# Bridging artificial intelligence and physics-based docking for better modelling of biomolecular complexes
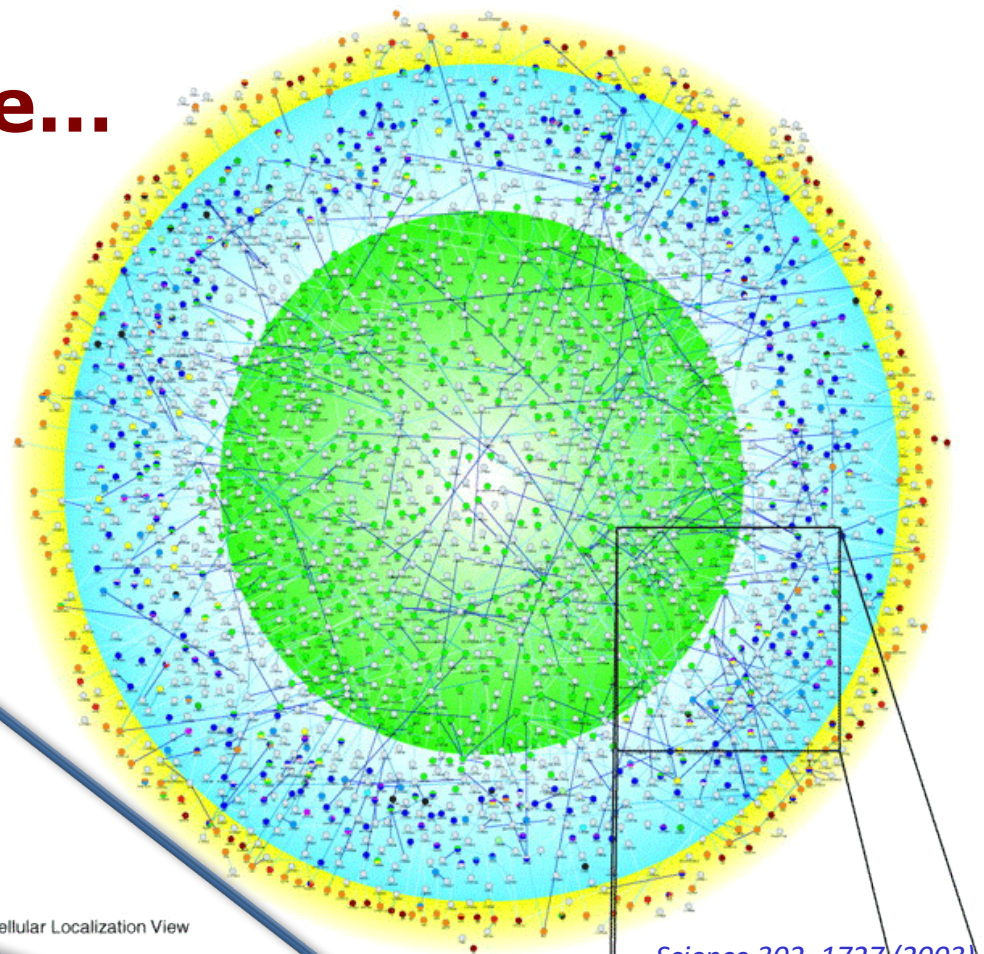
## Li Xue

**Computational Structural Biology Lab**
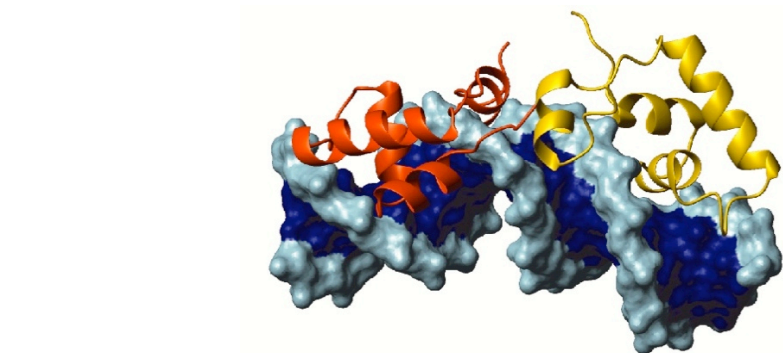
**Utrecht University, the Netherlands**

# The molecular machines of life

chromosome

DNA

University ... istry

# The network of life...



B

Sub-Cellular Localization View

Extracellular
Extracellular Matrix
Plasma Membrane
Synaptic Vesicle
Mitochondria
Endoplasmic Reticulum
Golgi
Lysosome
Cytoplasm
Cytoskeleton
Peroxisome
Ribosome
Centrosome
Nucleus
Unknown

Nuclear Proteins

Cytoplasmic Proteins

Membrane and
Extracellular Proteins

Interaction Ratings
0.9 - 1.0
0.8 - 0.9
0.65 - 0.8
< 0.65

*Science 302, 1727 (2003)*

Chemistry

**Universiteit Utrecht**

**Structures/Models**

- **Provide structural insight** into protein **function** and regulation.

- Can **guide experimental studies**

- Can help **rationalize** the effect of **genetic defects**

- Can function as **drug target**

**BUT: The number of experimentally solved protein structures are greatly lagged behind of sequences**

108,857,716

128,749

protein-only structures in PDB

protein sequences in UniProt

**This calls for complementary computational methods**

Universiteit Utrecht

[Faculty of **Science Chemistry**]

# Protein-Protein Docking & its challenges

**Protein A**

**Protein B**

Docking Program

Docking generates *vast numbers* of docked conformations

**Challenge 1 – Huge sampling space:** How to generate native-like (correct) models for *flexible* proteins?

**Challenge 2 – Scoring:** How to *select* the native-like models?

> 100 scoring functions published

PatchDock
2002

- ZRank   SwarmDock
- pyDock
  2007

interEvScore
2013

```
1990        1995        2000        2005      2010        2015        2020
```

2003
- HADDOCK
- RosettaDock

2006
ClusPro

2009
LZERD

2011
- DECK
- SIPPER
- IRAD
- SPIDER

1992
the first FFT-based docking
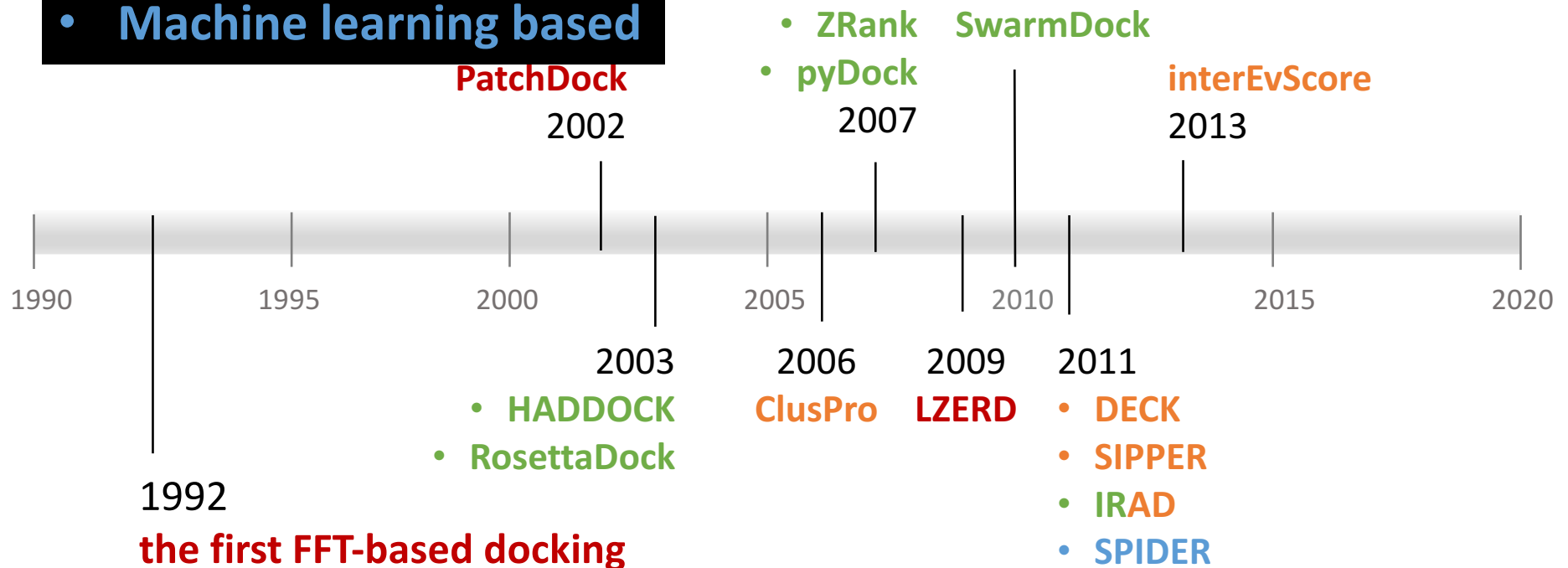
# Scoring in the past 25 years

> 100 scoring functions published

- **Shape complementarity**
- **physics based**
- **Statistical potentials**
- **Machine learning based**

- **ZRank**
- **pyDock**

**SwarmDock**

**PatchDock**

2002

2007

**interEvScore**

2013

1990    1995    2000    2005    2010    2015    2020

2003

- **HADDOCK**
- **RosettaDock**

2006

**ClusPro**

2009

**LZERD**

2011

- **DECK**
- **SIPPER**
- **IRAD**
- **SPIDER**

1992

**the first FFT-based docking**

7

# Scoring in the past 25 years

> 100 scoring functions published

- **Shape complementarity**
- **physics based**
- **Statistical potentials**
- **Machine learning based**

- ZRank    SwarmDock
- pyDock    interEvScore

**PatchDock**    2007    **2013**
2002

1990    1995    2000    2005    2010    2015    2020

2003    2006    2009    2011

- **HADDOCK**    **ClusPro**    **LZERD**    - **DECK**
- ~~**RosettaDock**~~    - **SIPPER**
    - **IRAD**
    - **SPIDER**

1992
**the first FFT-based docki**

14 years unchanged and still one of the best.

> 100 scoring functions published

- More data accumulated
- More computational power
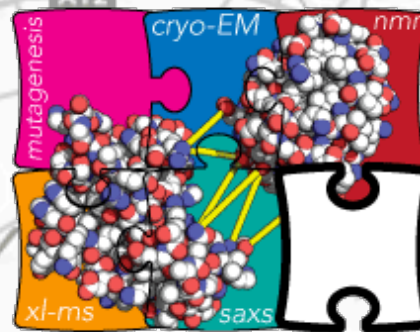- Better machine learning algorithms

# Can we do Better?

1990

2020
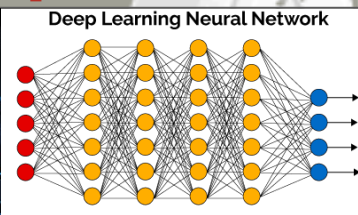
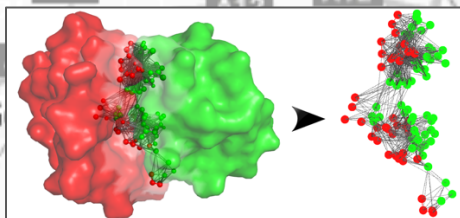1992
the first FFT-based docking

- IRAD
- SPIDER

**DeepRank**

**metaScore**

**iScore**

**iSEE**

$\Delta\Delta G$

Deep Learning Neural Network

The ensemble model

Forest output probability $p(c|\mathbf{v}) = \frac{1}{T}\sum_{t}^{T} p_t(c|\mathbf{v})$

11

# Our DeepRank project:
# Why Deep Learning, specifically convNets?



**A docked model
with its interface (red/pink) in a grid box**

**Human expert visually checking the interface of a docked model**
*@ CAPRI 38 competition – manual selection stage*

**Use convNets to scan the
interface of a docked model**

*https://github.com/DeepRank*

STAYTUNED

**iScore**

**iSEE**

$ΔΔG$
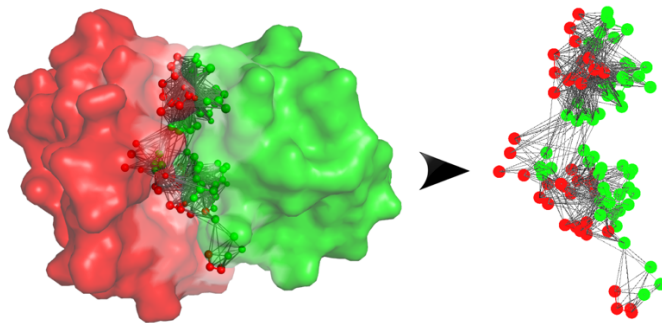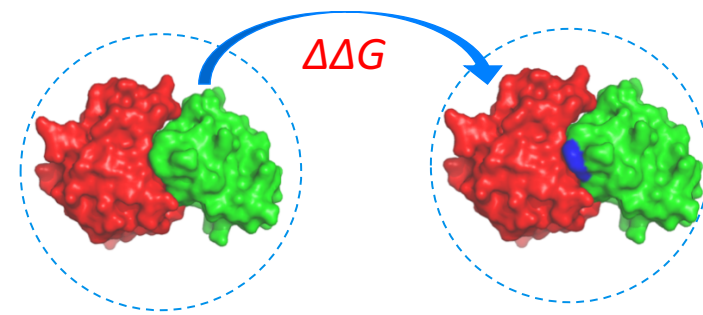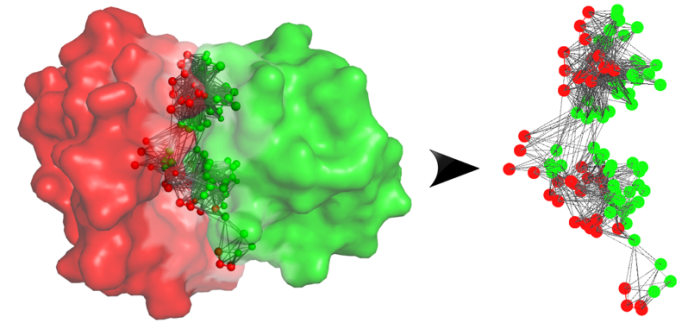
structure ➜ **interactions** ➜ **binding free energy** ➜ function

# iScore

interface graph based docking scoring function

**structure** ➛ <span style="color:darkred">**interactions**</span> ➛ **binding free energy** ➛ **function**

# iScore: a novel machine learning based scoring function

**I propose a novel approach that treat the scoring problem as a *graph comparison problem.***
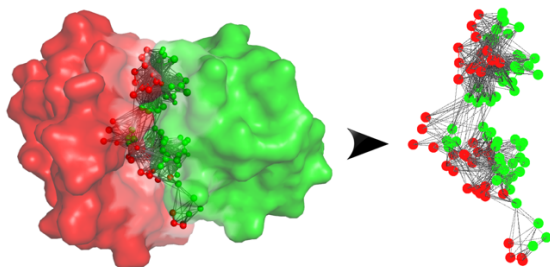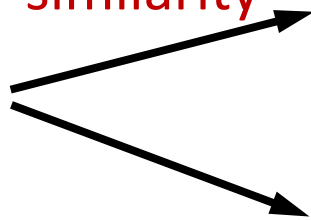
**A docked model and its interface graph**

<span style="color:red">Calculate graph similarity</span>

*Figure credit: Cunliang Geng*

**Training data**

**Positive data**

Interface graphs from *correct* protein-protein structures

**Negative data**

Interface graphs from *wrong* protein-protein models

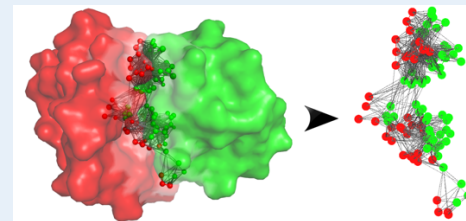# iScore: a novel machine learning based scoring function

*Strength #1*: **Much higher resolution of the input information is considered.**

Typical machine learning scoring function :

Uses **whole interface descriptors** *instead of* working at the resolution of atoms or residues.

iScore exploits:

- Network topology
- Atom/Residue level information (node label)
- Pairwise information (edge label)

*Strength #2*: Full profile of interface conservation can be exploited.



**Evolution information is critical for**
- **Protein recognition**
- **Protein folding**

Position
1
2
3
4
...

Conservation profiles (a PSSM)

# iScore v1.0: an evolution based scoring function

**iScore v1.0:**

- uses residue conservation only.

- Residue level resolution

Trained on 114 bound complexes, 114 wrong models.

Tested on 64 cases, 400 docked models per case.

# iScore vs. HADDOCK score

| HADDOCK score | iScore v1.0 |
|---|---|
| Atom resolution | Residue resolution |
| Interaction cutoff: 8.5 Å | Interaction cutoff: 6 Å |
| Linear function | Non-linear graph similarity based |
| Interaction energy based | Residue conservation based |
| Fast: 2 seconds per model | Slow: 0.05-0.5 second per graph pair , but many comparisons |

**A hit**: a correct (near-native) docked model (interface RMSD $\leq$ 4 Å)

**Success Rate:**

the percentage of cases that have at least one hit among the top $m$ conformations.

**Success Rate Plot**



Success Rate

Method 1
Method 2

Top m

20

**A hit**: a correct (near-native) docked model (interface RMSD $\leq$ 4 Å)

**Average Hit Rate:** the percentage of hits that are included among the top $m$ conformations.



Hit Rate Plot

Hit Rate

Top m

Method 1
Method 2

21

# iScore vs. HADDOCK score



**The Combined Score is the average of iScore and Haddock Score.**

- iScore is an effective graph based scoring function.
- iScore v1.0 (*based on conservation only*) outperforms HADDOCK scoring function on the majority of dimer cases of docking benchmark 4.0.

iScore complements HADDOCK score:

HADDOCK score = 1.0 Evdw + 0.2 Eelec + 1.0 Edesolv

$\Delta\Delta G$

# iSEE

**i**nterface **S**tructure, **E**nergy and **E**volution based $\Delta\Delta G$ predictor

**structure** → **interactions** → <span style="color:red">**binding free energy**</span> → **function**

# Why are we interested in mutations?

- Coding variants -> phenotype
- Carcinogens
- Errors in reading tRNAs at the ribosome
- The generation of antigen-binding CDR loops of antibodies
- Protein engineering
- Directed evolution

# iSEE

iSEE aims to model:

*Binding affinity change upon mutation (ΔΔG) = f (features of the mutation site)*



*Geng et al. iSEE: Interface Structure, Evolution and Energy-based random forest predictor of binding affinity changes upon mutations, submitted to Bioinformatics.*

# 10 times 10-fold cross validations

**iSEE:** interface **S**tructure, **E**nergy and **E**volution based ΔΔG predictor.

# iSEE on two blind test datasets

### The NM dataset (19 mutations)

| Method | RMSE* | PCC |
|--------|-------|-----|
| iSEE | 1.37 | 0.73 |
| **BindProfX** | **1.11** | **0.81** |
| **FoldX** | **1.15** | **0.72** |
| **ZeMu** | **1.28** | **0.70** |
| CC/PBSA[a] | 1.33 | 0.60 |
| pred2 | 1.44 | 0.48 |
| pred1 | 1.49 | 0.39 |
| BeAtMuSiC | 1.70 | 0.24 |
| mCSM | 1.83 | 0.16 |

### The MDM2-p53 dataset (33)

| Method | RMSE* | PCC |
|--------|-------|-----|
| iSEE | 0.77 | 0.65 |
| BindProfX | 1.17 | 0.51 |
| FoldX | 1.36 | 0.50 |
| BeAtMuSiC | 1.02 | 0.48 |
| mCSM | 0.97 | 0.22 |

**\*RMSE in kcal mol$^{-1}$**

[a] The CC/PBSA predictor was trained on a dataset including the NM dataset.

**iSEE competes favorably with state-of-the-art ΔΔG predictors on two blind test datasets.**

# Full mutation scanning on MDM2-p53 complex

MDM2-p53 is a prime target for cancer therapeutics



p53

MDM2

PDB ID: 1YCR

**MDM2-p53 dataset**

➢ **33 mutations: MDM2(16) and p53(17).**

➢ **With corresponding experimental $\Delta\Delta G$ data.**

# Full mutation scanning on MDM2-p53 complex



+ **Experimental hot-spots**
* **Predicted hot-spots**

**iSEE enables high-throughput scanning.**

# Summary of iSEE



1. **iSEE is a machine learning-based ΔΔG predictor.**

2. **iSEE competes with state-of-the-art predictors.**

3. **iSEE can be used for high-throughput applications.**

https://github.com/haddocking/iSee

**iScore**

**iSEE**

$\Delta\Delta G$

**structure** → **interactions** → **binding free energy** → **function**

# Artificial Intelligence Toward protein Interaction Prediction

**Evolution information**



**Secondary structures**



**Physico-chemical attributes**



**Interface propensity**



**A novel interface potential score**

$$e(w_P, w_R) = \min_{r \leq 5\text{Å}} e(A_n, R_n, r) +$$

$$\sum_{i \neq n} \min_{j, r \leq 10\text{Å}} e(A_i, R_j, r)$$

*Improved 3D modelling of protein interaction*



Interface

# Acknowledgement

**Thank You !**

What can we learn from 3D structures (models) of complexes?



◎**Models provide structural insight into function and mechanism of action**

◎**Models can drive and guide experimental studies**

◎**Models can help understand and rationalize the effect of disease-related mutations**

◎**Models provide a starting point for drug design**

- Docking
- Binding affinity estimation
- Hot-spot identification (ddG)
- Interaction design
- Interactome prediction

**structure → interactions → binding free energy → function**

# Evaluation - melquiplot

**y-axis**:
Decoy Ranking
(the smaller the better)

400

200

1

Scoring method 1          Scoring method 2

# iScore vs. HADDOCK score

# Computational methods for predicting $\Delta\Delta G$

**Rigorous methods, e.g., FEP**
Accurate

*However,*
Requires extensive
conformation sampling

**Empirical energy based methods**
Fast but not accurate

**Machine learning based**

Fast and can integrate
heterogonous data



40

# Current leading ΔΔG predictors

**BindProfX (Y Zhang et al., 2017)**

Interface structural profile + FoldX physics potentials

**mCSM (D Pires et al., 2014)**

contacts + pharmacophore change on mutation position

**ZeMu (D Dourado et al., 2014)**

MD simulations for flexibility of mutation zone + FoldX potentials

**BeAtMuSiC (Y Dehouck et al., 2013)**

Statistical potentials

**CC/PBSA (A Benedix et al., 2009), pred1 and pred2 (A Panchenko et al., 2014)**

MD simulations +  PBSA for solvation effect

**FoldX (R Guerois et al., 2002)**

10 force field based energy terms + rotamers for sidechain flexibility

# Binding affinity change upon mutation (ΔΔG)



$$\Delta G_{wildtype}$$

$$\Delta G_{mutant}$$

$$\Delta\Delta G = \Delta G_{mutant} - \Delta G_{wildtype}$$
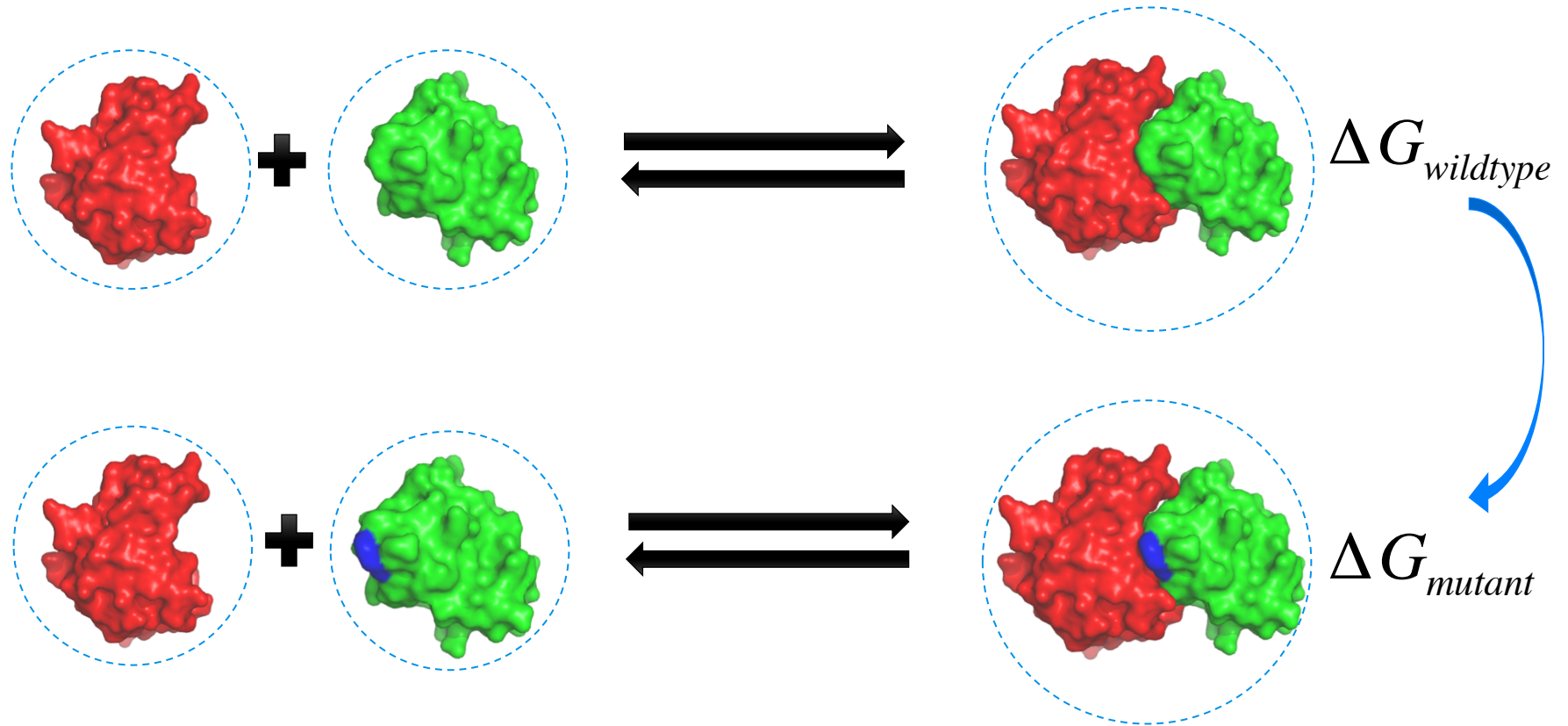
# iSEE uses random forest to predict $\Delta\Delta G$

iSEE aims to model:

*Binding affinity change upon mutation ($\Delta\Delta G$) = f (features of the mutation site)*

A random forest is a forest of decision trees.



**A conceptual decision tree.**

**INPUT**: a feature vector of the mutation site

**OUTPUT**: predicted $\Delta\Delta G$

$\Delta$ Eelec < 50

no — $\Delta$ Evdw < 100

yes — $\Delta\Delta G$ = 3 kcal/mol

$\Delta\Delta G$ = 0.3 kcal/mol

$\Delta\Delta G$ = -2 kcal/mol

# Methods



SKEMPI-based training dataset

Structural, energetic and evolutional features

Random Forest model

iSEE

NM dataset

MDM2-p53 dataset

**Training**

**Testing**

# Training and Evaluation

**10 times 10-fold cross validations are used for training.**

**Root Mean Square Error (RMSE):**

$$RMSE = \sqrt{\frac{\sum_1^n(\Delta\Delta Gexp - \Delta\Delta Gpred)^2}{n}}$$

**Pearson's Correlation Coefficient (PCC):**

$$PCC = \frac{cov(\Delta\Delta Gexp, \Delta\Delta Gpred)}{\sigma_{\Delta\Delta Gexp} \cdot \sigma_{\Delta\Delta Gpred}}$$

# Training and test datasets

➢ **Wildtype protein-protein complexes with crystal structures and experimental $\triangle\triangle$G values**

➢ **Only *single point mutations* in the *interface* of dimer complexes are considered.**

**Training dataset:**

**1102 mutations in 57 complexes from the SKEMPI database**

**Blind test datasets:**

**1. NM dataset**

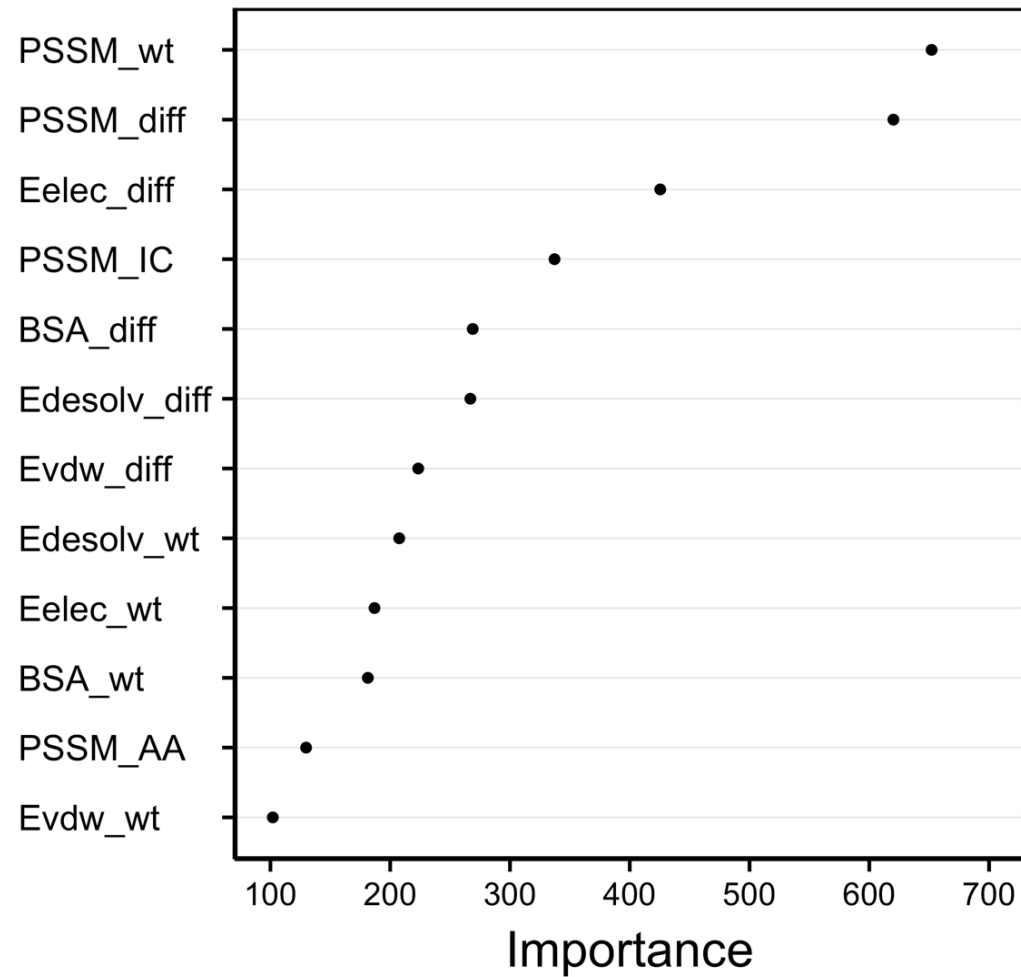   **37 mutations in 2 complexes, 1BXI(18) and 1IAR(19).**

**2. MDM2-p53 complex dataset**

   **33 mutations in the novel MDM2-p53 complex. MDM2(16) and p53(17).**
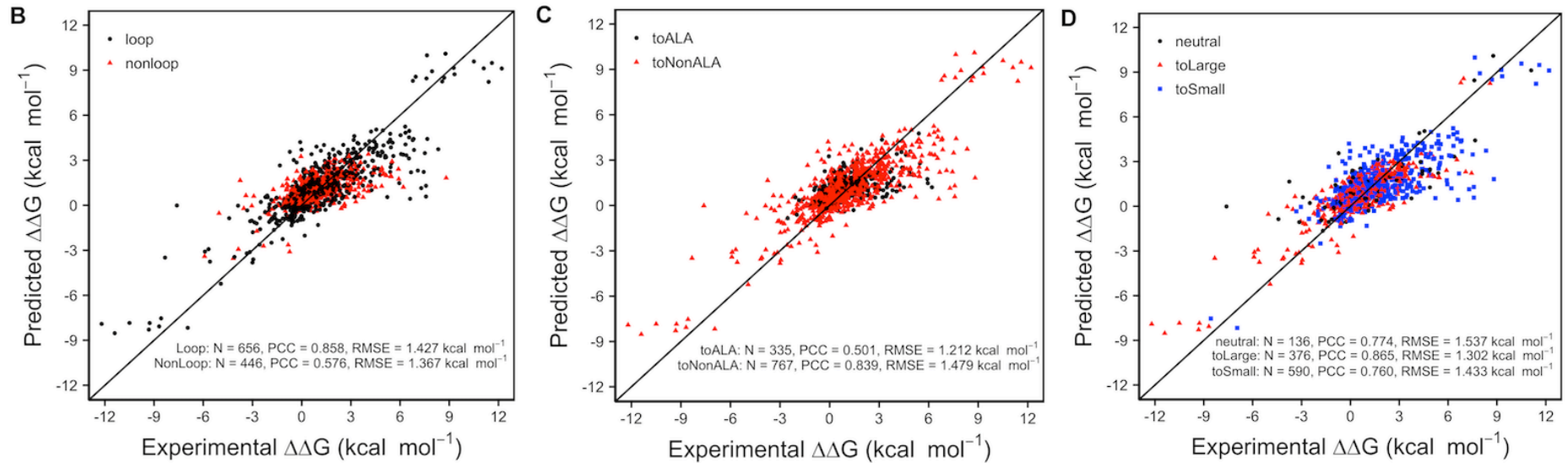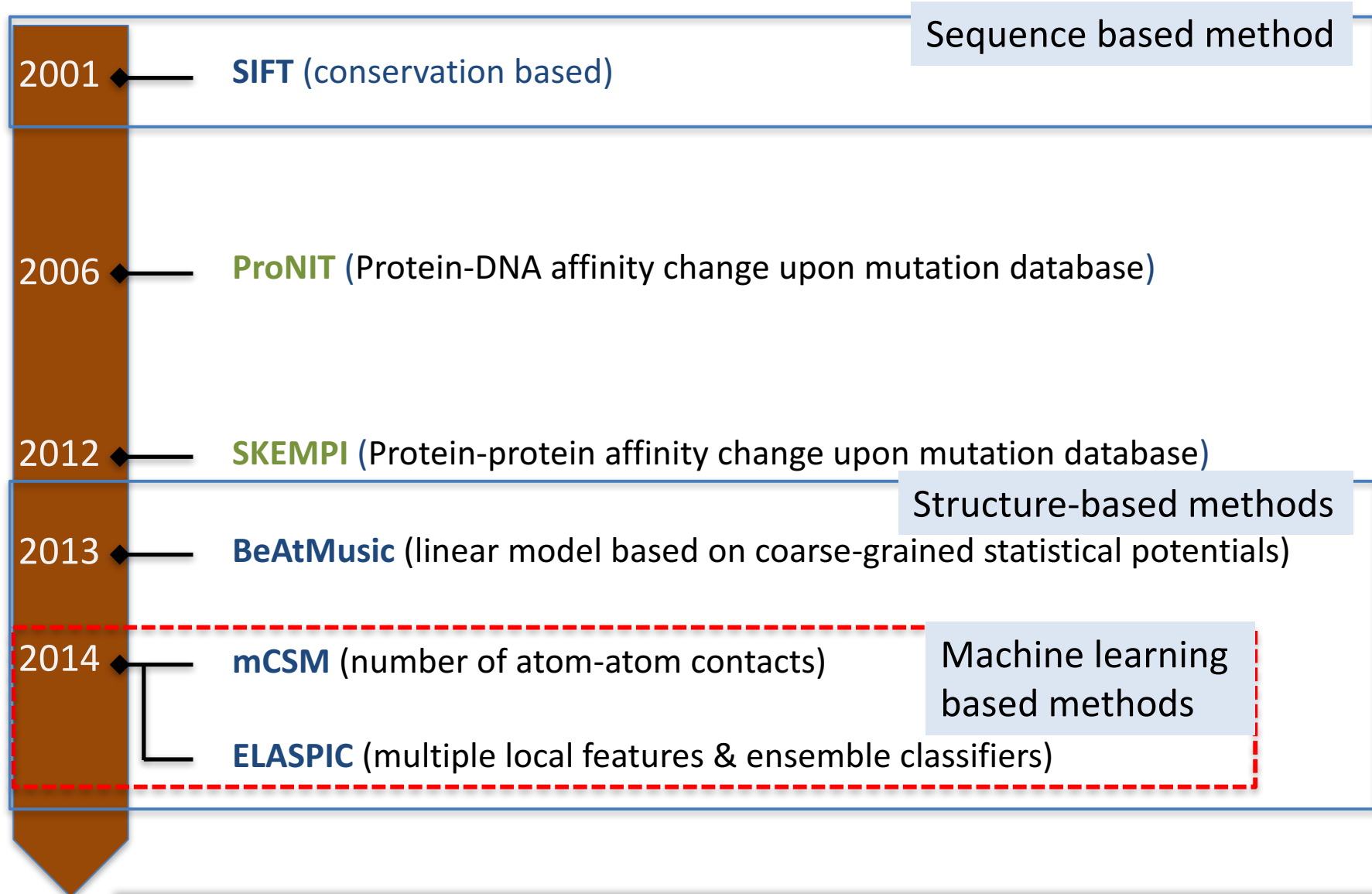
# Feature importance



**Evolutionary features ranked top.**

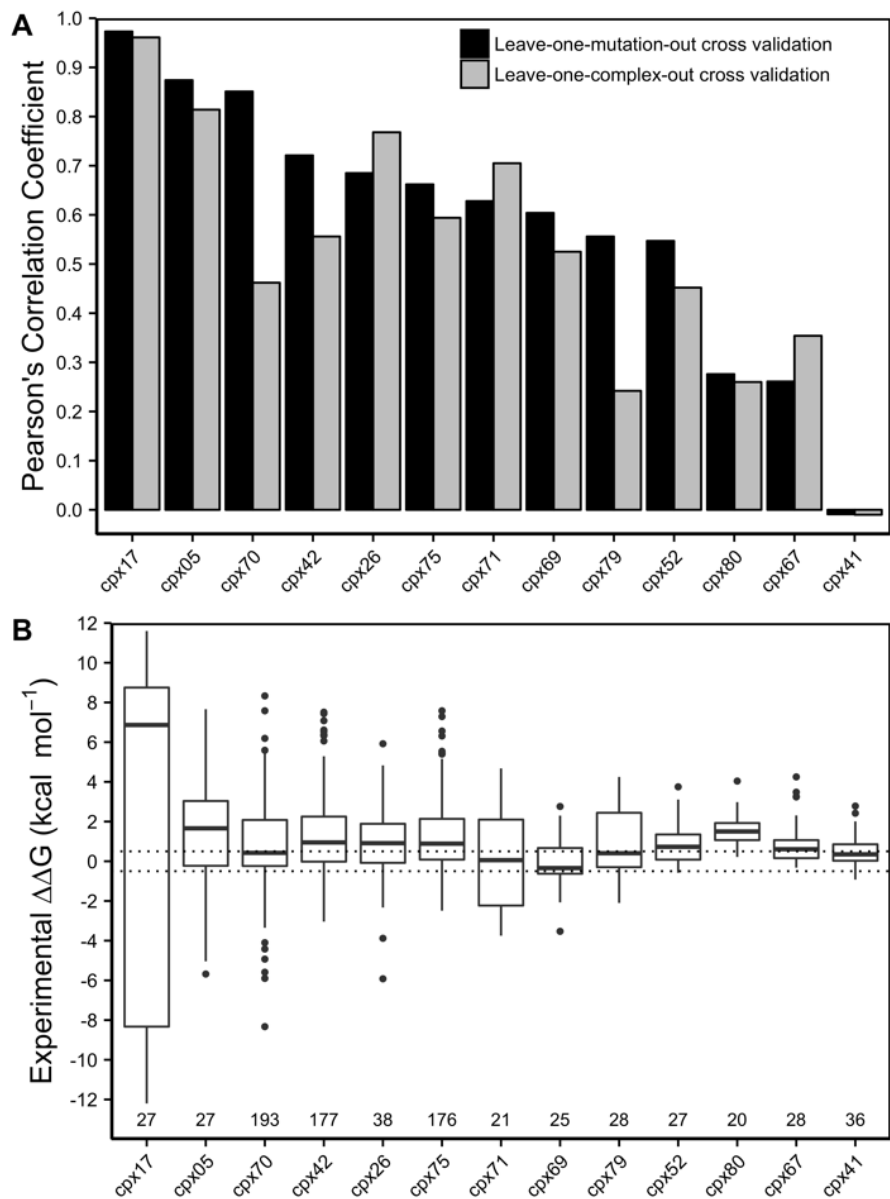# iSEE performance on different types of mutations



**iSEE highly and consistently performed on different types of mutations**
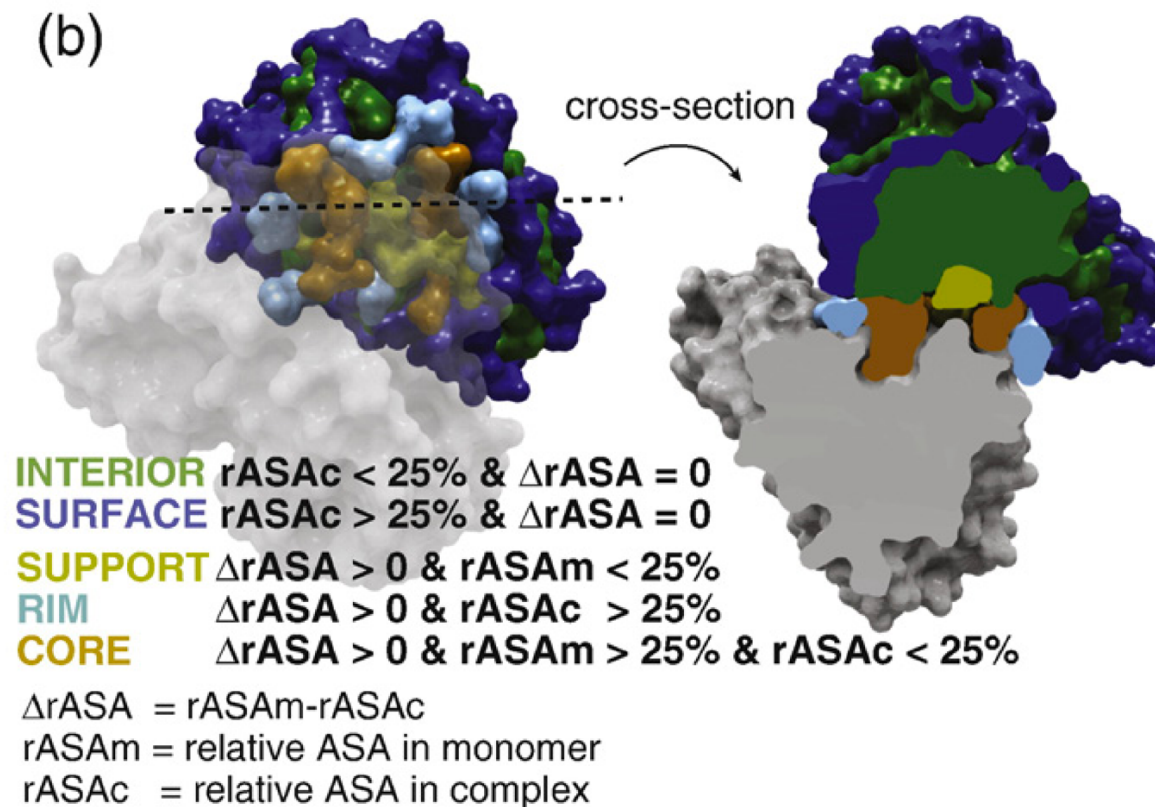
# Machine Learning based prediction methods

**2001** — **SIFT** (conservation based)

Sequence based method

**2006** — **ProNIT** (Protein-DNA affinity change upon mutation database)

**2012** — **SKEMPI** (Protein-protein affinity change upon mutation database)

Structure-based methods

**2013** — **BeAtMusic** (linear model based on coarse-grained statistical potentials)

**2014** — **mCSM** (number of atom-atom contacts)

Machine learning based methods

**ELASPIC** (multiple local features & ensemble classifiers)

Machine learning based methods are significantly better than other methods

# Interface positions



(b)

cross-section

INTERIOR $rASAc < 25\%$ & $\Delta rASA = 0$
SURFACE $rASAc > 25\%$ & $\Delta rASA = 0$
SUPPORT $\Delta rASA > 0$ & $rASAm < 25\%$
RIM $\Delta rASA > 0$ & $rASAc > 25\%$
CORE $\Delta rASA > 0$ & $rASAm > 25\%$ & $rASAc < 25\%$

$\Delta rASA$ = $rASAm$-$rASAc$
$rASAm$ = relative ASA in monomer
$rASAc$ = relative ASA in complex

# Optimization of RF models



**Achieved highest performance at ntree=80 and mtry=7**