# Mathematics is the champion of biomolecular data challenges

**Zixuan Cang, Duc Nguyen, Bao Wang, Chengzhang Wang, Menglun Wang, Guo-Wei Wei, Kedi Wu, and Zhixiong Zhao**
**Department of Mathematics, Michigan State University**

## Feature functional theory (FFT)

**Goal:** To prediction microscopic and macroscopic relationships in biomolecular data

**Basic assumptions:**

- **Representability assumption**: there exists a microscopic feature vector that can uniquely characterize, and distinguish one molecule from another

$$\mathbf{v}_i = \left(\mathbf{x}_i; \mathbf{o}_i\right) = \left(x_{i1}, x_{i2},..., x_{in}; o_{i1}, o_{i2},...., o_{il}\right)$$

microscopic features; macroscopic features

- **Similarity assumption**: molecules with similar microscopic features have similar macroscopic features.

- **Feature-function relationship assumption**: the macroscopic features, i.e., solvation and binding free energies, of molecule A are functionals of microscopic feature vectors:

$$\Delta G_{\mathrm{A}} = f_{\mathrm{A}}\left(\mathbf{x}_{\mathrm{A}}, \mathbf{v}_1, \mathbf{v}_2,..., \mathbf{v}_n\right)$$

## Feature functional theory (FFT)
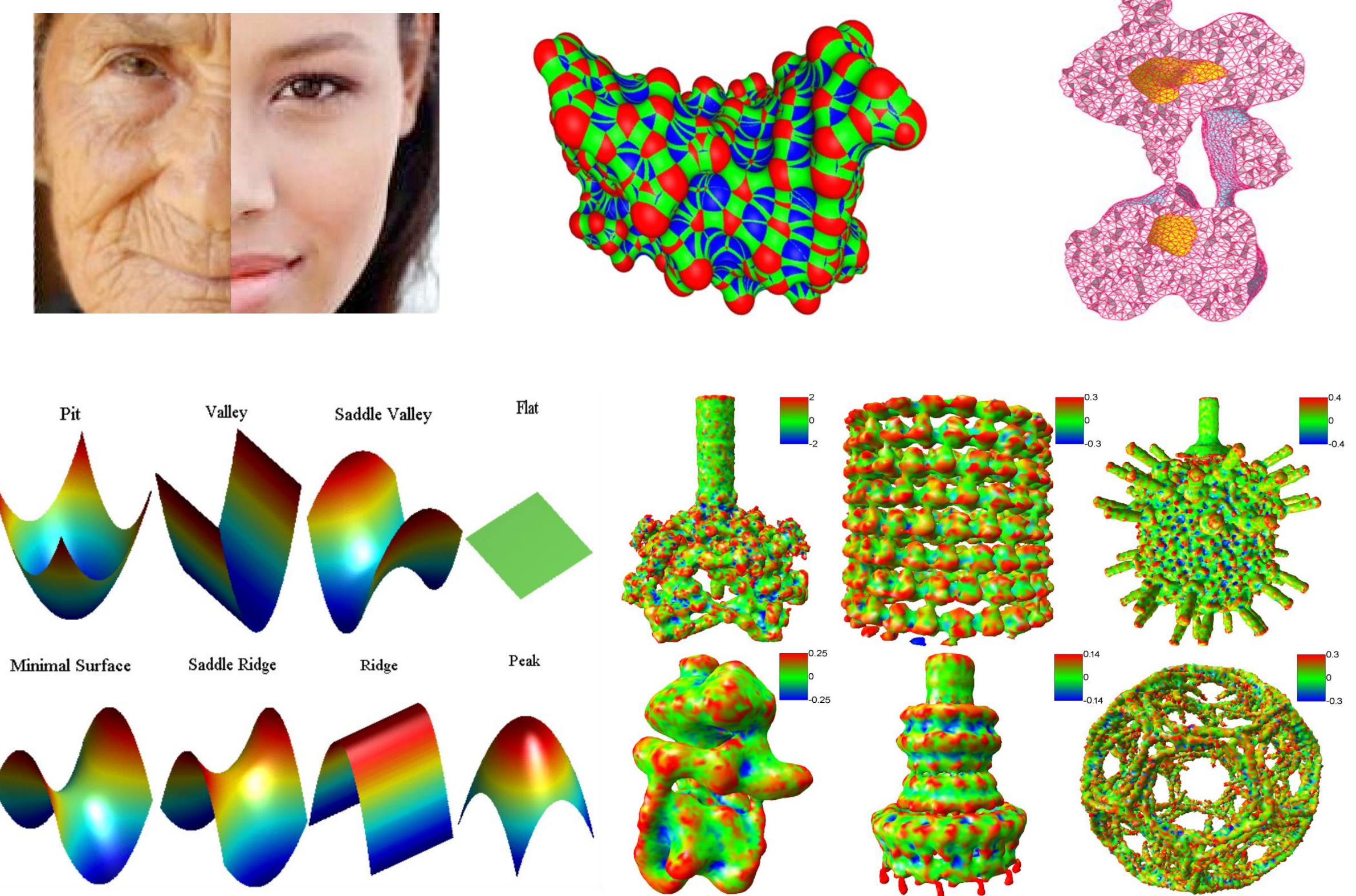
**Microscopic features:**

- Geometric: atomic surface areas, volume & curvatures
- Topological: Betti numbers
- Graph theory: discrete Laplacian, rigidity & flexibility
- Electrostatic: atomic charges, dipoles, quadrupoles & reaction filed energies (Poisson-Boltzmann equation)
- van der Waals: Lennard-Jones potentials

**Macroscopic features:**

- Protein-ligand binding affinities
- Protein mutation energy changes (stability changes)
- Drug partition coefficients
- Drug solvation free energies
- Protein-DNA/RNA binding energies
- Protein-protein binding affinities
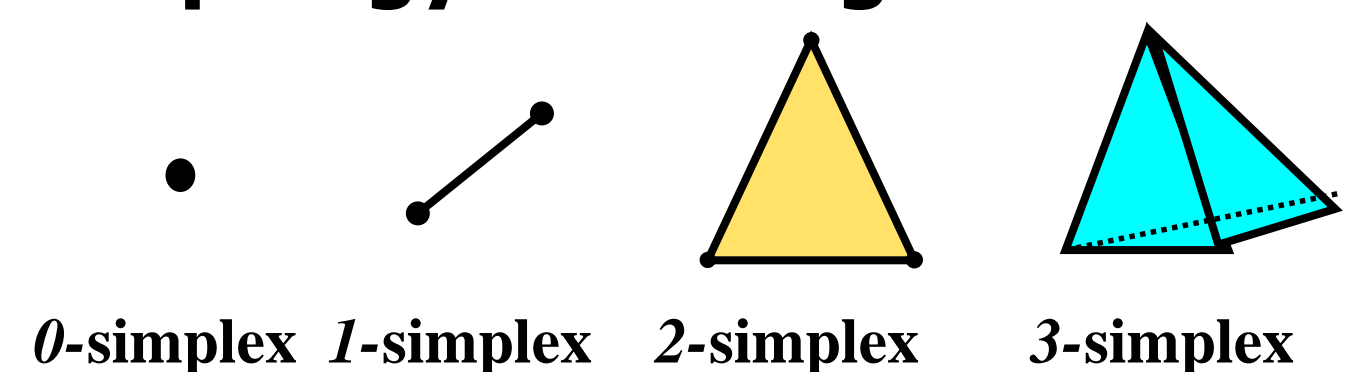
## Geometric modeling
### Besides surface area and volume, curvature matters



Feng, Xia, Tong and Wei, IJNMBI,2012; JCP, 2013

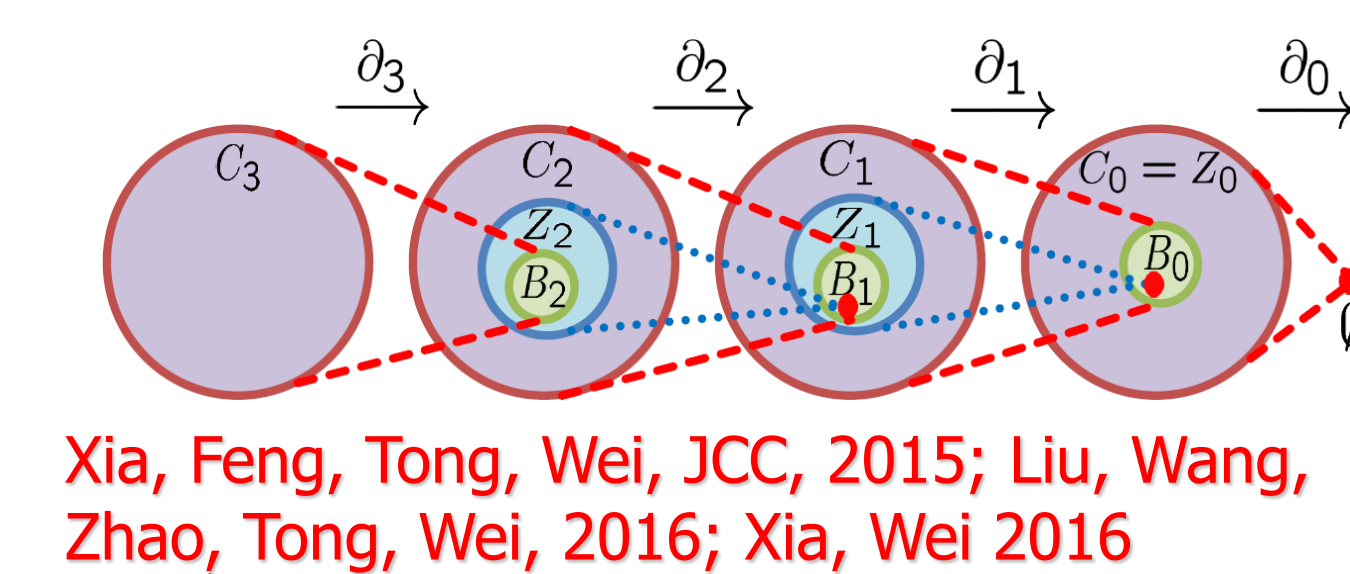## Persistent topology modeling

**Simplicial complex:**

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

**k-chain:** $\sum_i c_i \sigma_i^k$

**Chain group:** $C_k(K, Z_2)$

Xia, Feng, Tong, Wei, JCC, 2015; Liu, Wang, Zhao, Tong, Wei, 2016; Xia, Wei 2016

**Boundary operator:**

$$\partial_k \sigma^k = \sum_{i=0}^{k} (-1)^i \{v_0, v_1,..., \widehat{v_i},..., v_k\}$$

$Z_k = \mathrm{Ker}\, \partial_k$

$B_k = \mathrm{Im}\, \partial_{k+1}$

$H_k = \dfrac{Z_k}{B_k}$

$\beta_k = \mathrm{Rank}\left(H_k\right)$

2D persistence of a unfolding protein:

Quantitative topology:

$\beta_0$
$\beta_1$
$\beta_2$

## Graph theory modeling
### Weighted graph Laplacian, Flexibility rigidity index (FRI)
**FRI is about 20% more accurate than Gaussian network model (GNM) in B-factor prediction, based on 364 proteins.**

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

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

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

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

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

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

(Opron, Xia and Wei, JCP, 2013; JCP 2014; JCP, 2015)

## Physical modeling

**Explicit solvent models (*Molecular dynamics, QM/MM, MC*)**
- Atomistic modeling of both solvent and solute molecules.
- Accurate but time consuming and subjects to force field errors.

**Integral equation models (*Ornstein-Zernike, Percus-Yevick and hypernetted-chain equations, RISM, LDFT, etc.*)**
- Continuous function modeling of solvent molecules, while atomistic modeling of the solute.
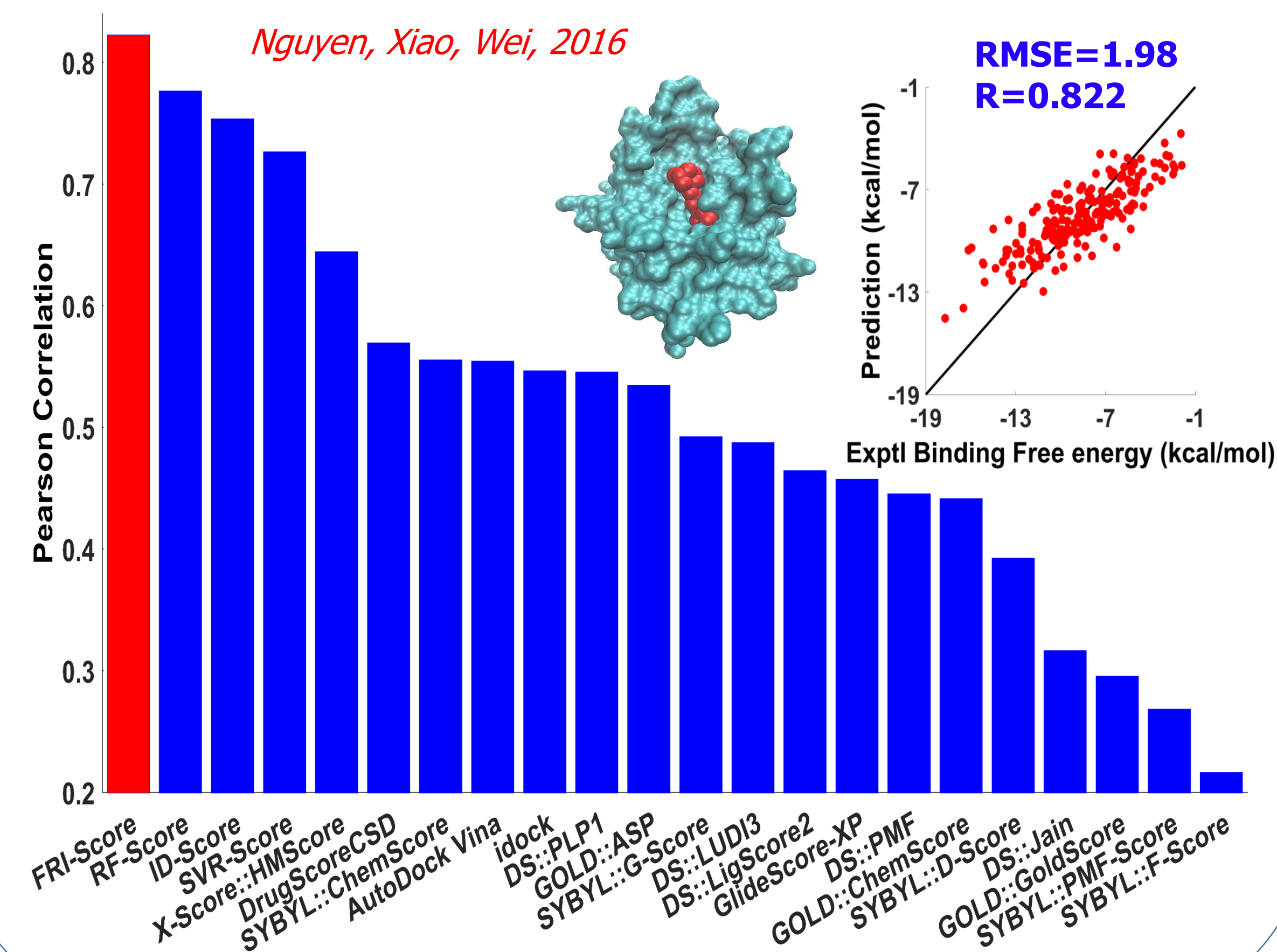- Still accurate but less time consuming.

**Implicit solvent models (*Image charge, Generalized Born, Poisson-Boltzmann, Polarizable Continuum*)**
- Dielectric continuum modeling of solvent molecules, while atomistic modeling of the solute.
- A good trade off between accuracy and efficiency.

**Variational multiscale models (*nonpolar, polar and QM*)**
- Couple polar and nonpolar components by variational surfaces.
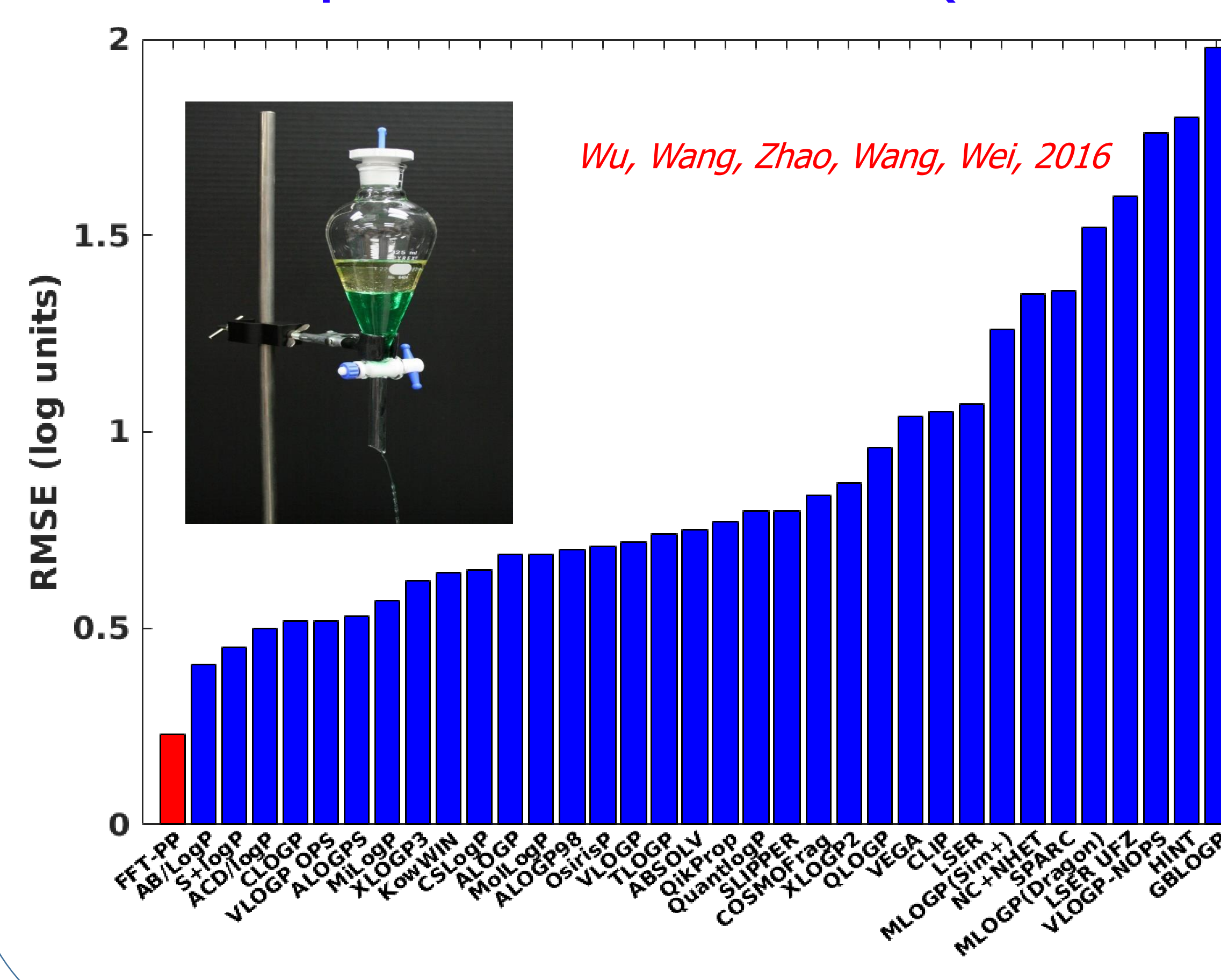- Self-consistent surface, charge, polarization and energy.

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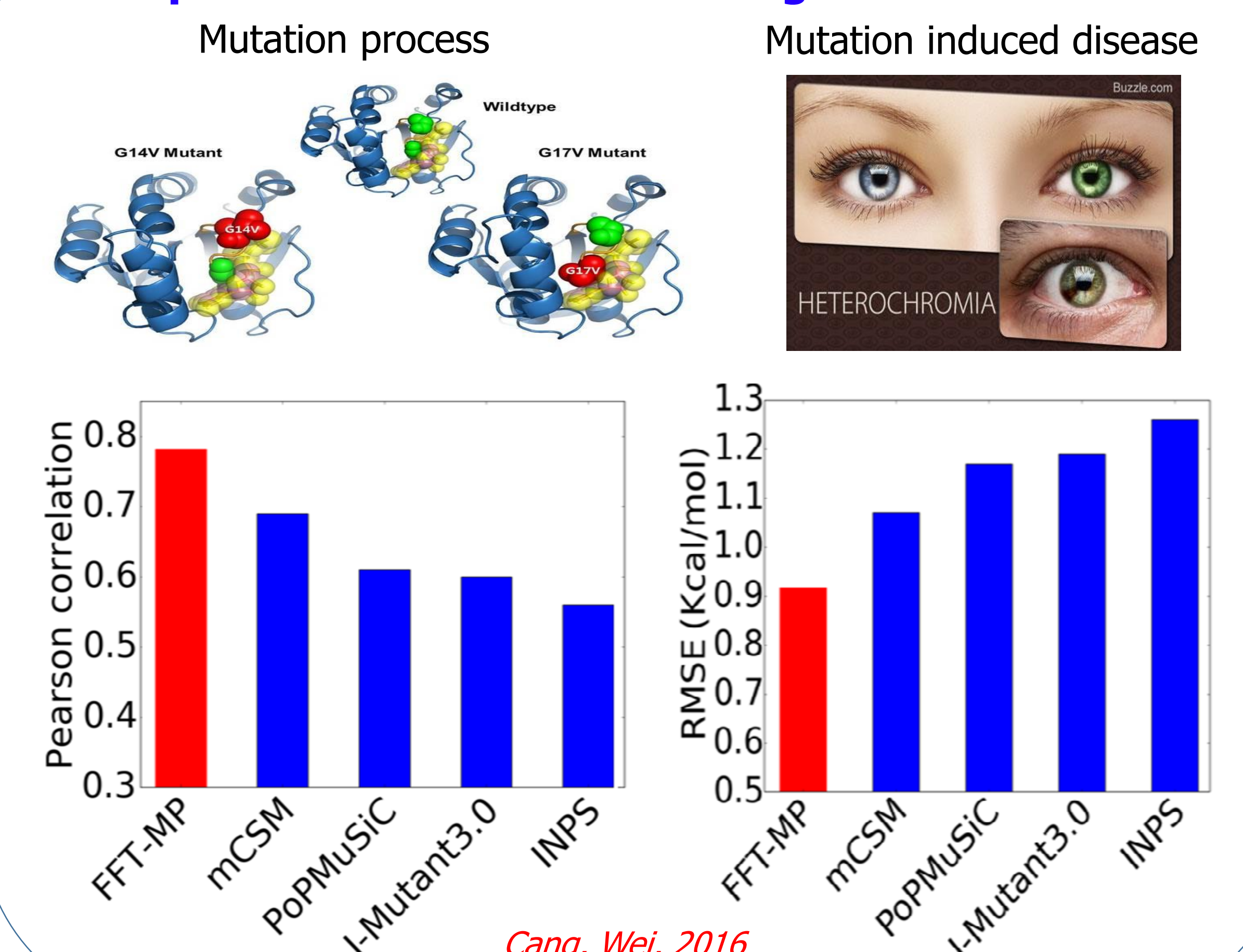


Nguyen, Xiao, Wei, 2016

RMSE=1.98
R=0.822

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



Wu, Wang, Zhao, Wang, Wei, 2016

## Blind prediction of mutation energies of 2648 dataset

Mutation process

Mutation induced disease

HETEROCHROMIA



Cang, Wei, 2016

## Acknowledgement:

http://users.math.msu.edu/users/wei/