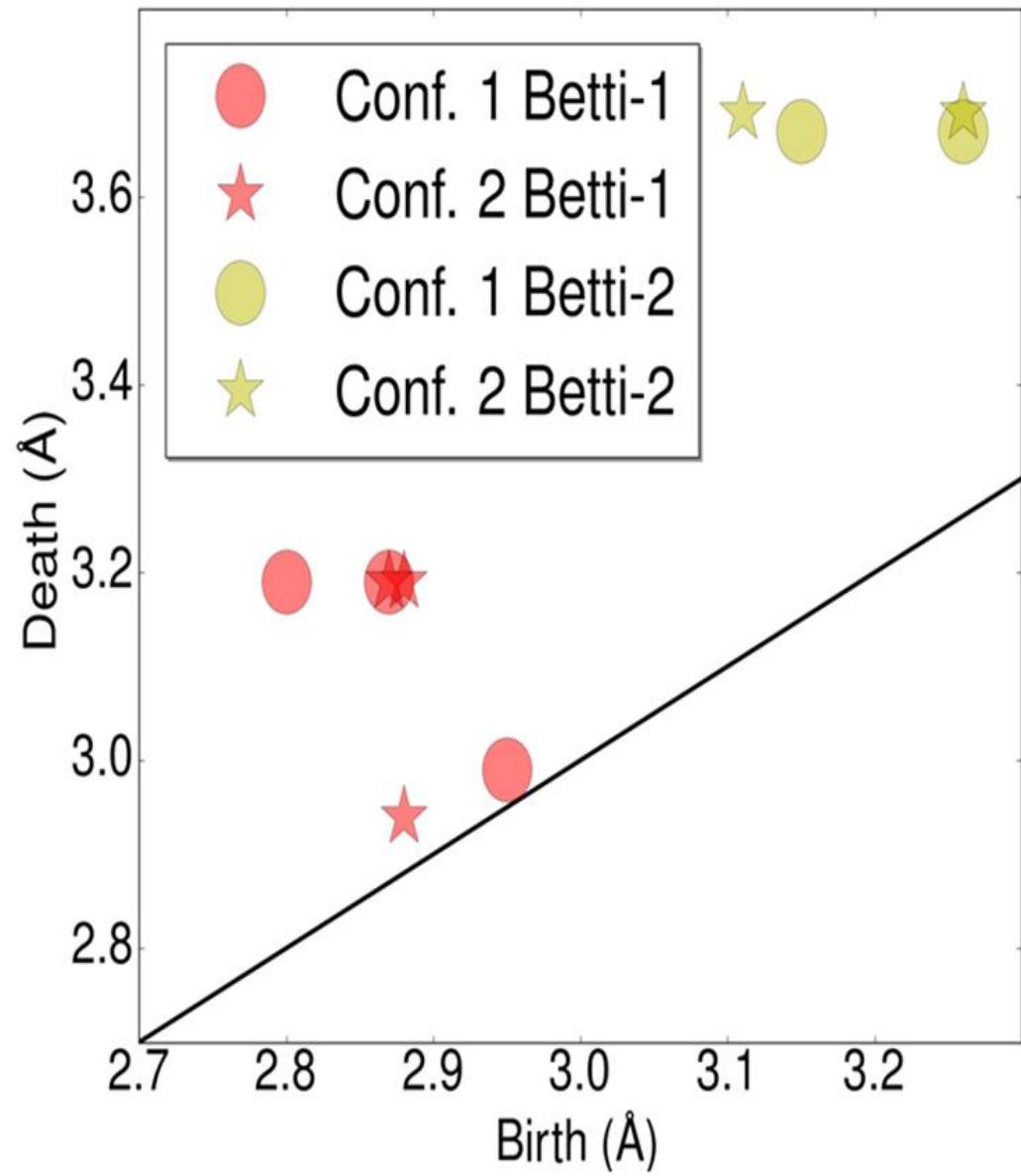
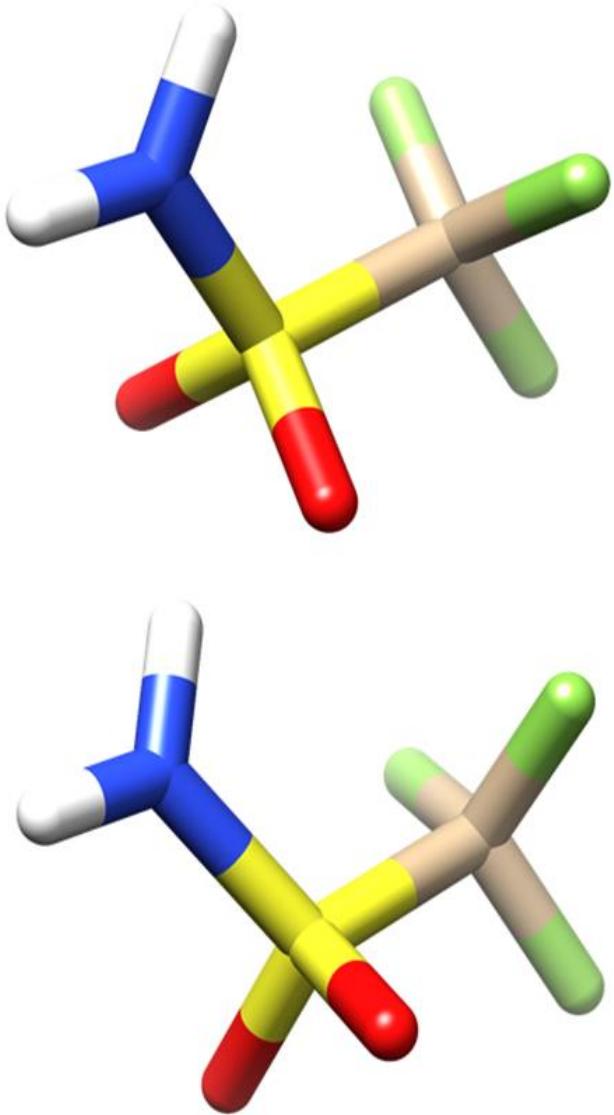
Bioassays:

LD50 is defined as the **lethal dose** at which 50% of the population is killed in a given period of time.

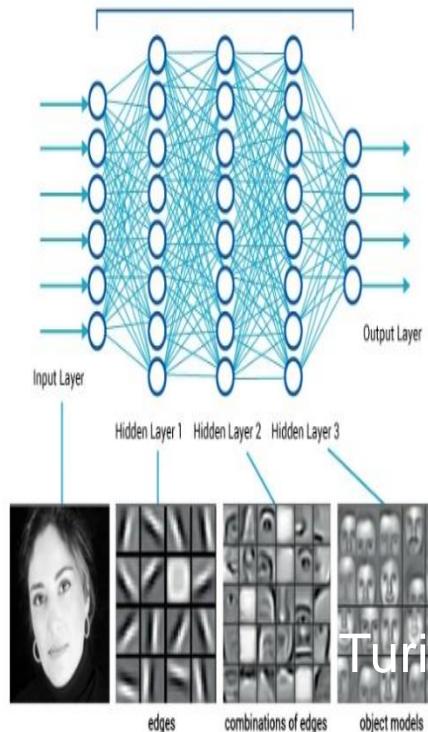
LC50 is the **lethal concentration** required to kill 50% of the population. The LC50 is a measure, *e.g.* in mg/l, of the concentration of the toxin whereas a dose is a more general term.



Topological fingerprints of stereoisomers (rotamers)



Artificial intelligence and machine learning



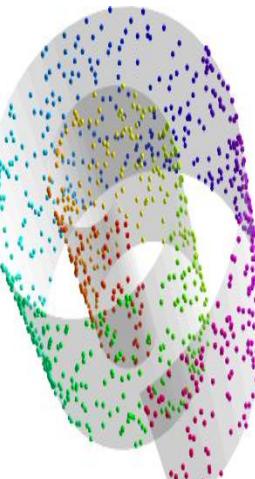
The logo for the 2013 Turing Award. It features a large blue sphere in the center, surrounded by several smaller white spheres of varying sizes. To the left of the sphere, the word "Turing" is written in a white serif font. To the right, the word "Award" is written in a white sans-serif font. Below the sphere, the year "2013" is written in a white sans-serif font. The background is black.



Turing Award 2019



Geoffrey Hinton, Yann LeCun, and Yoshua Bengio



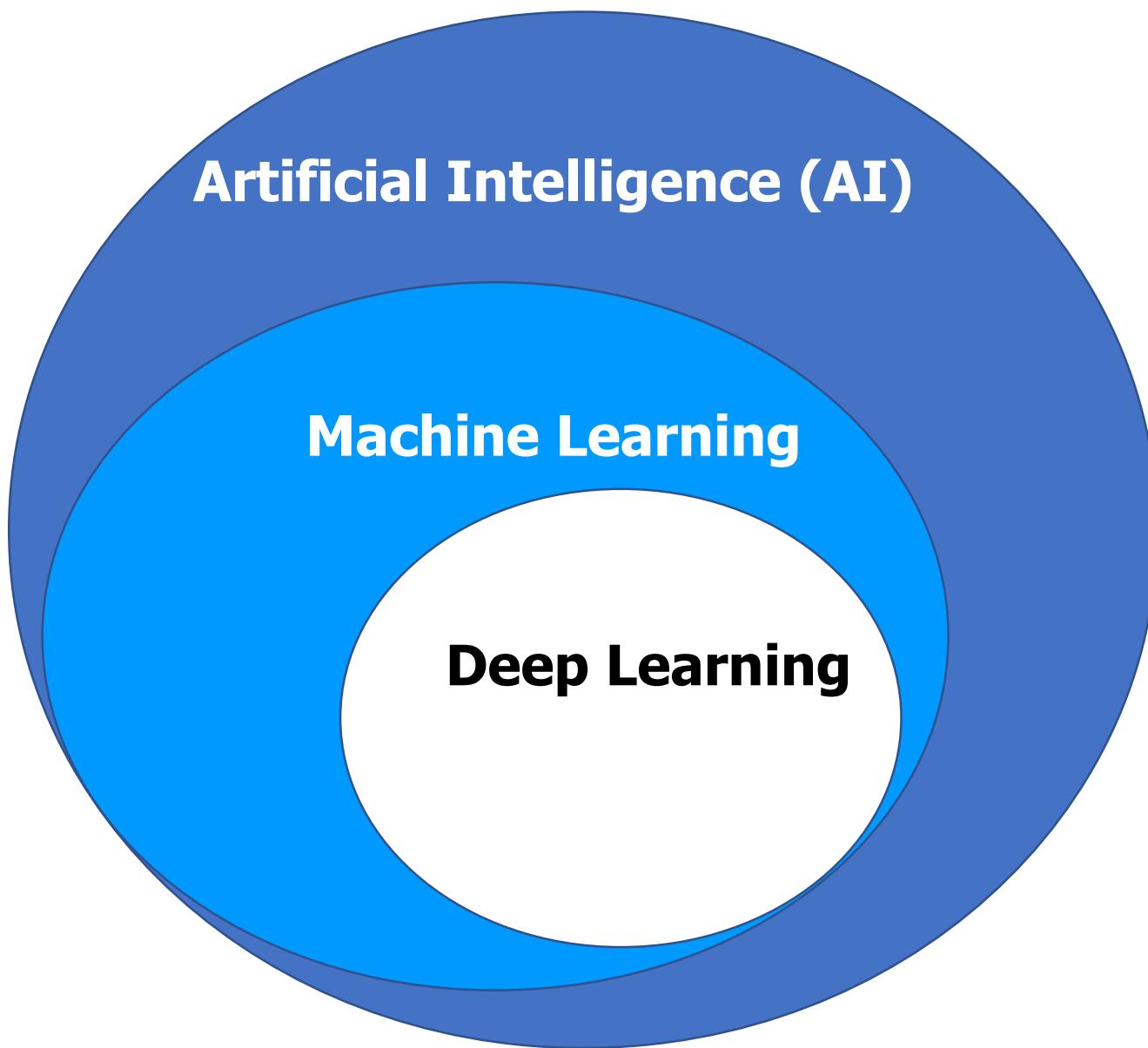
data statistical

Learning classification ensure increasingly structures able get

ns clustering unimportant important



AI and Machine Learning



Basic math needed in Machine Learning

- Calculus
- Linear algebra
- Probability
- Statistics

Advanced math involved in Machine Learning

- Topology
- Geometry
- Analysis
- Algebra
- Graph theory
- Dynamical system
- ODE and PDE
- Combinatorics

Machine Learning Methods and Data

- Unsupervised learning (clustering):

$$\text{Data: } \mathcal{D} = \{\mathbf{x}^{(i)} \mid \mathbf{x}^{(i)} \in \mathbb{R}^n\}_{i=1}^M$$

- Supervised learning (classification):

$$\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}) \mid \mathbf{x}^{(i)} \in \mathbb{R}^n, y^{(i)} \in \{C_1, C_2, \dots, C_K\}\}_{i=1}^M$$

- Supervised learning (regression):

$$\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}) \mid \mathbf{x}^{(i)} \in \mathbb{R}^n, y^{(i)} \in \mathbb{R}\}_{i=1}^M$$

- Semi-supervised learning (regression):

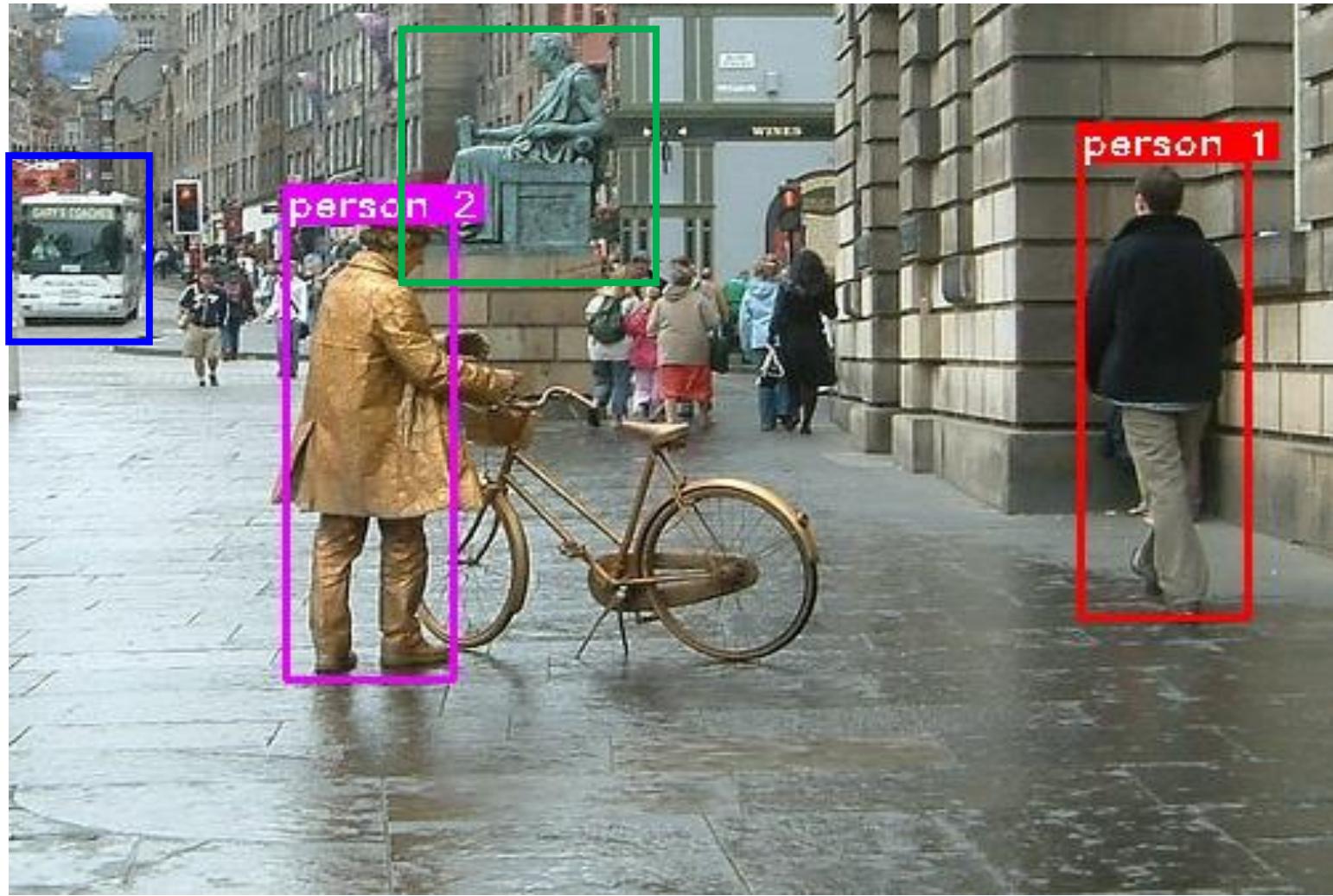
$$\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}) \mid \mathbf{x}^{(i)} \in \mathbb{R}^n, i = 1, 2, \dots, M; y^{(i)} \in \mathbb{R}, i = 1, 2, \dots, M' < M\}$$

- Semi-supervised learning (Classification):

$$\mathcal{D} = \left\{ (\mathbf{x}^{(i)}, y^{(i)}) \middle| \begin{array}{l} \mathbf{x}^{(i)} \in \mathbb{R}^n, i = 1, 2, \dots, M; \\ y^{(i)} \in \{C_1, C_2, \dots, C_K\}, i = 1, 2, \dots, M' < M \end{array} \right\}$$

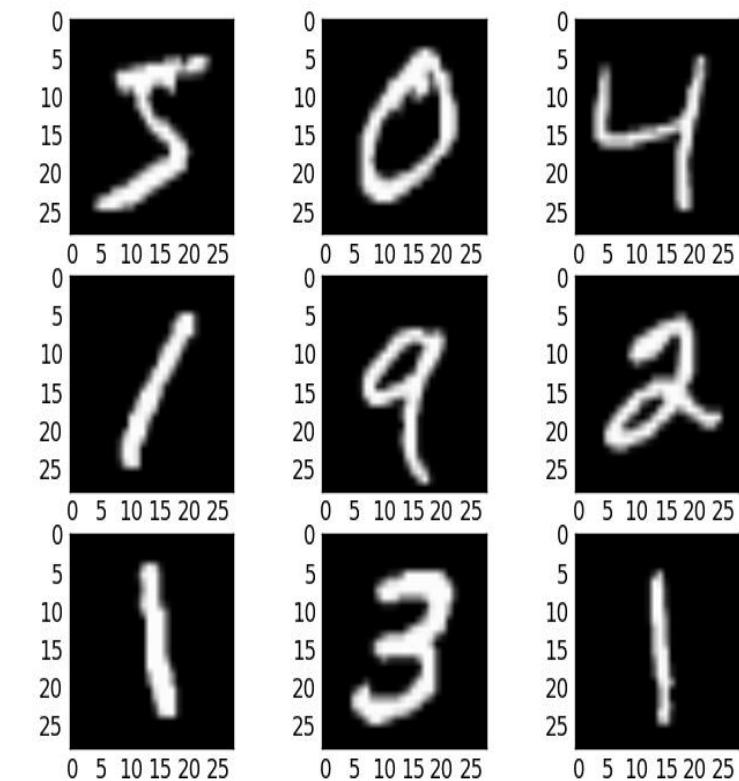
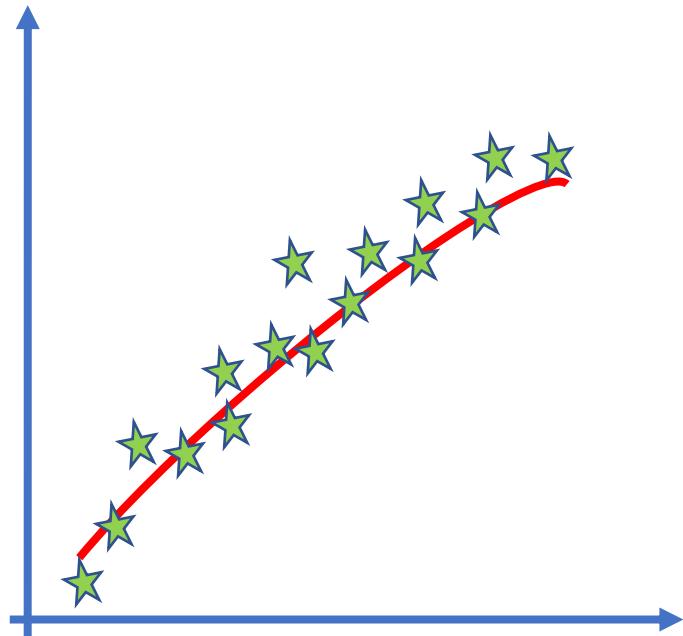
Unsupervised Learning

- K-means
- Generative Adversarial Networks
- Reinforced learning,



Supervised Learning

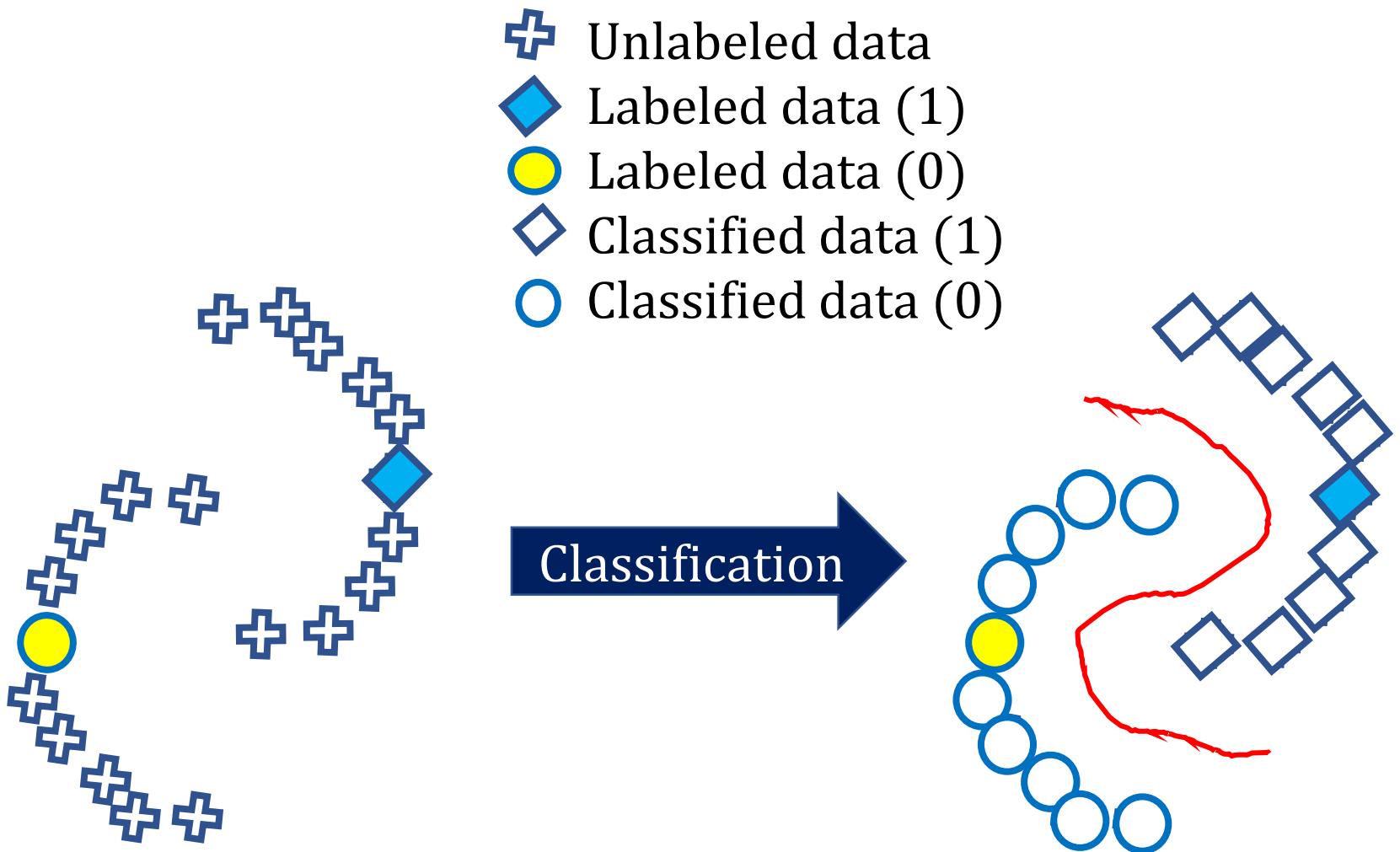
- Regression: Linear regression, Logistic regression, SVM, Random Forest, GBDT, CNN, RNN, etc.
- Classification: K-NN, SVM, Random Forest, CNN, RNN, etc.



Semi-supervised Learning (classification)

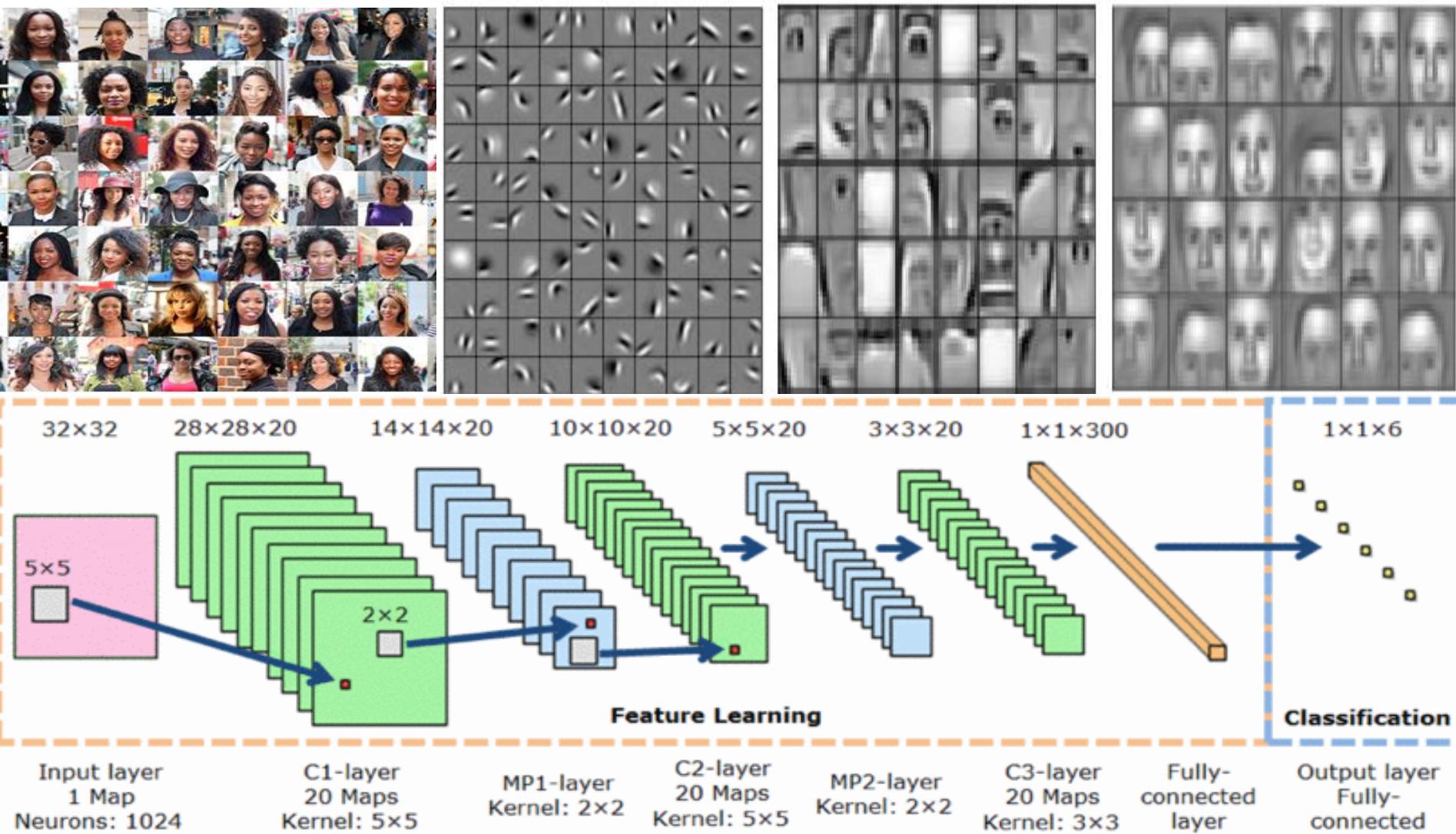
Predict unlabeled data from labeled information

- Manifold learning
- Kernel SVM



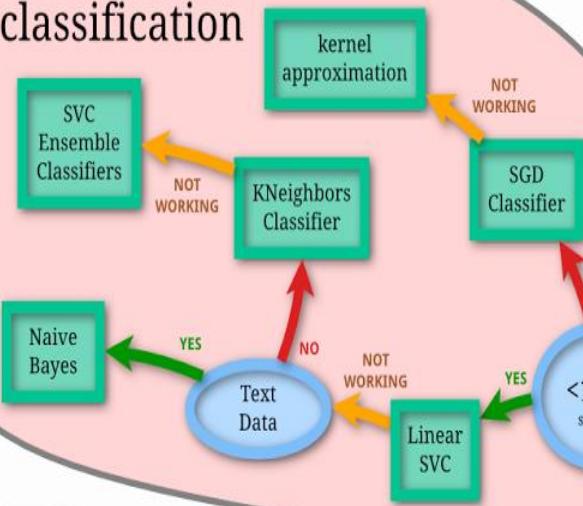
Artificial Intelligence & Deep learning

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, ...

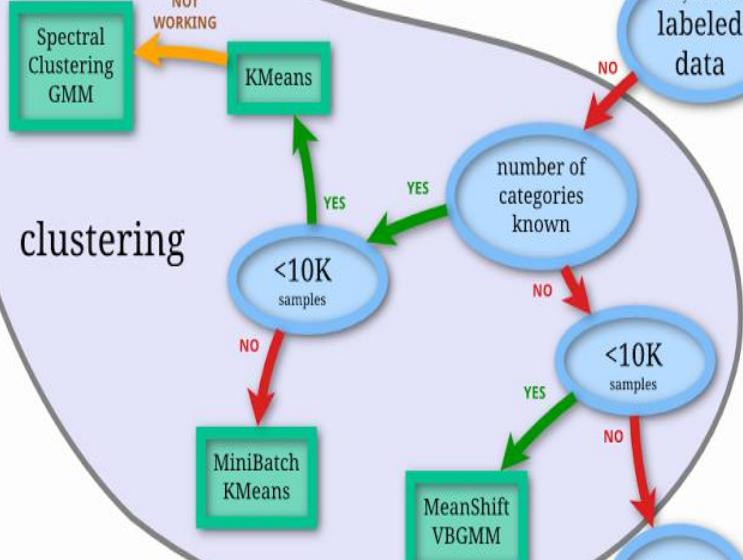


Machine Learning Summary

classification



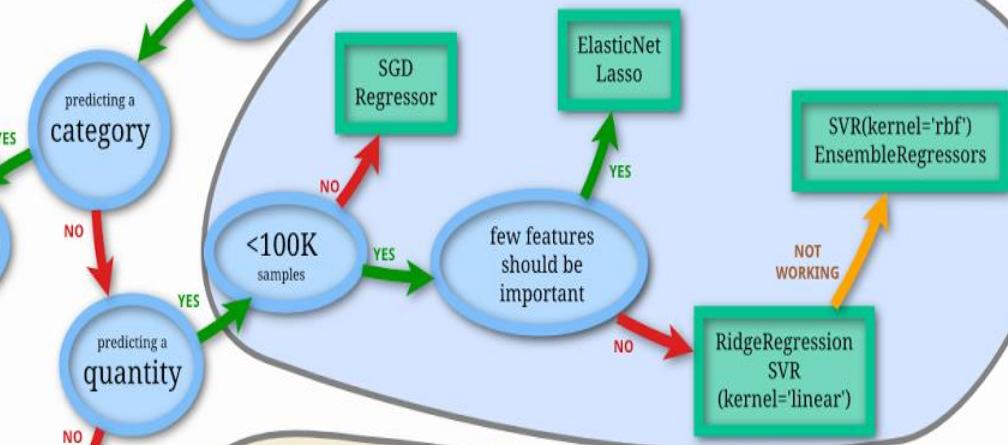
clustering



scikit-learn
algorithm cheat-sheet



regression

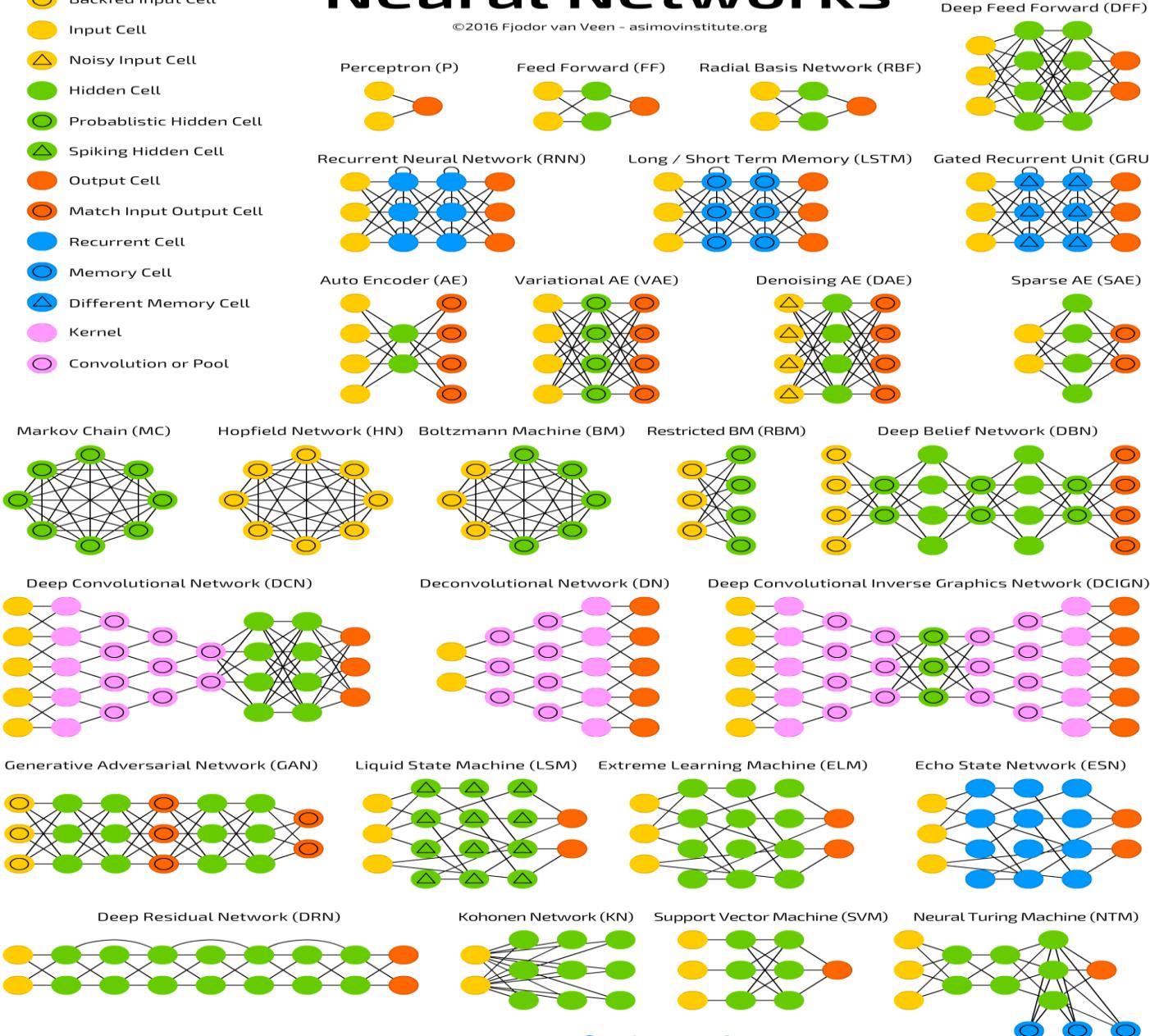


dimensionality
reduction

A mostly complete chart of
Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool



[Stefan Kojouharov](#)

Potential mathematics:
 Persistent homology
 Topological complexity
 Poset
 Directed graph
 Spectral graph
 Hypergraph
 Algebraic geometry
 K-theory?
 Exterior algebra?
 Homotopy
 Matroid
 Differential geometry

Computational issues:
 Linear algebra
 Matrix analysis
 Sophistical analysis
 Gradient descend
 ...

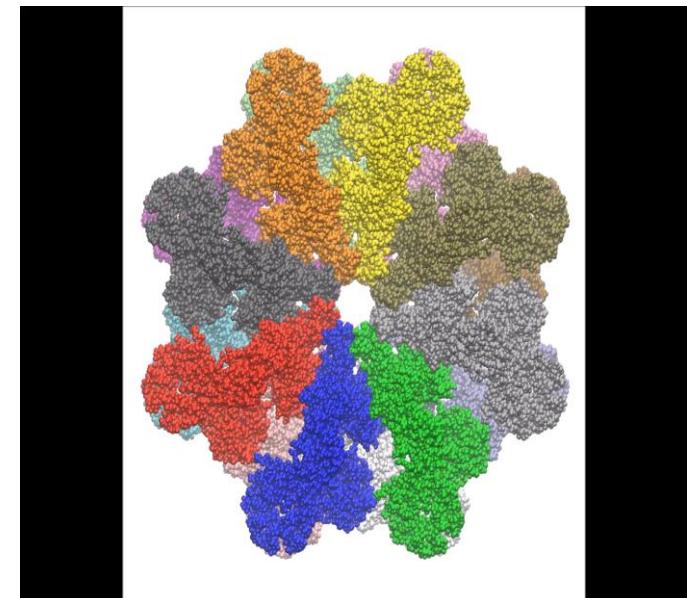
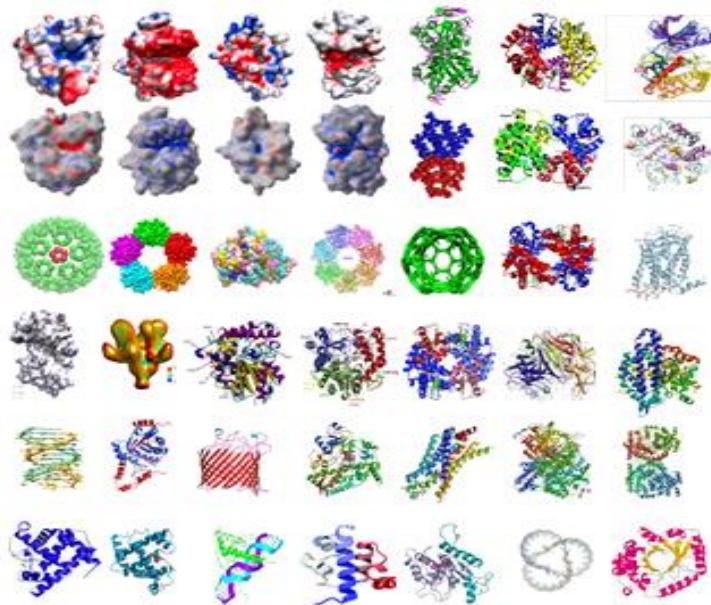
How to do deep learning for 3D biomolecular data?

Obstacles for deep learning of 3D biomolecules:

- **Geometric dimensionality:** \mathbb{R}^{3N} , where $N \sim 5000$ 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

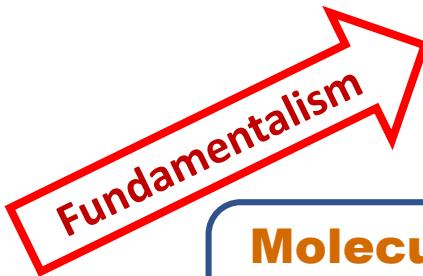
Solution:

- **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



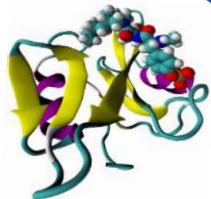
Molecular Mechanics
 \mathbb{R}^{3N}

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

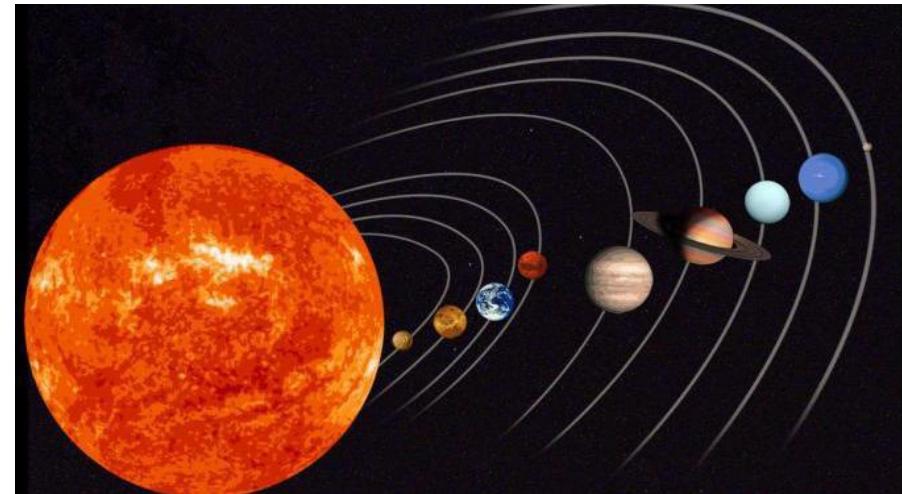
Quantum Mechanics
 \mathbb{R}^{3Nn+3N}

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

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



Differentiable Manifold
 \mathbb{R}^2

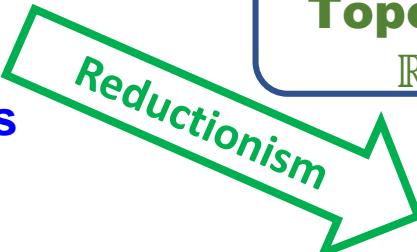


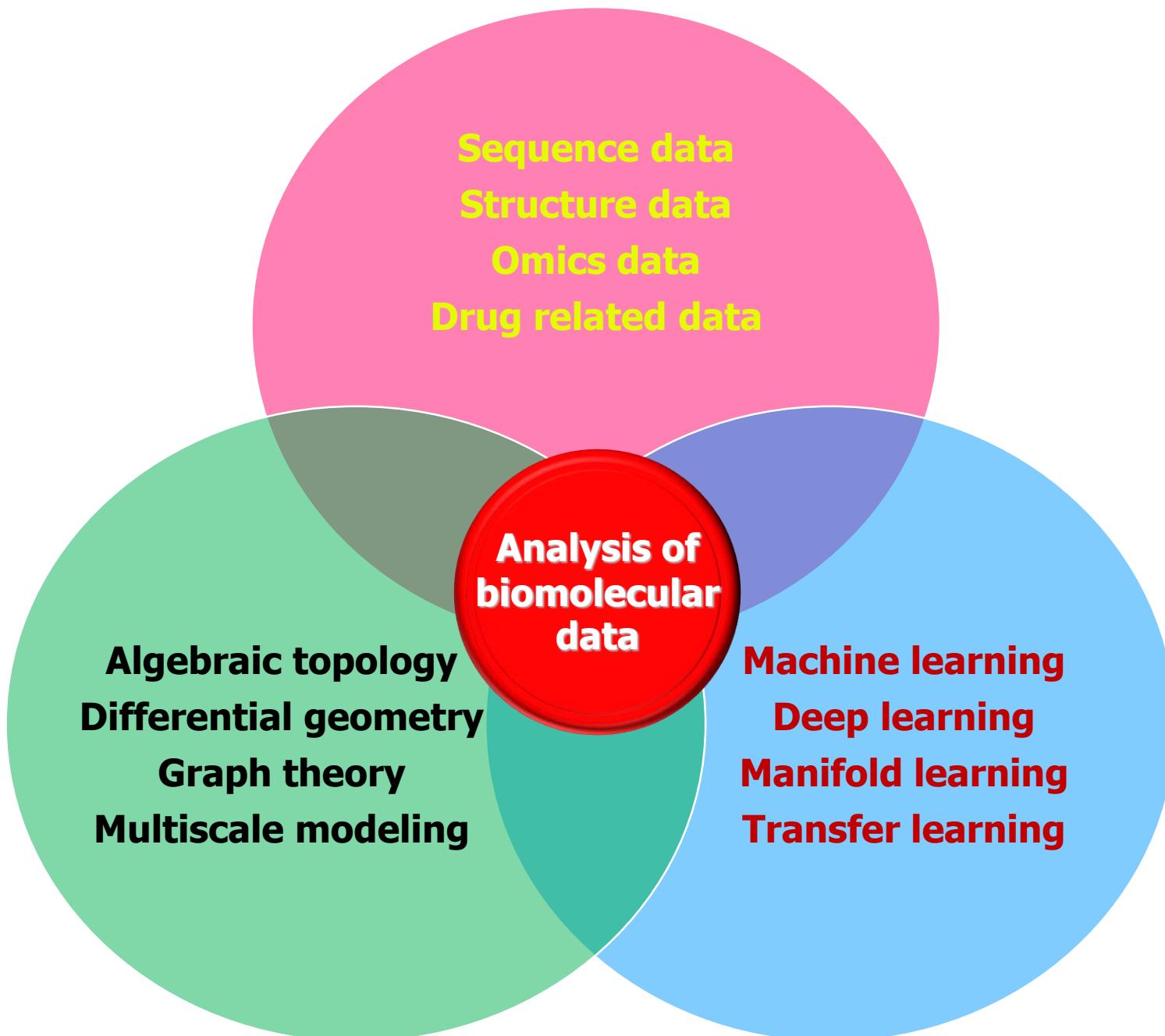
Algebraic Topology
 \mathbb{R}^1

Graph Theory
 \mathbb{R}^0

Geo-Top Indices
 \mathbb{R}^0

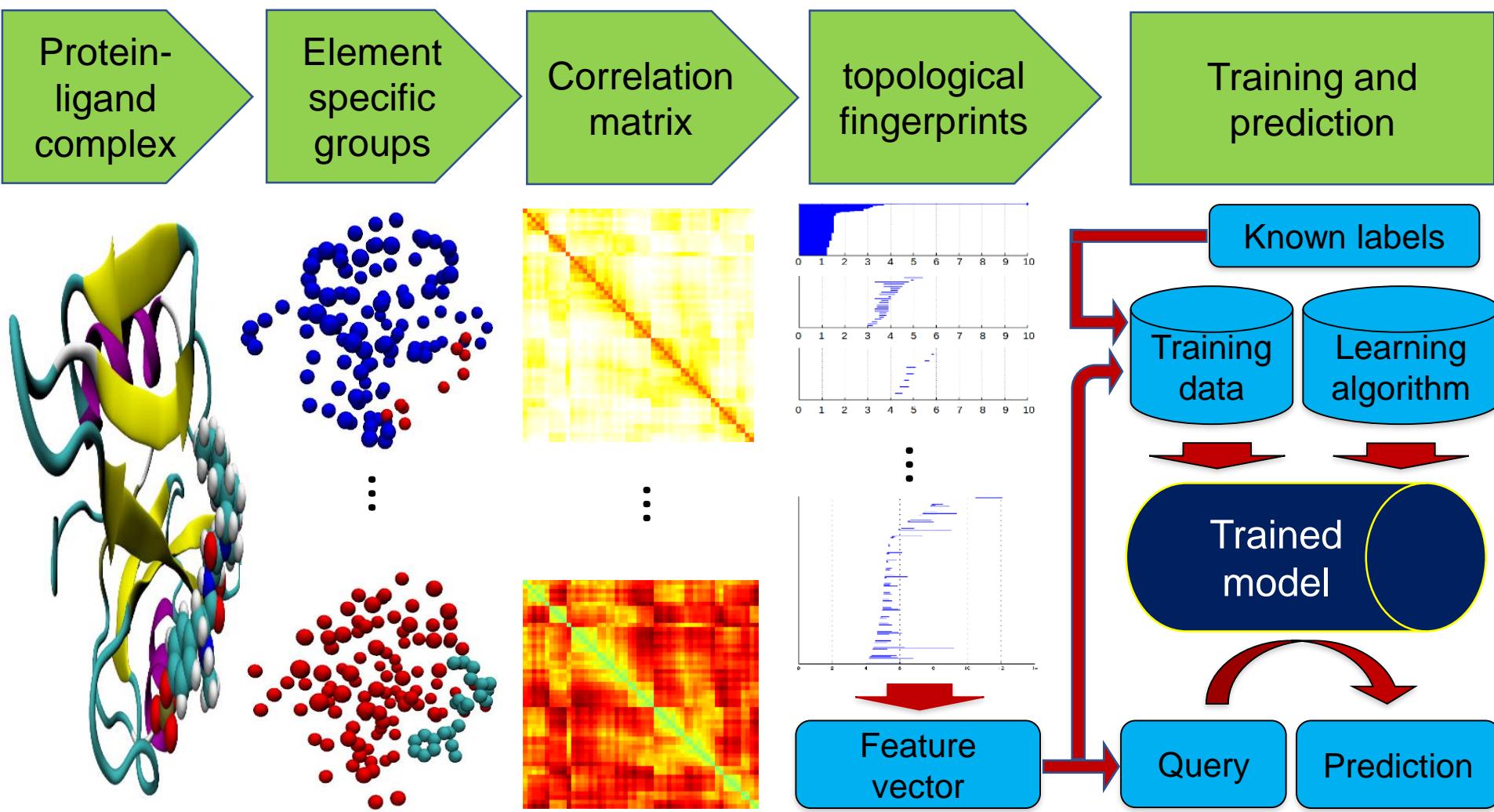
Basic hypothesis:
Intrinsic physics lies on low-dimensional manifolds in a high dimensional space



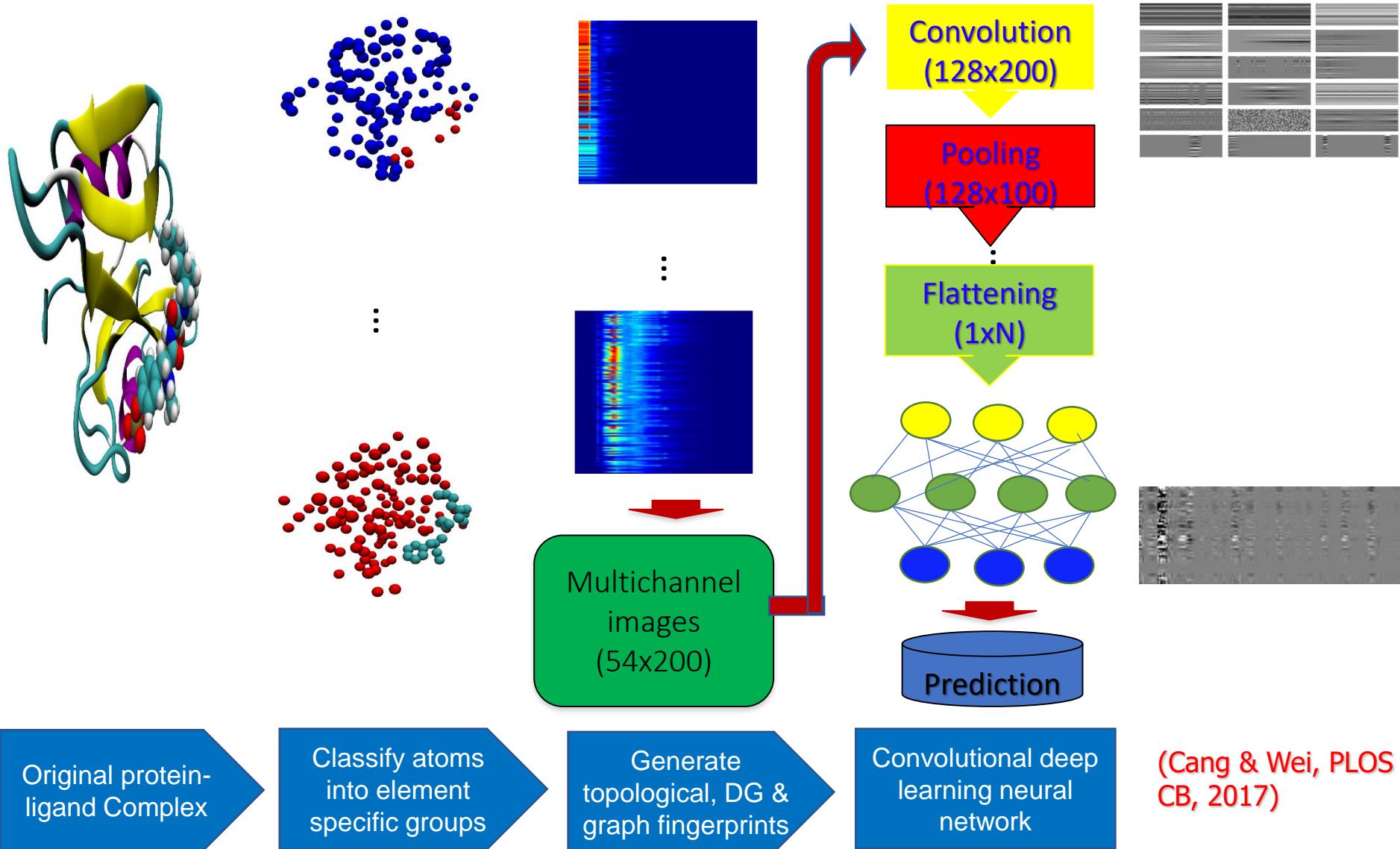


Topology based learning architecture

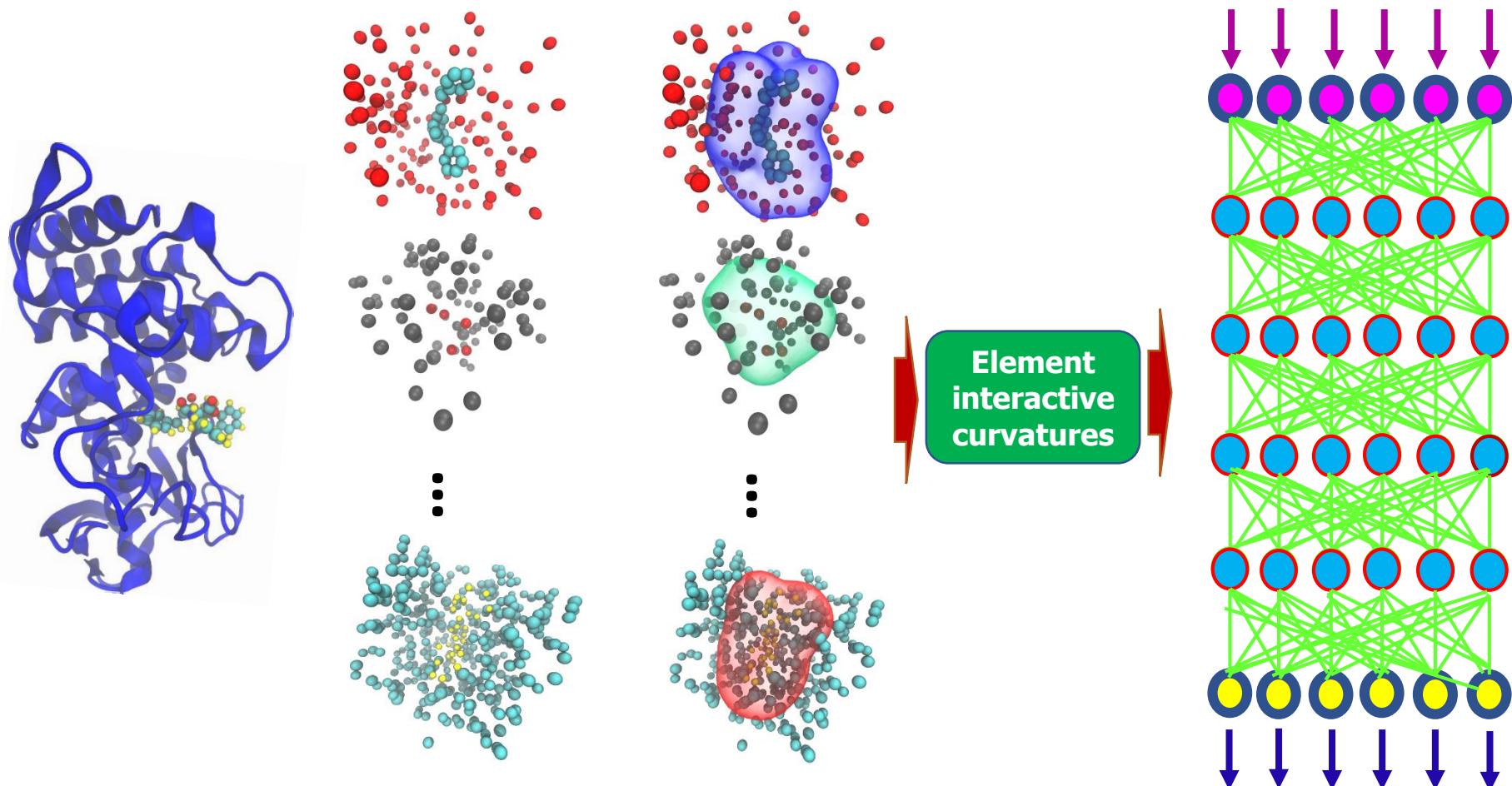
(Cang & Wei, IJNMBE, 2017)



Topological convolutional deep learning architecture



Differential geometry based deep learning



Protein-ligand complex

Element specific groups

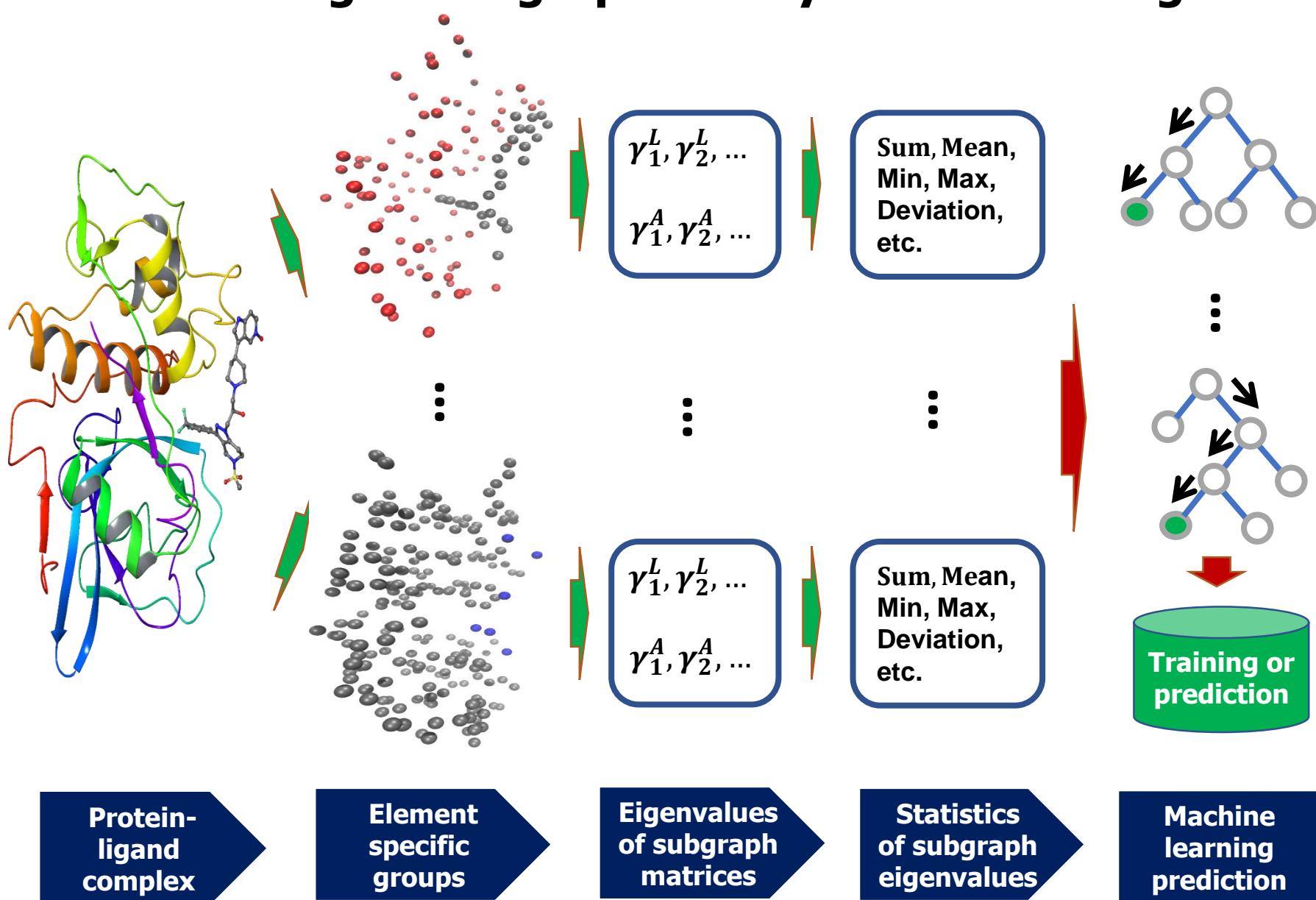
Element interactive manifolds

Differential geometry features

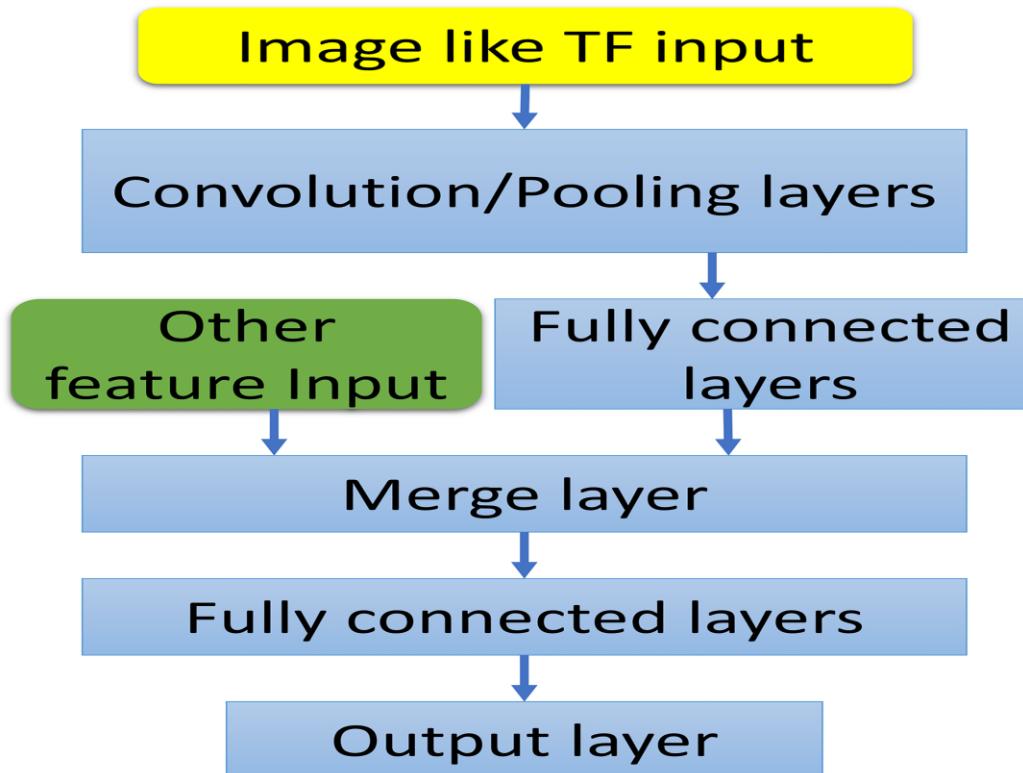
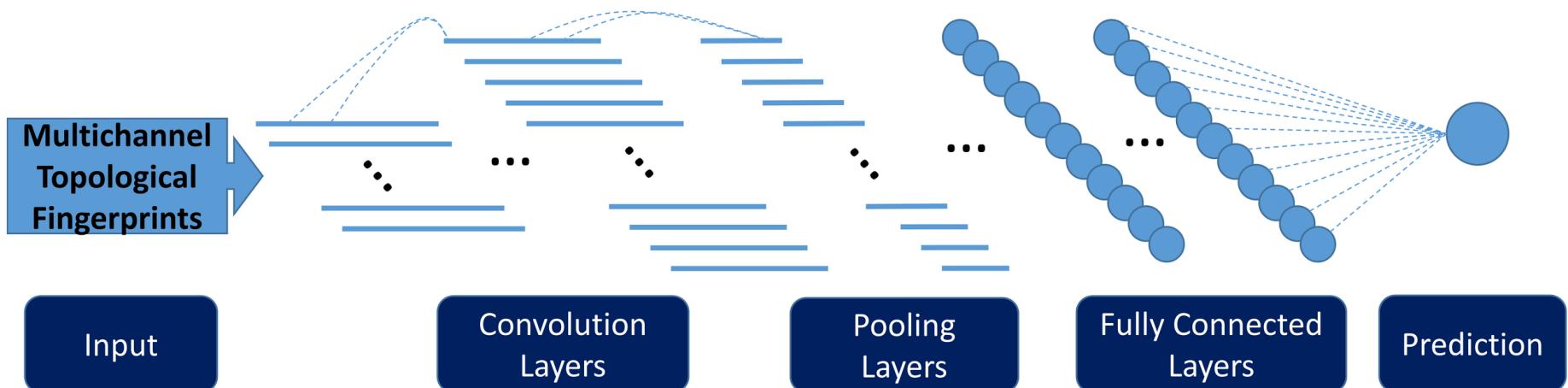
Machine learning prediction

(Nguyen & Wei, 2018)

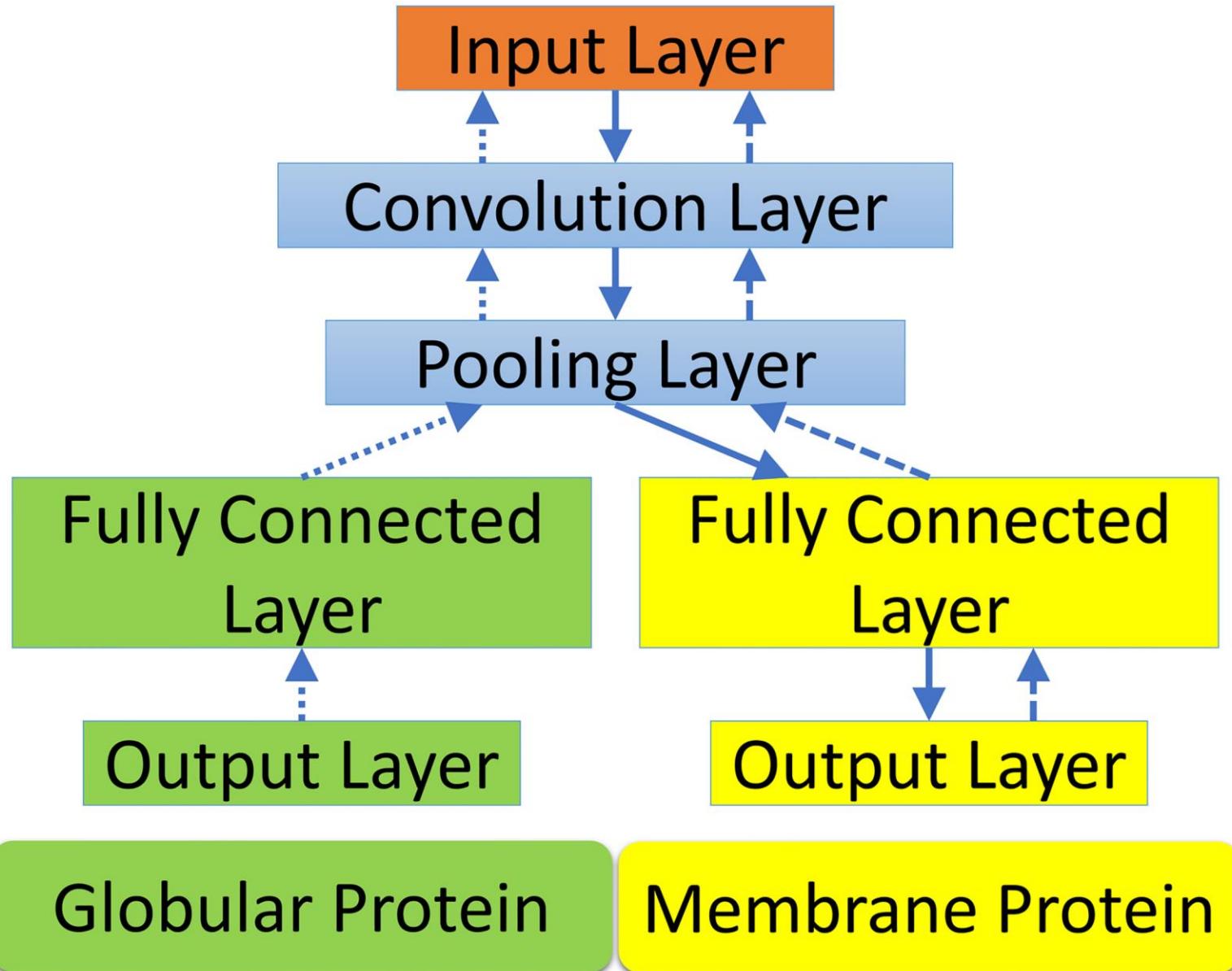
Algebraic graph theory based learning



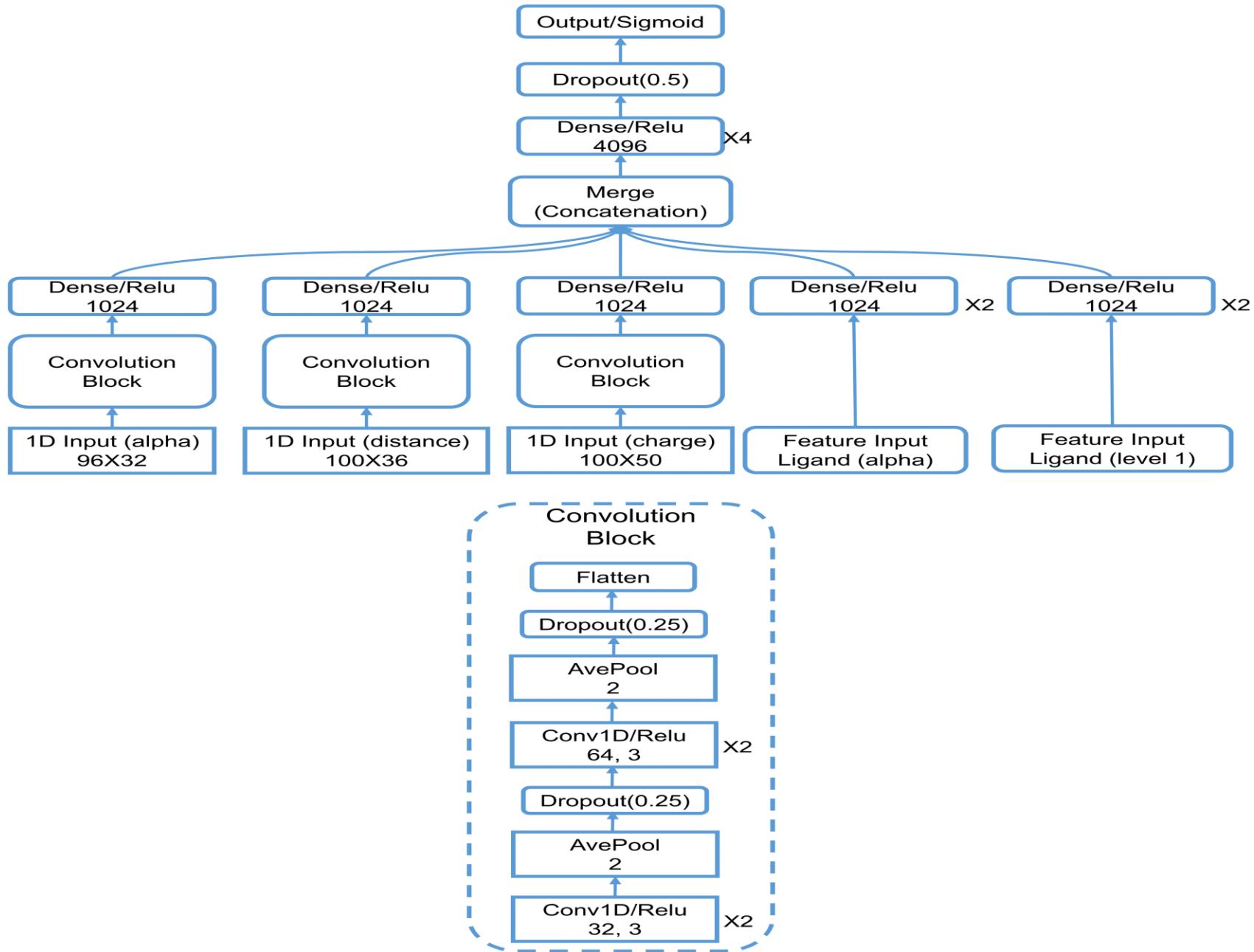
1D convolutional neural networks for mixed features



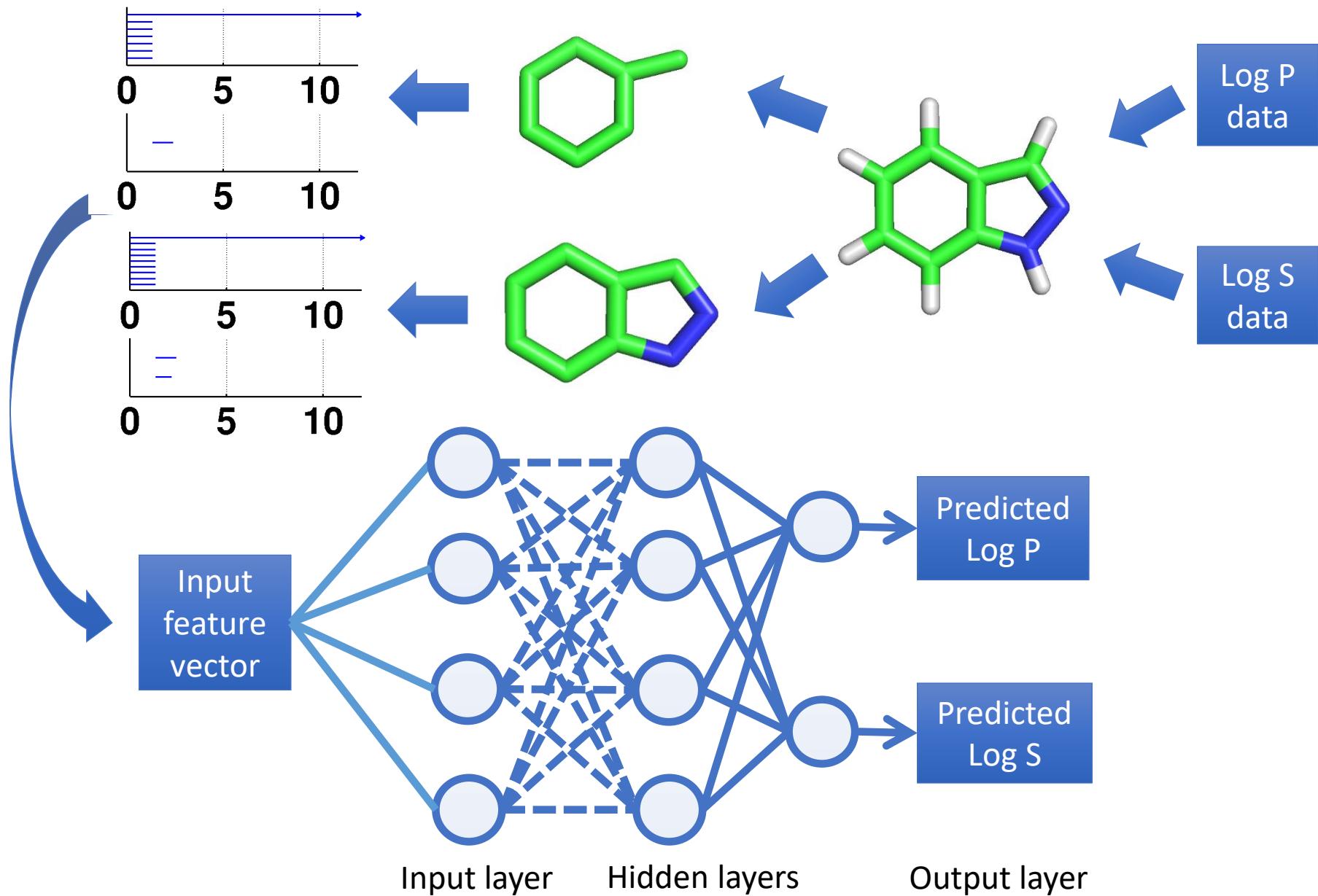
Multitask convolutional neural networks



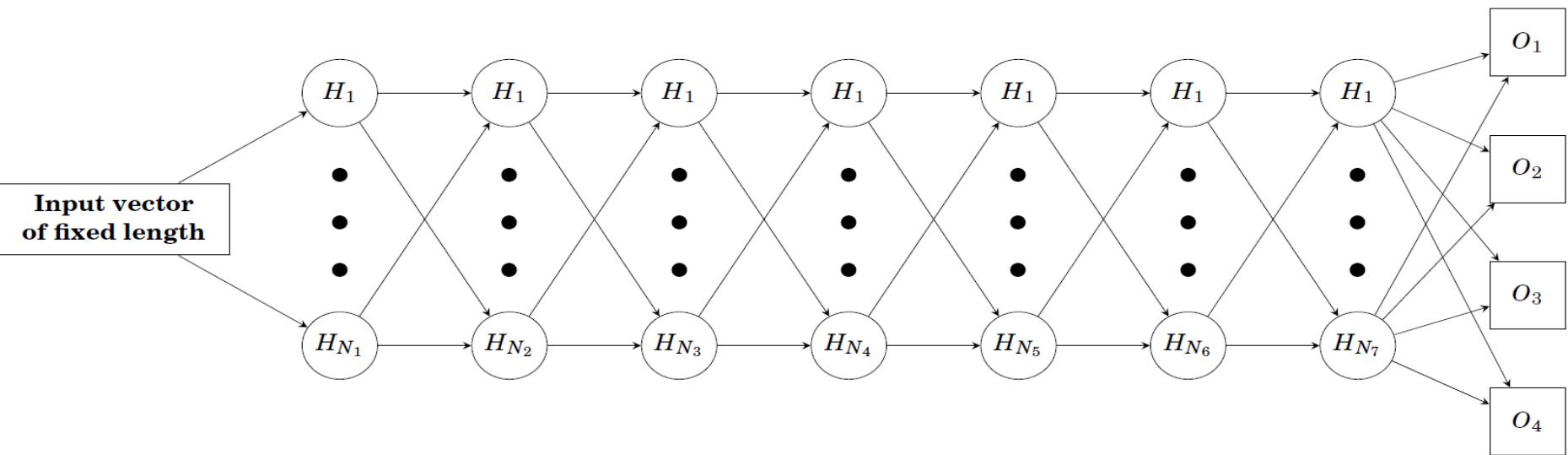
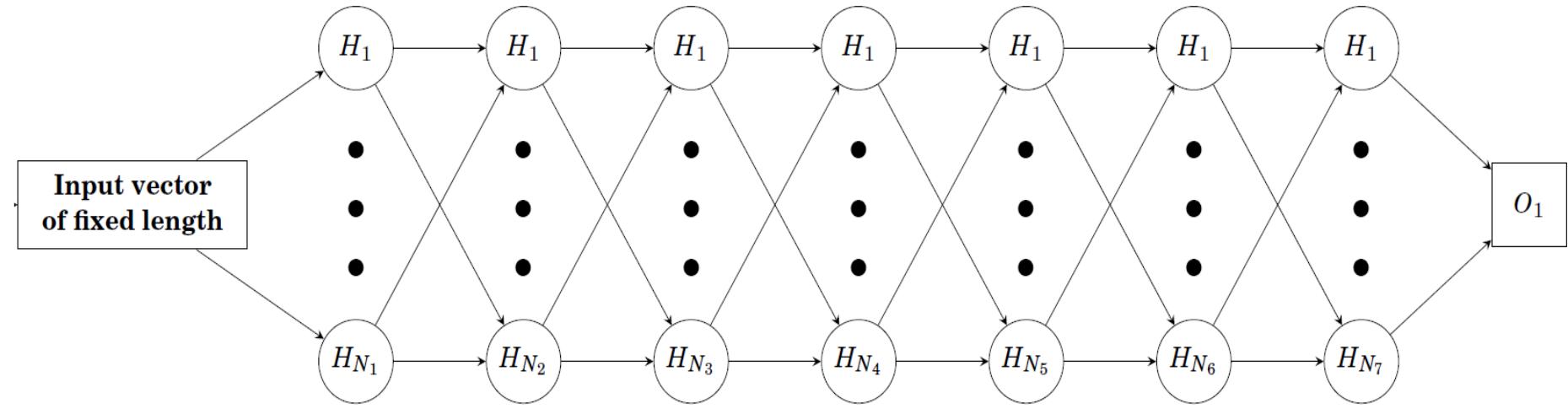
Convolutional neural networks for mixed features



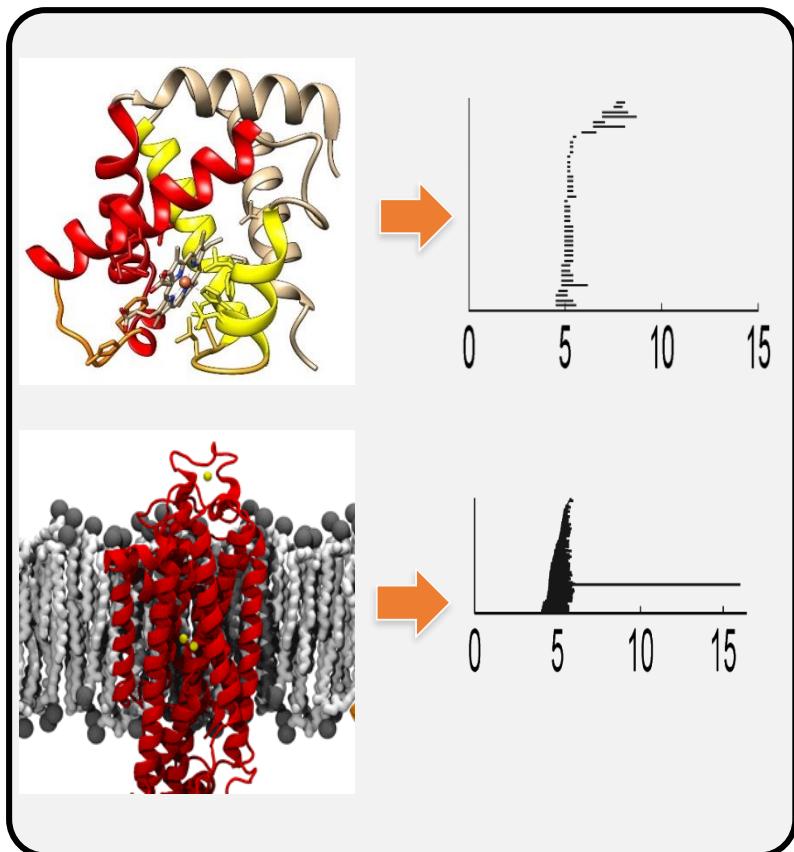
Artificial neural networks for multi-task solubility and partition coefficient predictions



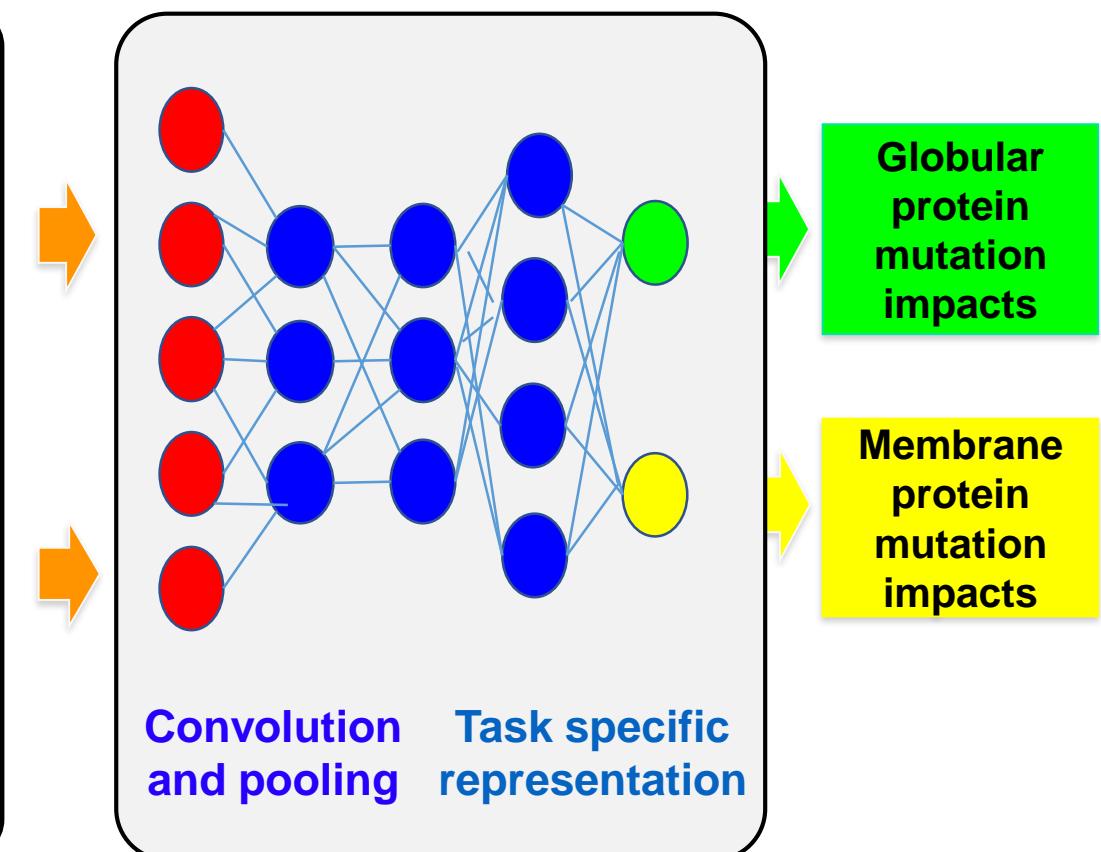
Artificial neural networks for single-task and multi-task toxicity predictions



Topological Multi-Task Deep Learning



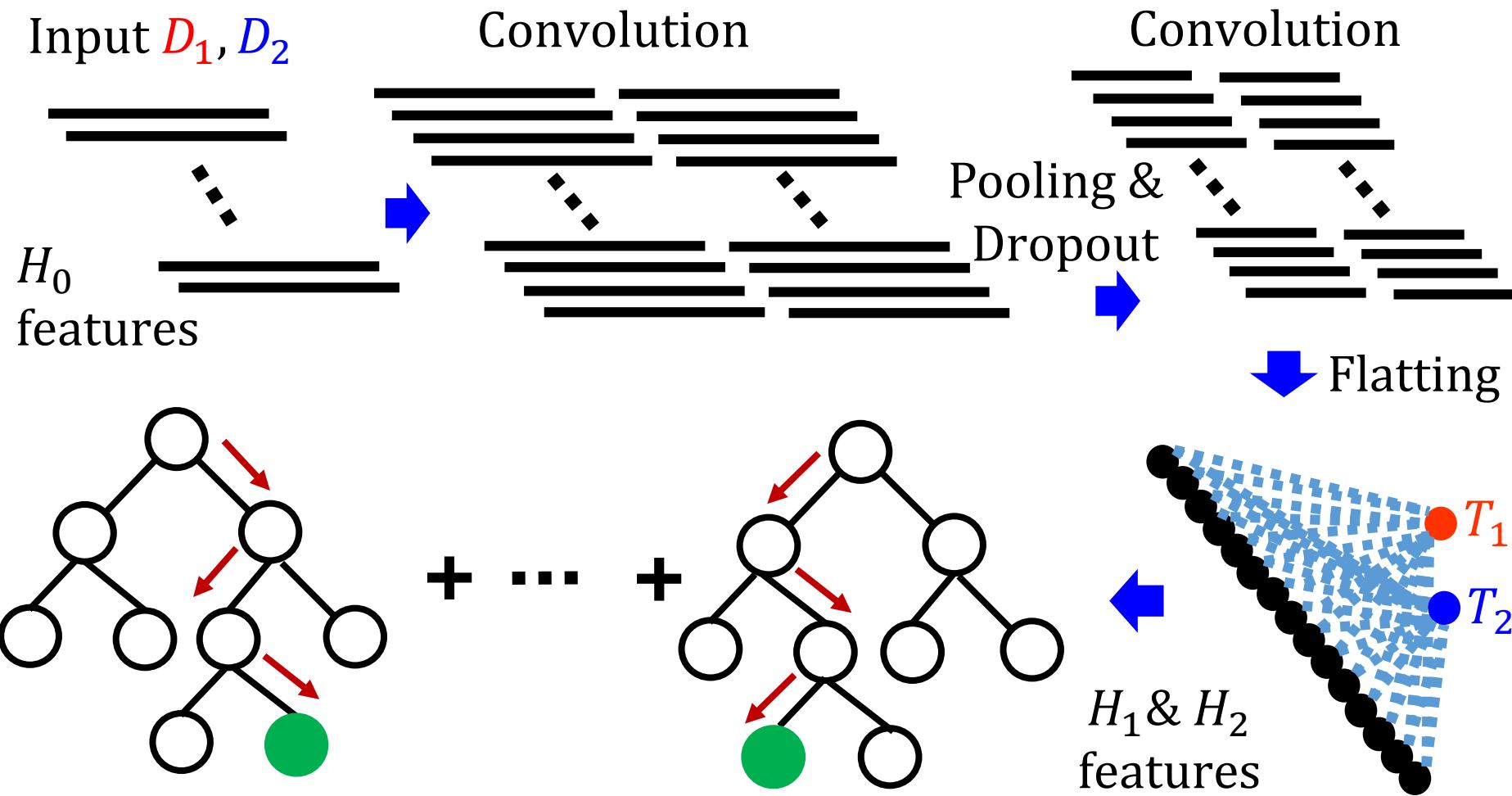
Topological feature extraction



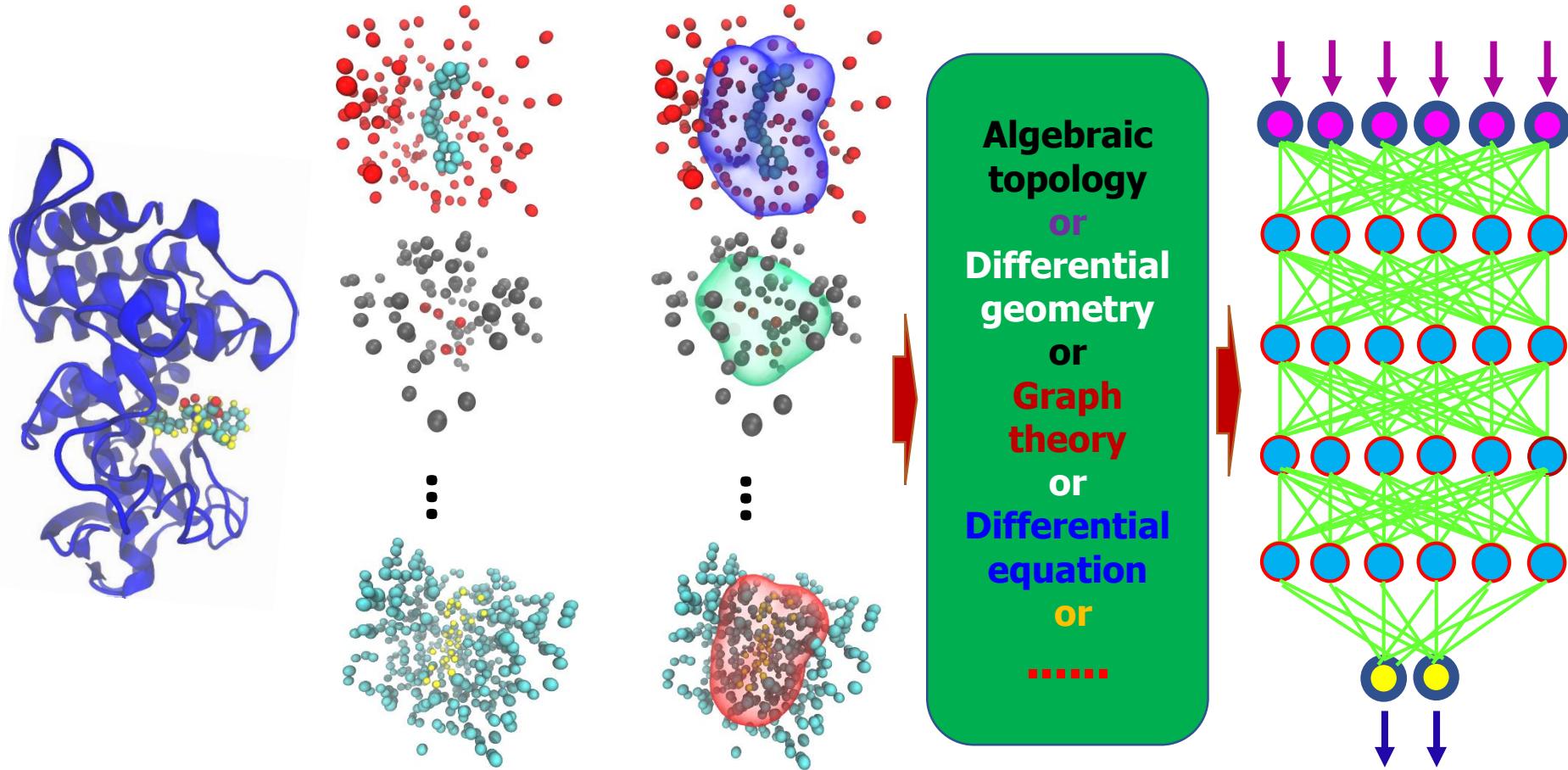
Multi-task topological deep learning

(Cang & Wei, PLOS CB, 2017)

Convolutional neural network assisted multitask gradient boosting trees



Mathematical deep learning



Protein-
ligand
complex

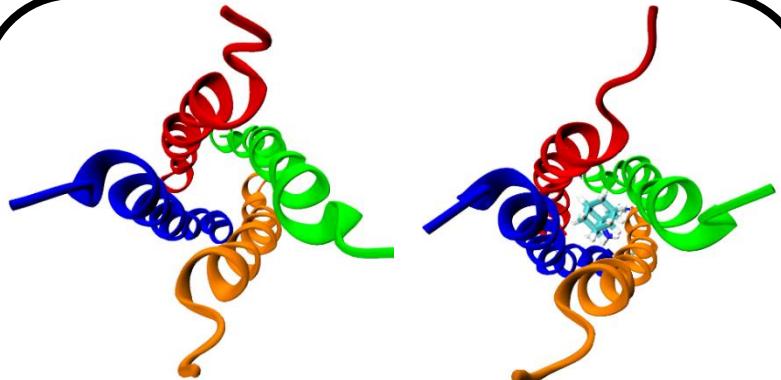
Element
specific
groups

Element
interactive
manifolds

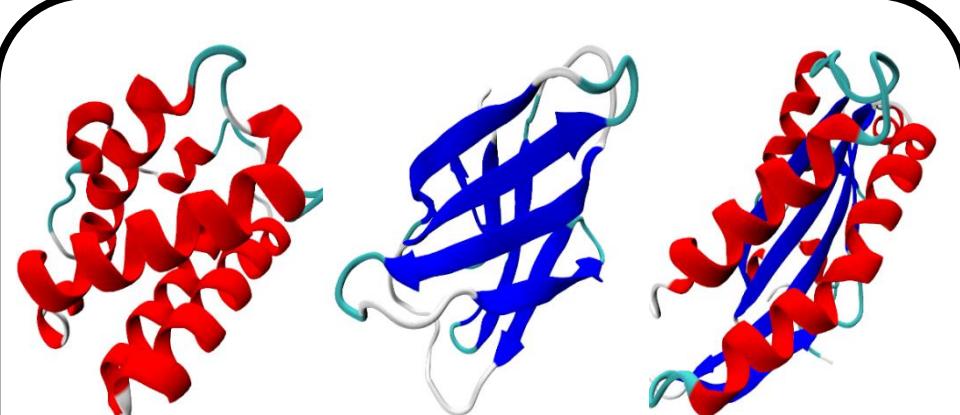
Various
Mathematical
features

Machine
learning
prediction

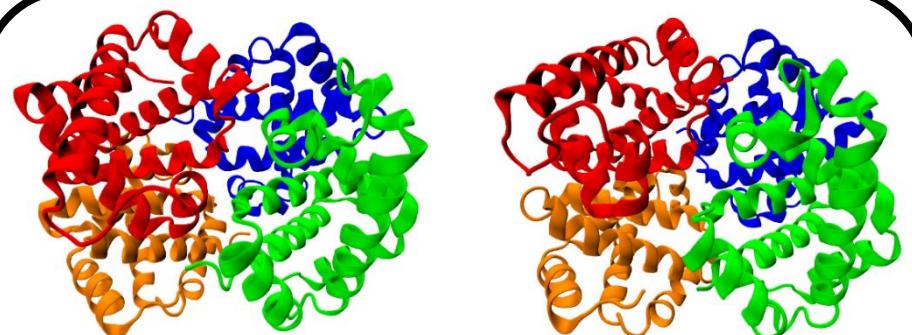
Topological fingerprint based machine learning (SVM) for the classification of 2400 proteins



Influenza A virus drug inhibition: 96% Accuracy



Protein domains: 85% Accuracy
(Alzheimer's disease)

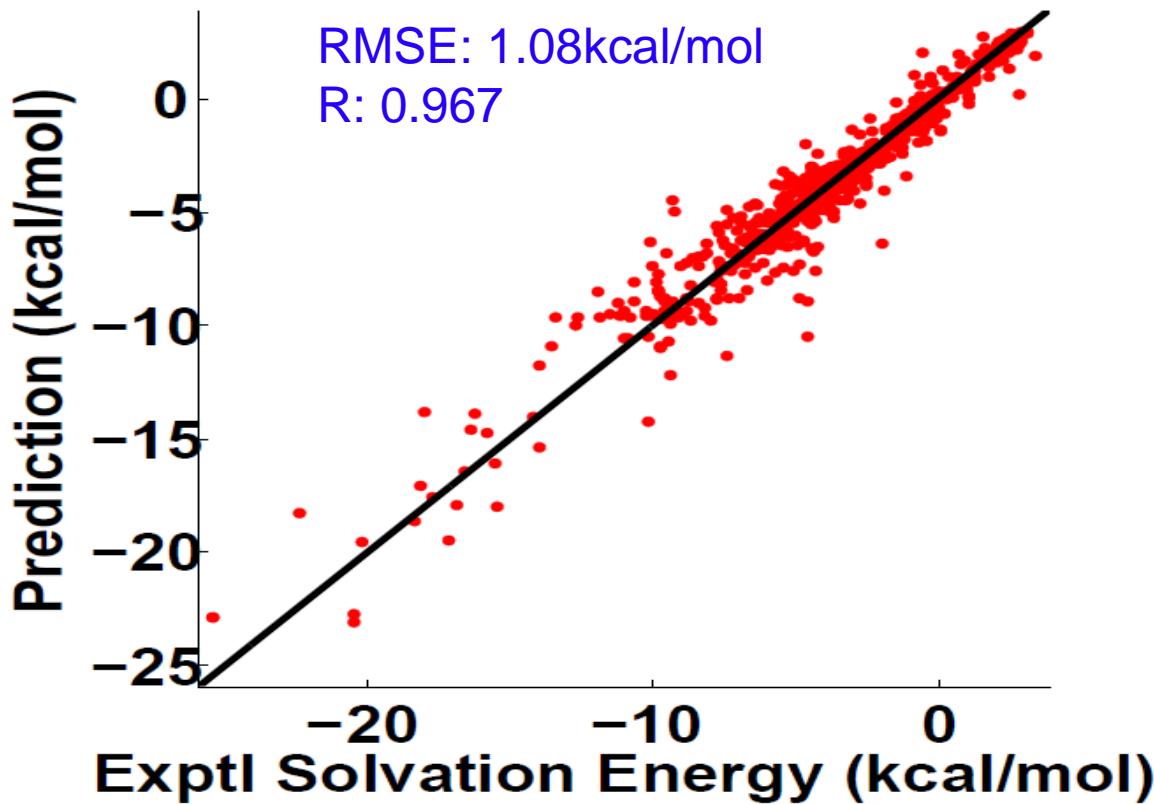


Hemoglobins in their relaxed and taut forms: 80% accuracy

(Cang et al, MBMB, 2015)

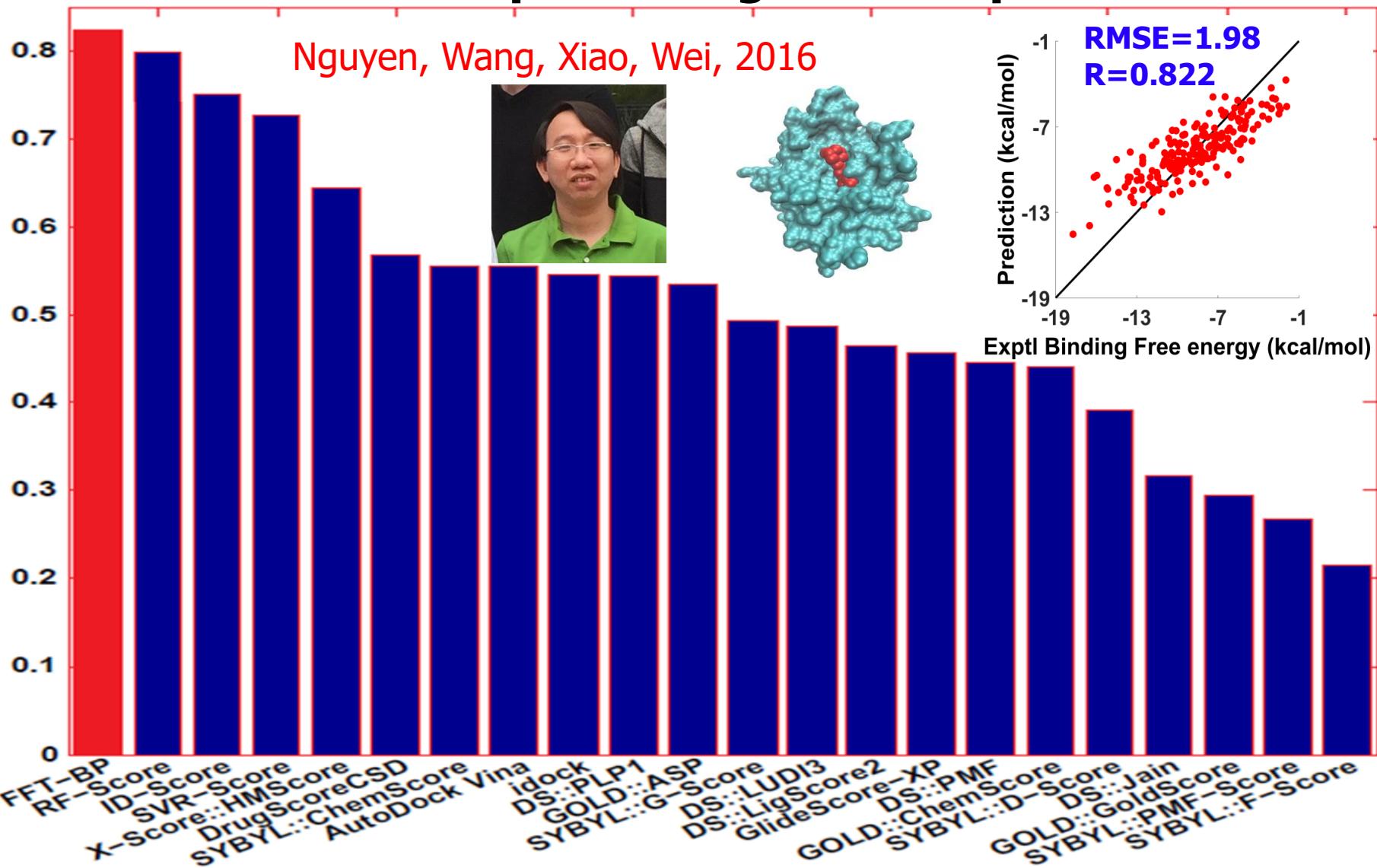
55 classification tasks of protein superfamilies over 1357 proteins from Protein Classification Benchmark Collection: 82% accuracy

Leave-one-out prediction for solvation free energies of 668 molecules

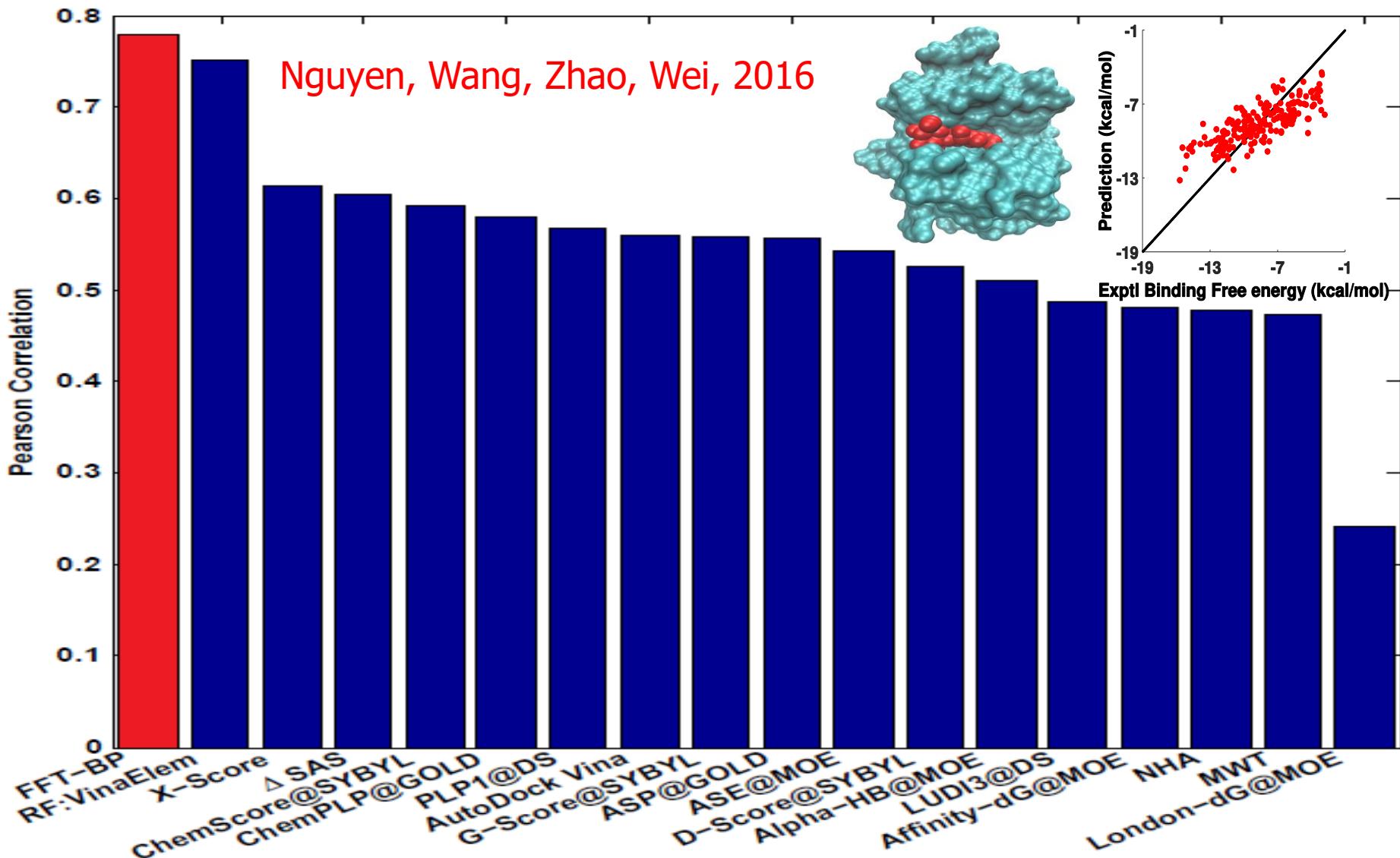


Wang, Zhao, Wei, JCP 2016

Blind binding affinity prediction of PDDBBind v2007 core set of 195 protein-ligand complexes

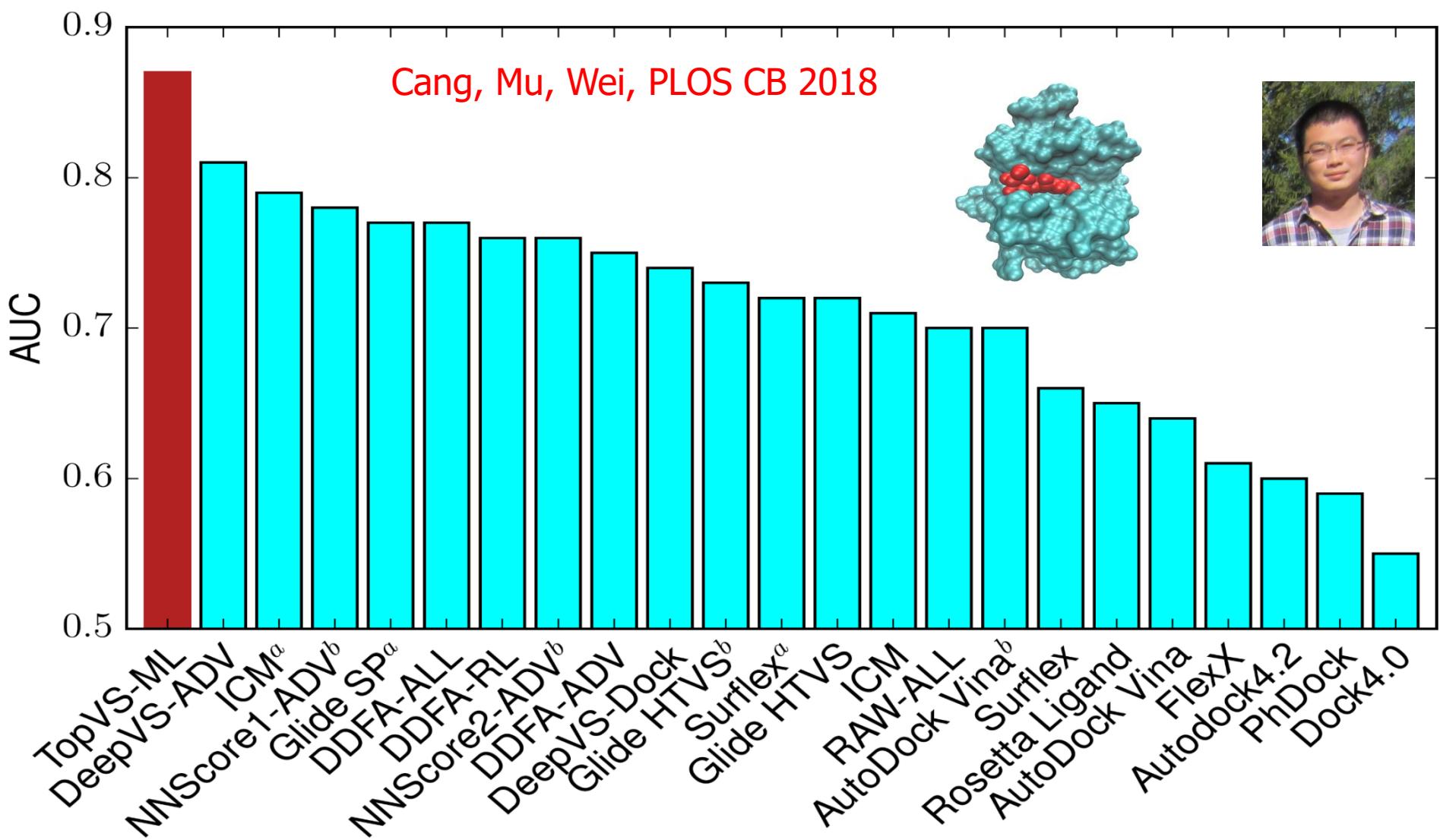


Blind binding affinity prediction of PDBBind v2013 core set of 195 protein-ligand complexes

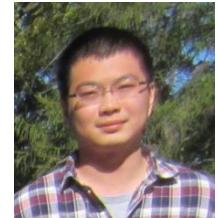
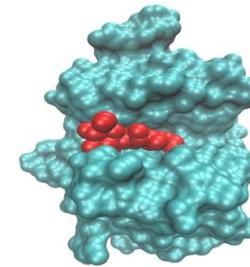


Directory of Useful Decoy (DUD)

Classification of 98266 compounds containing 95316 decoys and 2950 active ligands binding to 40 targets from six families

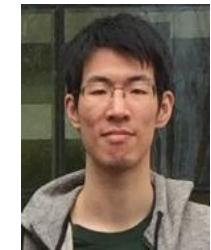
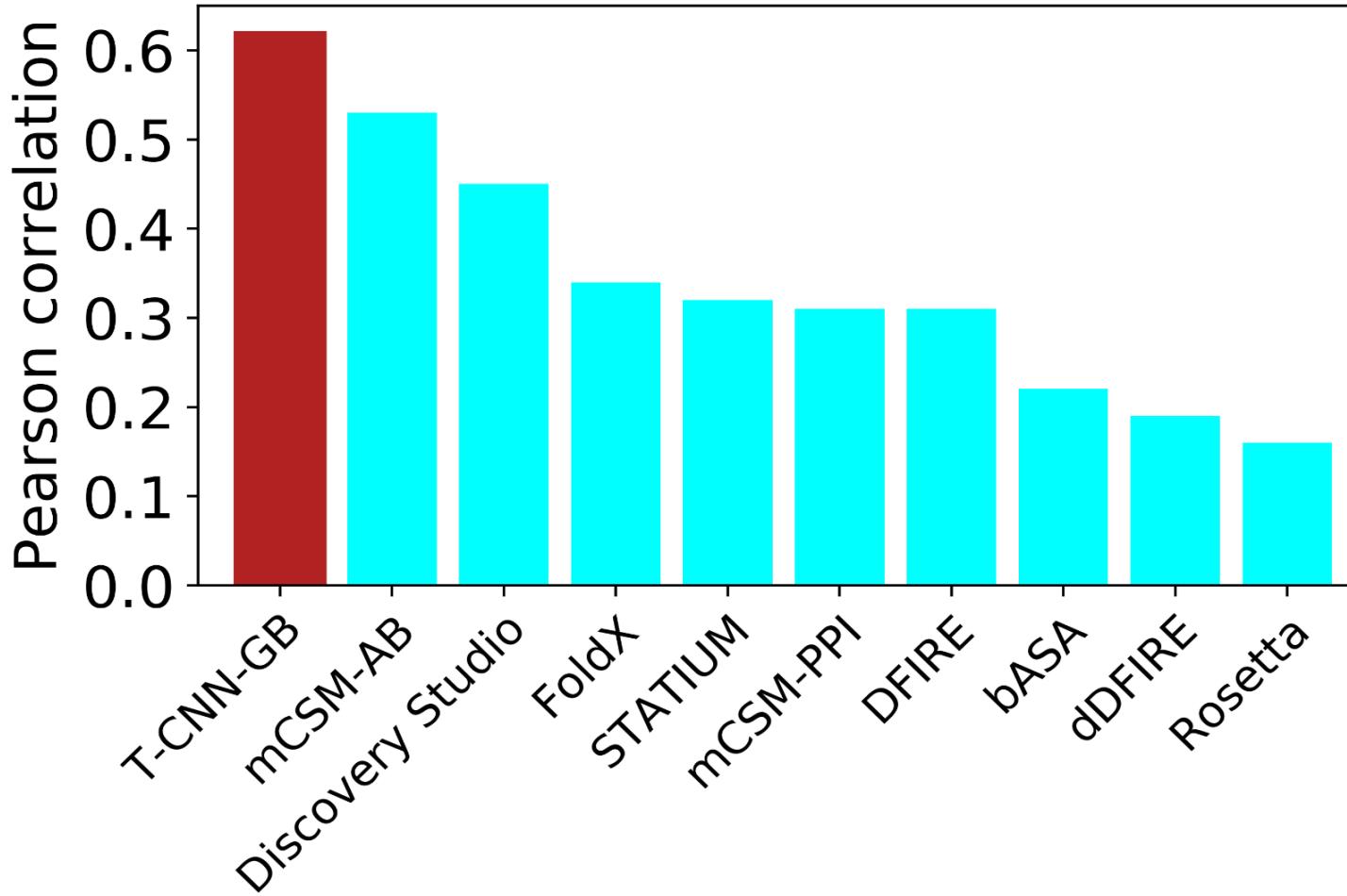


Cang, Mu, Wei, PLOS CB 2018



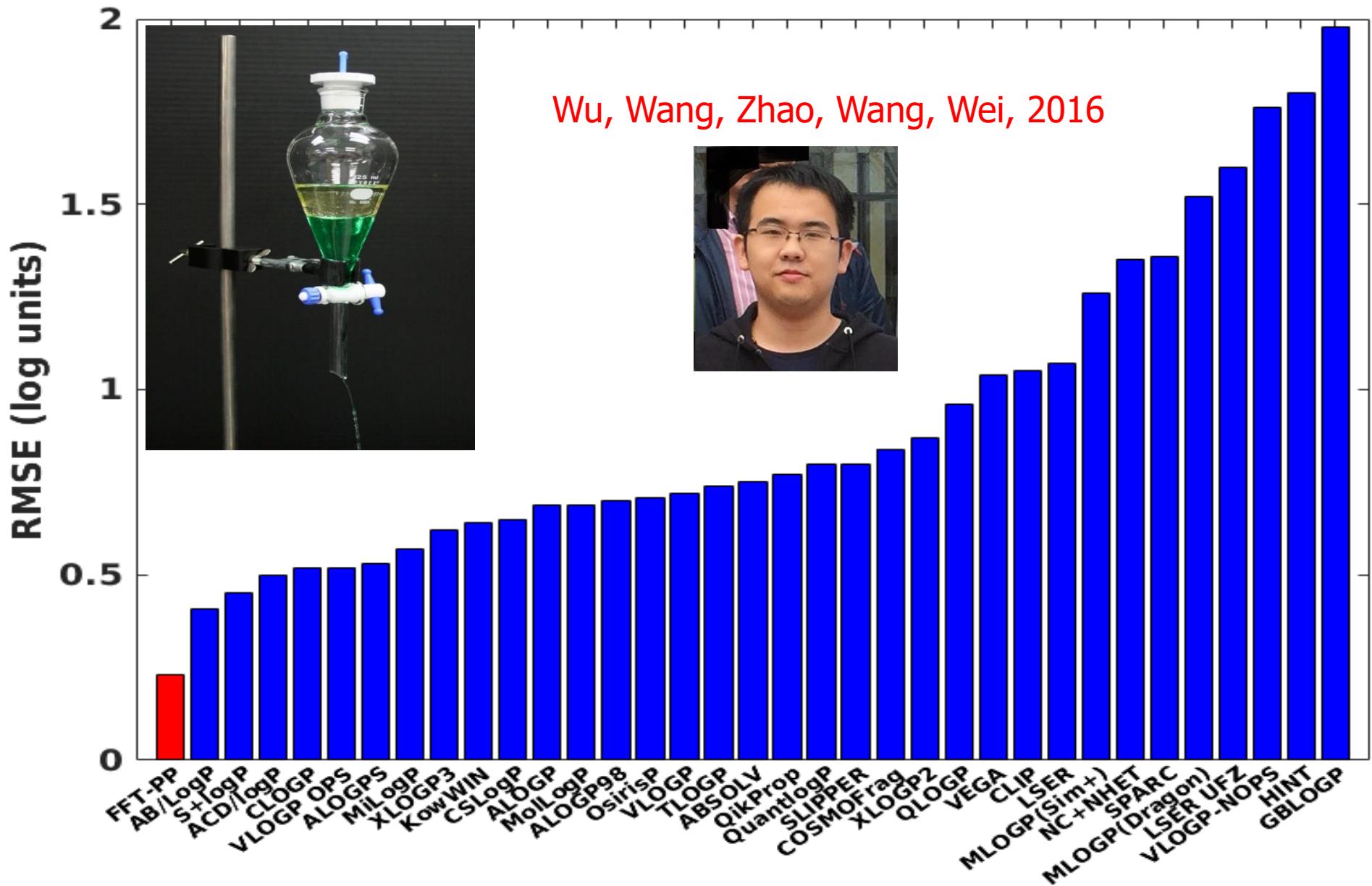
Antibody-antigen binding free energy changes upon mutation (AB-Bind dataset)

Wang, Cang, Wei, 2019

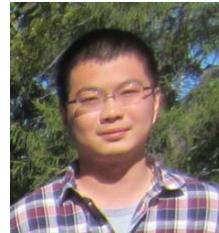
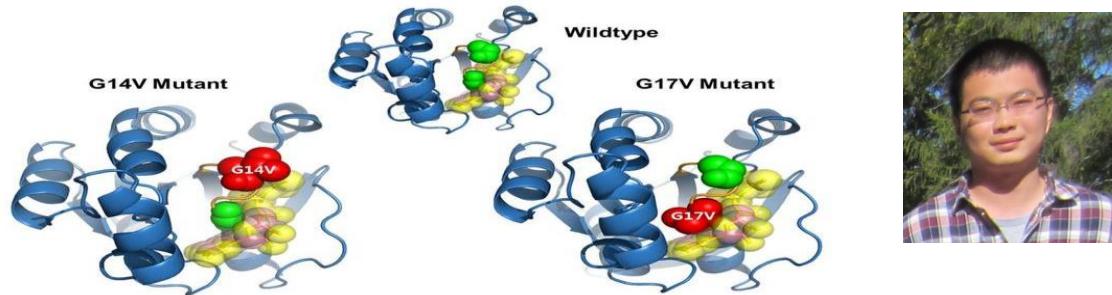


Menglun Wang

Prediction of partition coefficients: Star Set (223) molecules)

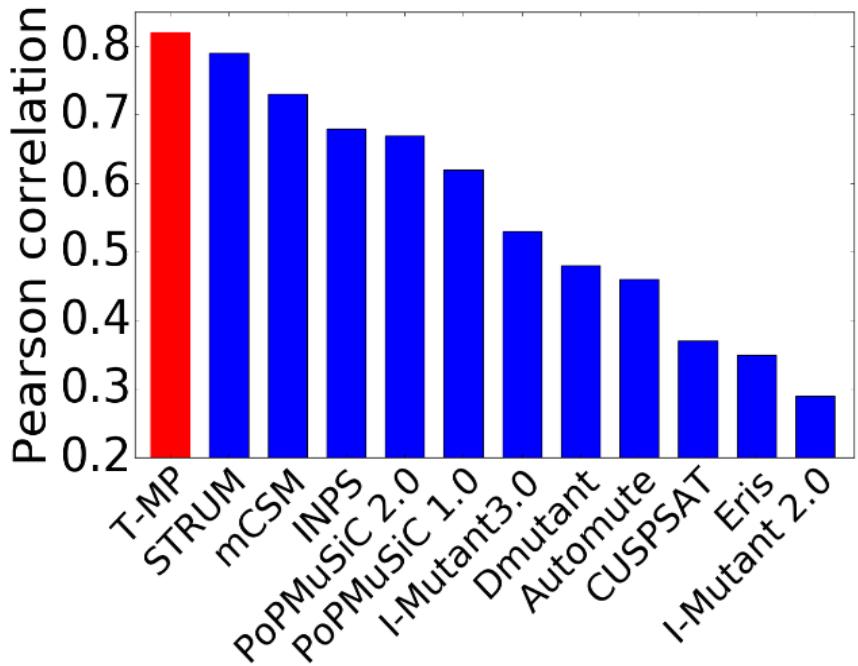


Blind prediction of mutation energies

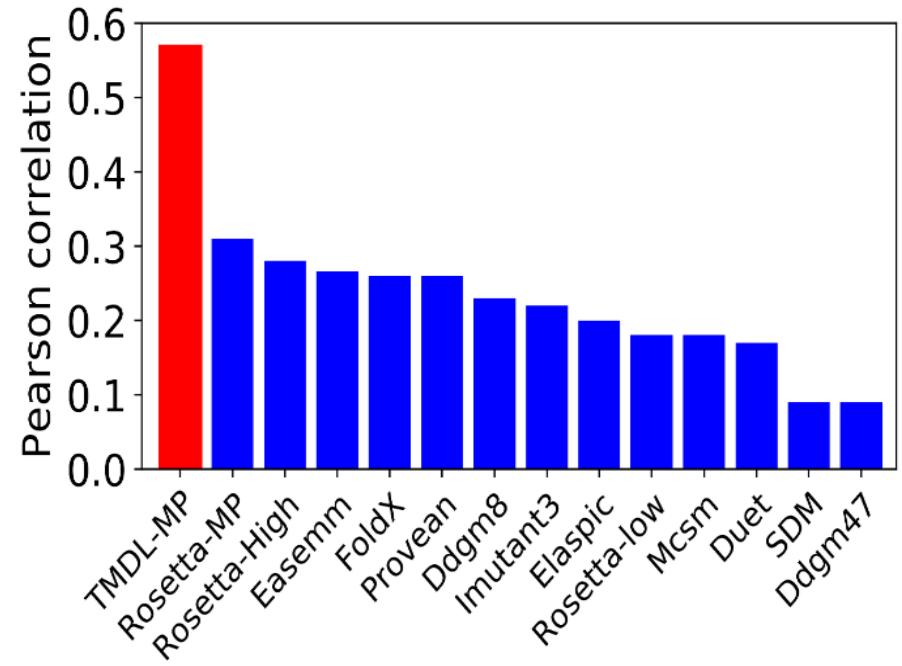


Cang, Wei, Bioinformatics 2017

Predicting mutations on
2648 globular proteins



Predicting mutations on 223 membrane proteins



Further topics and future directions

- Omics data analysis: proteomics, metabolomics, metagenomics, lipidomics, glycomics, transcriptomics, epigenomics, etc.
- Integration of molecular data and genomic data.
- Data fusion from molecule, subcellular organelle, cell, tissue, organism, to cross species.
- Using mathematics (geometry, topology, algebra, analysis, PDE, numerical analysis, etc.) to analyze the foundation of machine learning and AI.
- Creating new mathematical tools to reduce the dimensionality of massive, diverse and complex biological data.
- Developing and application of new machine learning algorithms for biomolecular data analysis and prediction.
- Gromov-Hausdorff and Gromov-Wasserstein distances for biomolecular data.
- Symplectic geometry based gradient descent.
- Predictability analysis.
- Biology inspired new mathematics to better understand the rule of life.



thank you