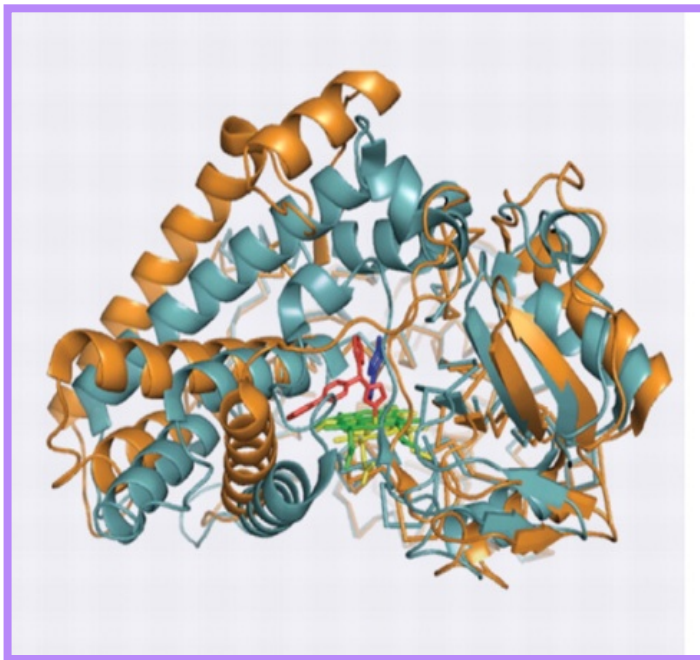# Intrinsically accessible motions enable
## Optimal binding of substrate or drugs



Conformational flexibility +
sequence variability mediates
**substrate selectivity**

- **Two conformations of P450-CYP2B4:**
  **open** (orange) with a large substrate (bifonazole, red), and
  **closed** (light blue) with the smaller substrate 4-(4-chlorophenyl) imidazole (blue)

See...

# Sequence evolution
## an information-theoretic approach

Residue index



correlated mutations

conserved

Information entropy (Shannon, 1951 )

$$S(i) = \sum_{x_i=1}^{20} P(x_i) \log \frac{1}{P(x_i)}$$

Mutual information (MI)

$$I(i,j) = \sum_{x_i=1}^{20} \sum_{y_j=1}^{20} P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(x_i) P(y_j)}$$

**for correlated mutations analysis (CMA)**

# Mutual Information
## without the influence of phylogeny

MIp - to eliminate random noise and phylogenetic components

$$\mathrm{MI_p}\,(i, j) = \mathrm{I}(i, j) - \mathrm{APC}$$

APC = Average product correction

$$= [\ \mathrm{I}(i, x)\ \mathrm{I}(j, x)\ ]\ /\ <\mathrm{I}(i, j)>$$

| R | | | | E | V | N | |
|---|---|---|---|---|---|---|---|
| E | | | | K | V | N | |
| K | | | | E | V | N | |
| R | | | | D | V | S | |
| D | | | | K | V | S | |
| D | | | | K | V | S | |
| E | | | | R | V | S | |

where $\mathrm{I}(i, x)$ is the mean mutual information of column $i = \sum_j \mathrm{I}(i, j)$

Dunn, Wahl and Gloor (2008) *Bioinformatics* **24**: 333-340

# HIV-1 protease correlated mutation analysis (CMA)

MSA of HIV-1 protease

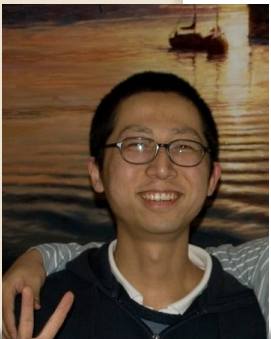MI matrix $\mathbf{I}_{ij} = \mathrm{I}\,(i,j)$

Shi and Malik (2000)

spectral clustering



residue index

reordered residue index

**Dr. Ying Liu**

# MDR mutations distinguished by CMA



MSA of HIV-1 protease
Stanford HIV Drug Resistance Database
http://hivdb.stanford.edu/

```
CTLVGTAIHEMMHALGFLHEQNREDRDDWVR
CDKFGIVVHELGHVVGFWHEHTRPDREDHVV
CFRFGTVIHEFIHALGFYHAQSAYTRDDYVL
NFTVGSLIHEIGHAFGLIHEHQRPDRDDYVI
CLTYGTPIHELMHALGFFHEQNRHERDSYVR
CDKFGIVVHELGHVVGFWHEHTRPDREKHVV
CDKFGVVVHELGHVVGFWHEHTRPDRNEFVG
CAYFGTIVHEIGHAIGFHHEQSRPDRDDYIN
CVYHGIIQHELSHALGFYHEHTRSDRNKYVR
CINSGTIIHEVLHALGVHHEQARADRDGYVT
```

untreated

reordered residue index

mobility profile

high

low

Drug-resistant cluster

```
CTLVGTAIHEMMHALGFLHEQNREDRDDWVR
CDKFGIVVHELGHVVGFWHEHTRPDREDHVV
CFRFGTVIHEFIHALGFYHAQSAYTRDDYVL
NFTVGSLIHEIGHAFGLIHEHQRPDRDDYVI
CLTYGTPIHELMHALGFFHEQNRHERDSYVR
CDKFGIVVHELGHVVGFWHEHTRPDREKHVV
CDKFGVVVHELGHVVGFWHEHTRPDRNEFVG
CAYFGTIVHEIGHAIGFHHEQSRPDRDDYIN
CVYHGIIQHELSHALGFYHEHTRSDRNKYVR
CINSGTIIHEVLHALGVHHEQARADRDGYVT
```
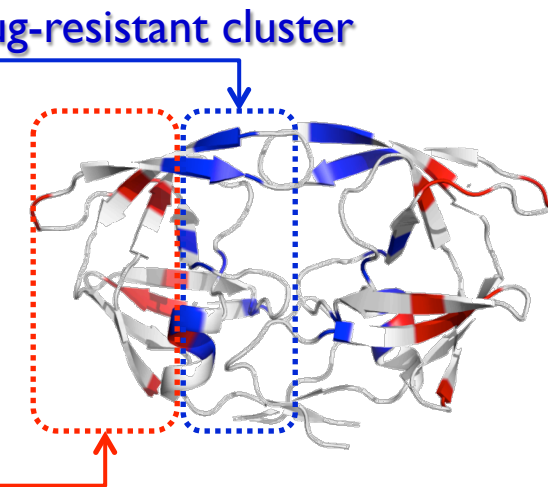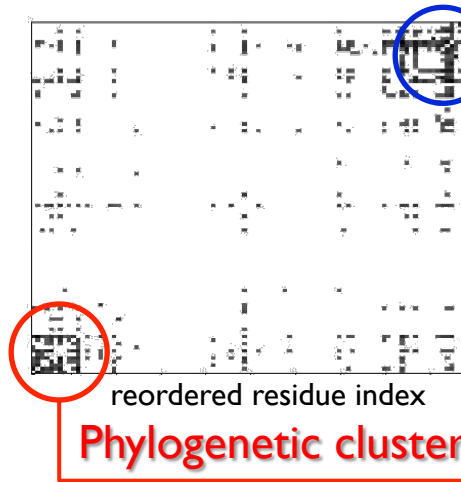
treated by at least
one drug

reordered residue index

Phylogenetic cluster

# Summary

- two groups of correlated mutation sites

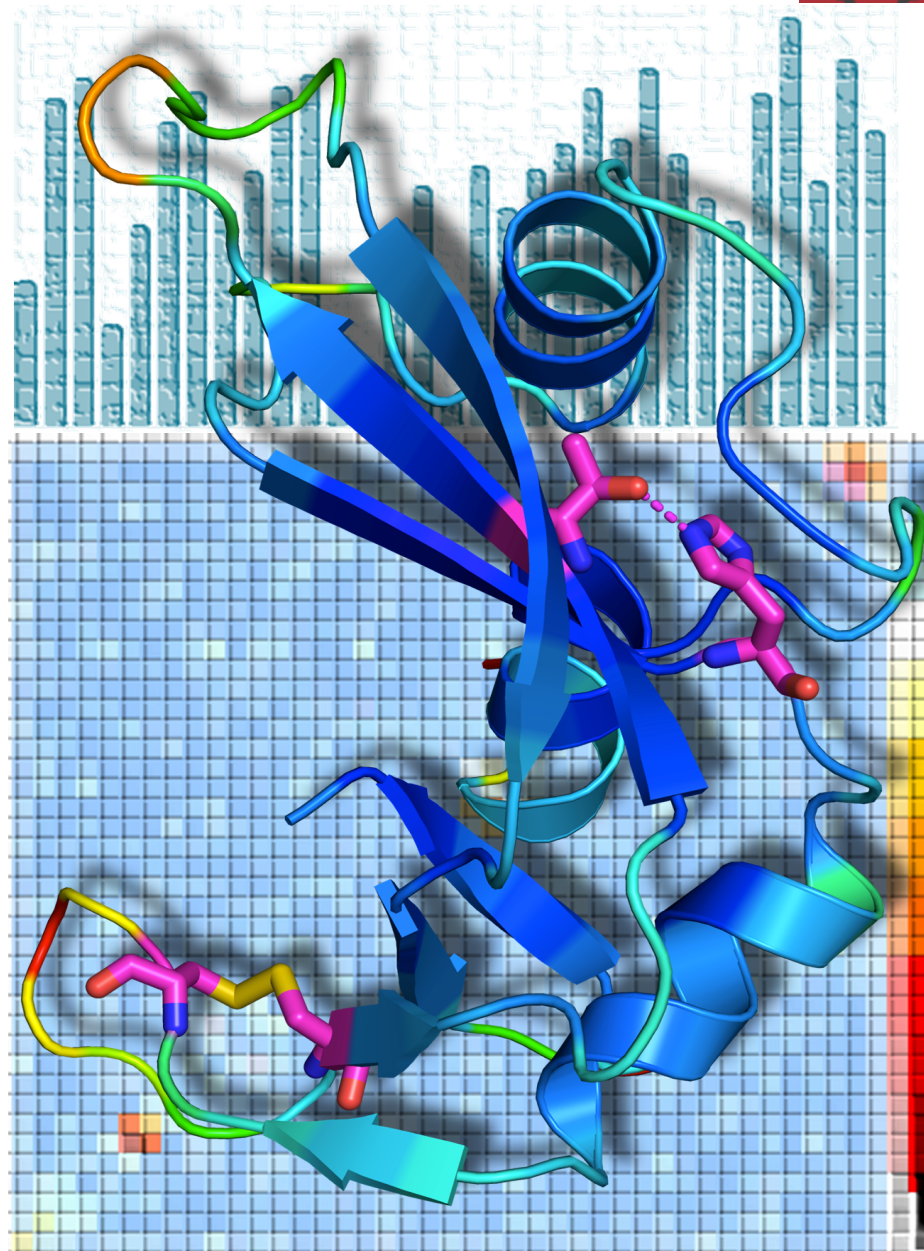| functional aspects | Structural location | structural dynamics |
|---|---|---|
| phylogenetic | exposed | mobile |
| multi-drug resistant | dimerization interface | restrained |

# Questions:

- Are key mechanical sites (e.g. hinges) conserved?

- Is there any correlation between sequence variability and structural dynamics?

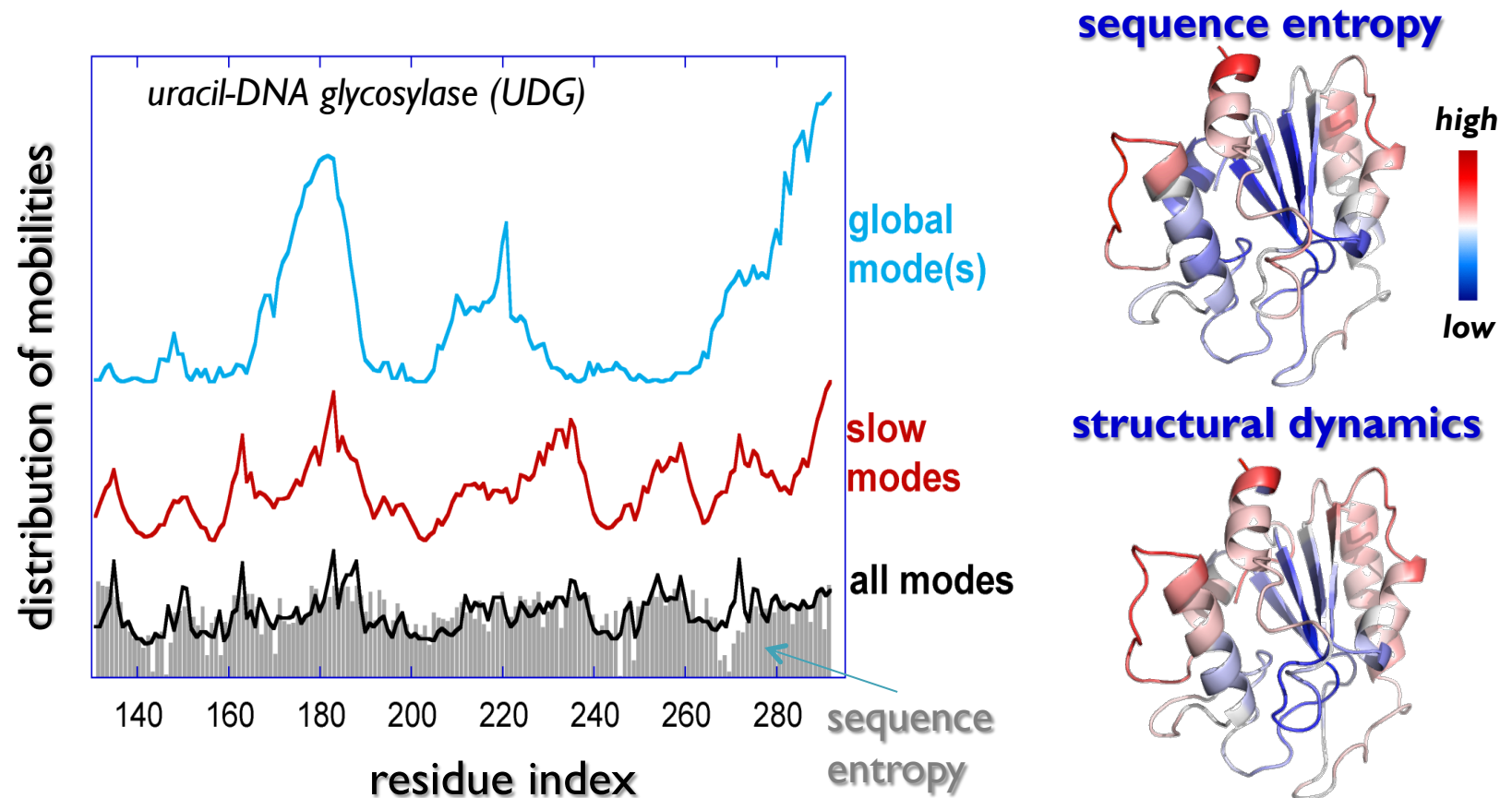- How does the structure ensure substrate specificity *and* conformational adaptability?

# A systematic study of a set of enzymes

# Evol

# Correlation between sequence entropy & conformational mobility



uracil-DNA glycosylase (UDG)

global mode(s)

slow modes

all modes

sequence entropy

distribution of mobilities

residue index

sequence entropy

high

low

structural dynamics

Liu Y, Bahar I (2012) "Sequence Evolution Correlates with Structural Dynamics" *Mol Biol Evol 9, 2253-63*
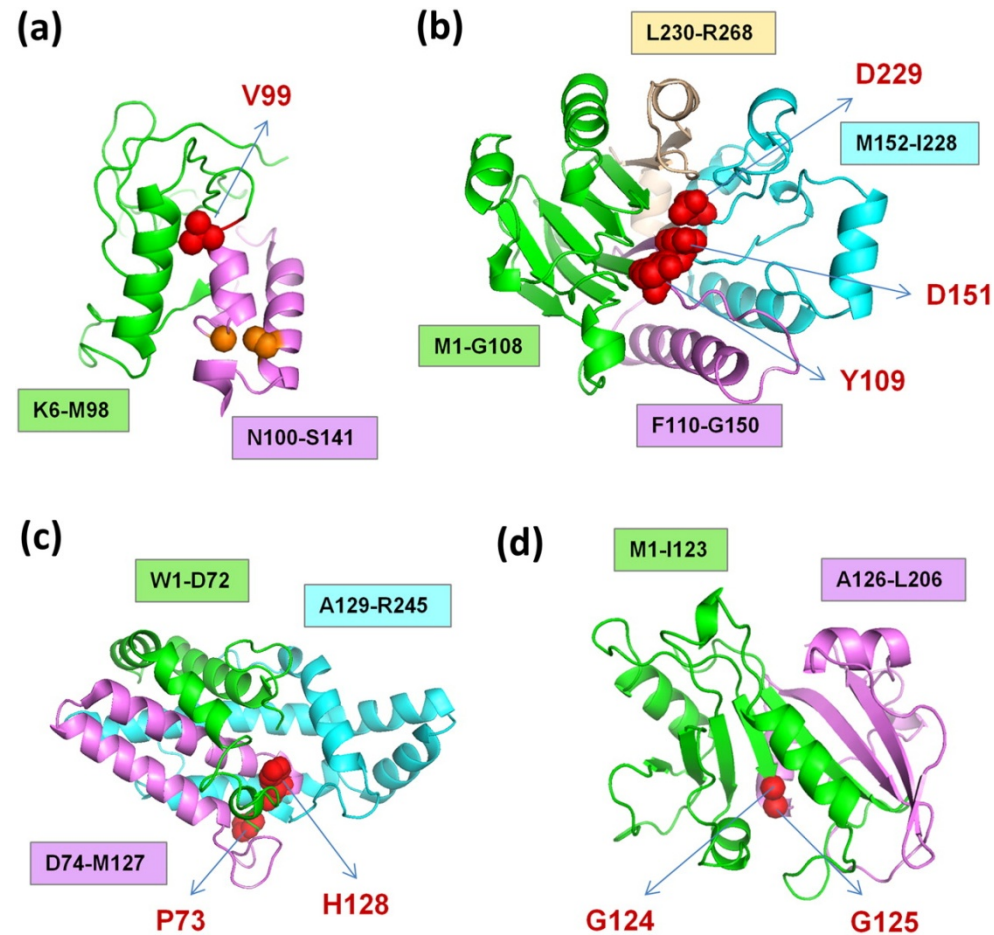
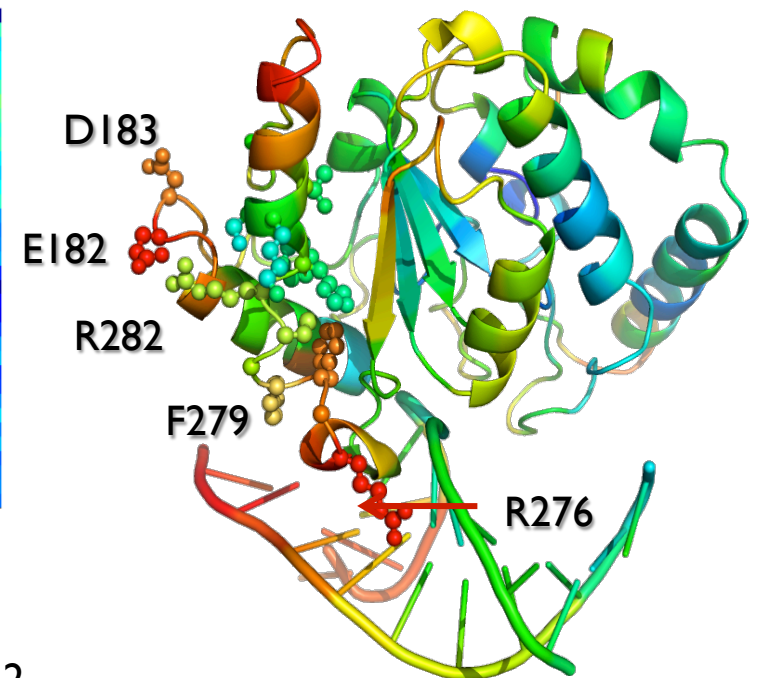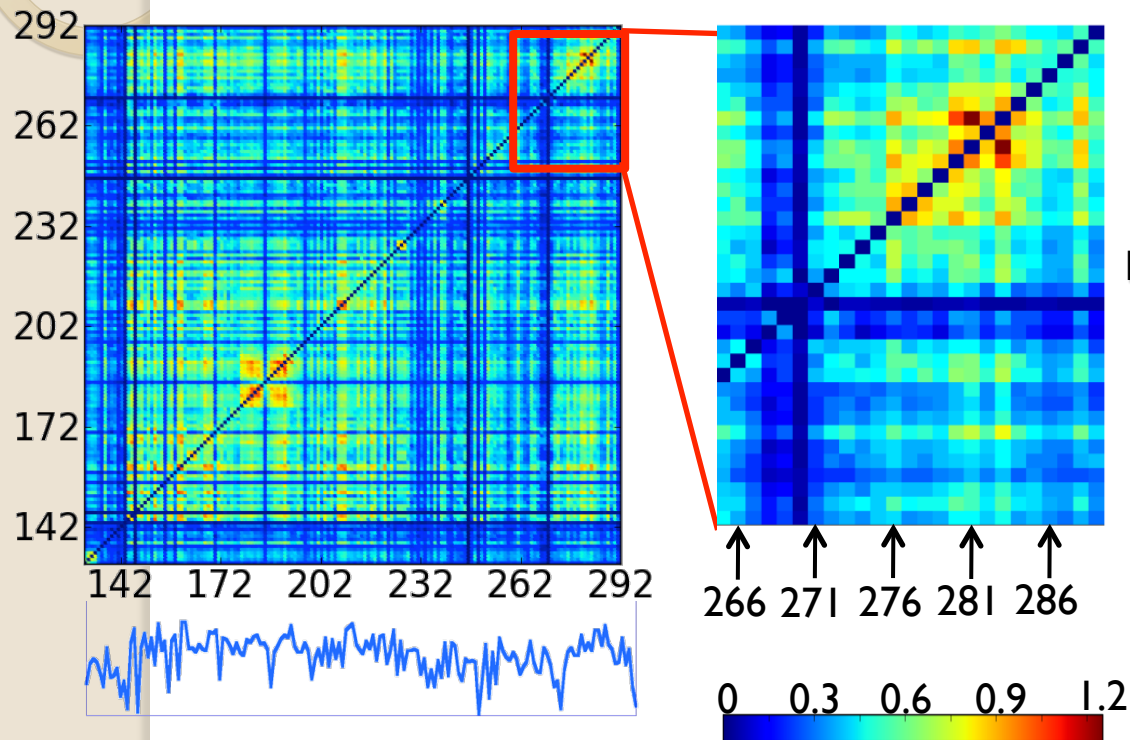# Mobility increases with sequence entropy

# Hinge sites are evolutionarily conserved

*despite their moderate-to-high exposure to environment*

Amino acids involved in intermolecular recognition are distinguished by **their co-evolution propensities**

# 3 Amino acids involved in intermolecular recognition are distinguished by **their high global mobility**



*cathepsin B*

substrate

residues involved in top 0.05% of I(*i*, *j*) values

$\langle M \rangle|_2$

MSF

residue index (*i*)

Liu Y, Bahar I (2012) "Sequence Evolution Correlates with Structural Dynamics" *Mol Biol Evol 9, 2253-63*

# **Summary**

Four types of functional sites

| Functional site | Mobility in global modes | Sequence evolution | Dominant Feature |
|---|---|---|---|
| Chemical (catalytic, ligand binding) | Minimal | Conserved | high fidelity, precision |
| Core | Minimal | Conserved | high stability |
| Hinge sites | Minimal | Conserved | rotational flexibility |
| Substrate recognition (specific) | High | High co-evolution propensity | adaptability |

# There are several methods for evaluating sequence co-evolution
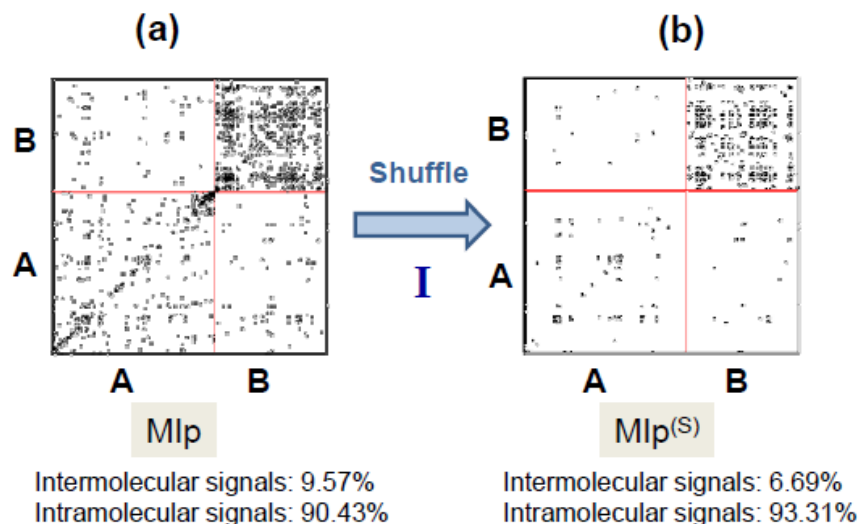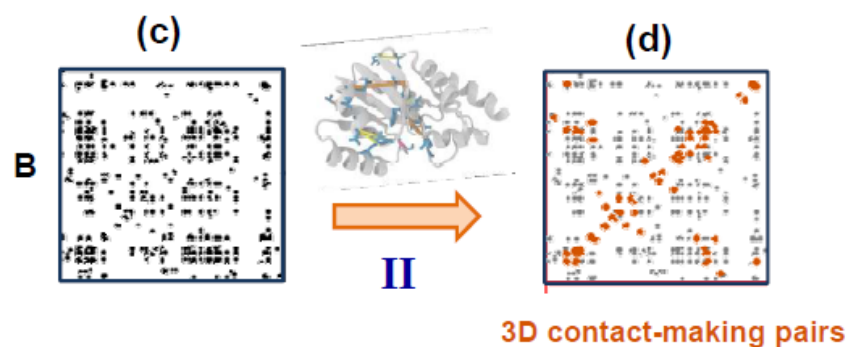
Four possible outcomes:

- True positive (TP) – correctly predicted to be a hit
- False positive (FP); predicted but it is a miss
- True negative (TN) –  correctly predicted to be a miss
- False negative (FN) – predicted as a miss, but is a hit

# Two criteria for assessing the performance of different methods



(a)

B

A

A B

Mlp

Intermolecular signals: 9.57%
Intramolecular signals: 90.43%

**Shuffle**

**I**

(b)

B

A

A B

Mlp(S)

Intermolecular signals: 6.69%
Intramolecular signals: 93.31%

(c)

B

**II**

(d)

**3D contact-making pairs**

- Minimizing false positives (signals between non interacting proteins)
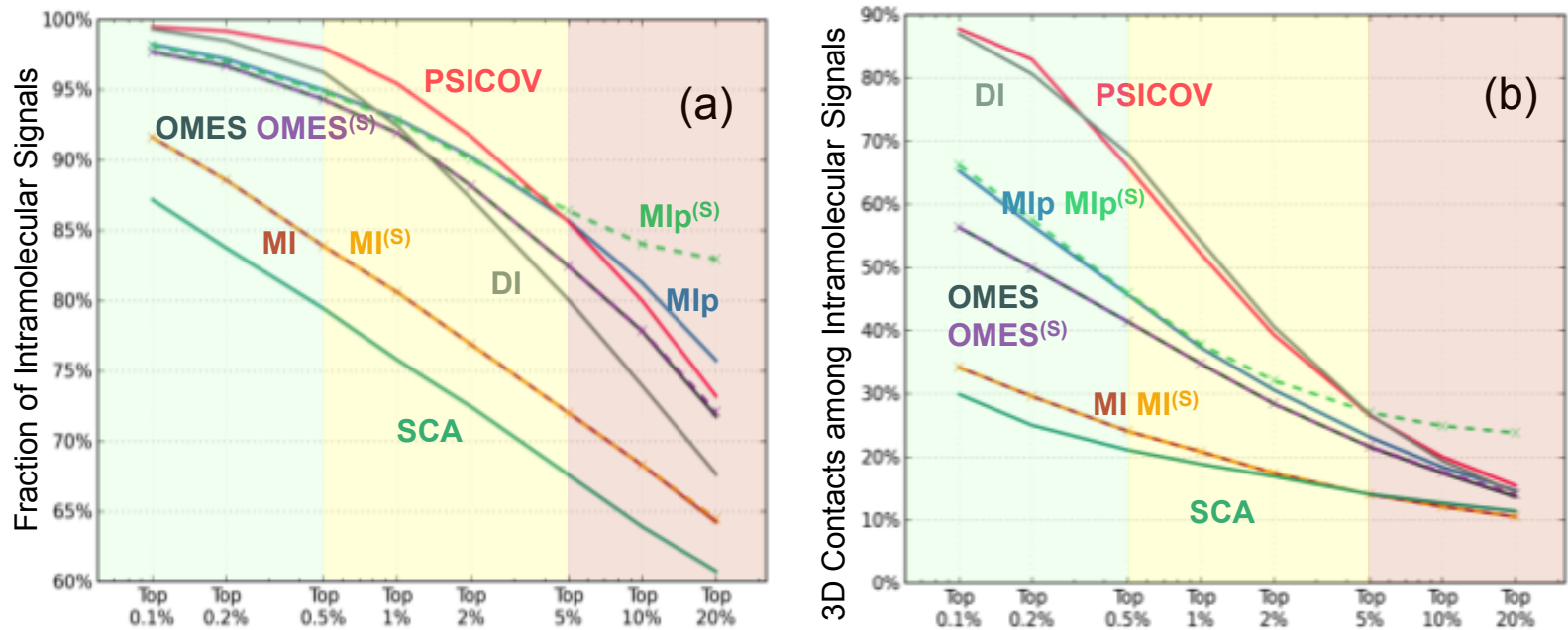
- Maximizing true positives (signals between contact making residues

# Screening of large databases

## For testing 9 methods, including

- observed-minus-expected-squared ((OMES) (Kass and Horovitz, 2002)
- statistical coupling analysis (SCA) (Halabi et al., 2009; Lockless and Ranganathan, 1999).
- Direct Coupling Analysis (DCA or DI for Direct Information) (Morcos et al., 2011; Weigt et al., 2009),
- Protein Sparse Inverse COVariance (PSICOV) (Jones et al., 2012),

# PSICOV and DI are the best



**Average performance of the nine methods based on two criteria, absence of intermolecular FPs (a), and fraction of 3D contact making pairs (b) among different subsets of top-ranking signals.** The signals are classified to 3 groups: strong coevolution signals (0.1-0.5%), intermediate signals (0.5-5%) and relatively weak signals (5-20%), which also refer to relatively small, intermediate, and high coverage of coevolving pairs. PSICOV and DI outperform other methods in the strong coevolution region. For the intermediate signal, OMES and MIp exhibit performances similar to PSICOV and DI in panel **a**. MIp[S] shows the best performance in the weak signal regime. SCA and MI (and its shuffled version) have lower performance compared to all others for both criteria over the whole range.

# Allosteric communication mechanisms explored by network models

- Diffusion of signal obeys a Markov process

- The structure is modeled as a network

- Network connectivity given by $\Gamma$

References

Laplacian based manifold methods (Belkin & Niyogi)

Chennubhotla & Bahar Mol Systems Biology (2006); PLoS Comp Bio (2007)

# Markov Model of Network Communication

$\Gamma = \mathbf{D} - \mathbf{A}$ where $\mathbf{A}$ = connectivity/affinity matrix and $\mathbf{D}$ = diagonal matrix of degrees

*A discrete-time, discrete-state* Markov process is defined by setting the conditional probability of signal transduction from residue *j* to *i* as

$$m_{ij} = a_{ij} / d_j$$

The conditional probability matrix $\mathbf{M} = \{m_{ij}\}$, also called the Markov transition matrix, is

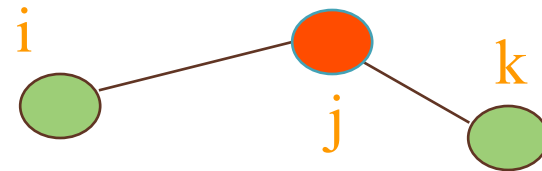$$\mathbf{M} = \mathbf{A} \, \mathbf{D}^{-1}$$

$\mathbf{M}$ completely defines the stochastics of information transfer over the network of residues.

# Hitting time: a measure of communication efficiency between two endpoints



Based on all possible pathways

| path | # of steps | Path Probability |
|---|---|---|
| $j \to i$ | 1 | 0.5 |
| $j \to k \to j \to i$ | 3 | $0.5^2$ |
| $j \to k \to j \to k \to j \to i$ | 5 | $0.5^3$ |

$$H(j,i) = 1 \times 0.5 + 3 \times 0.5^2 + \cdots = \sum_{j=1}^{\infty} (2j-1) \times 0.5^j, = 3.$$

| path | # of steps | Path Probability |
|---|---|---|
| $i \to j \to k$ | 2 | 0.5 |
| $i \to j \to i \to j \to k$ | 4 | $0.5^2$ |
| $i \to j \to i \to j \to i \to j \to k$ | 6 | $0.5^3$ |

$$H(k,i) = 2 \times 0.5 + 4 \times 0.5^2 + \cdots = 2\sum_{j=1}^{\infty} j \times 0.5^j = 4.$$

P(t) = M P(0), where M = AD$^{-1}$ is the conditional prob matrix for signal diffusion

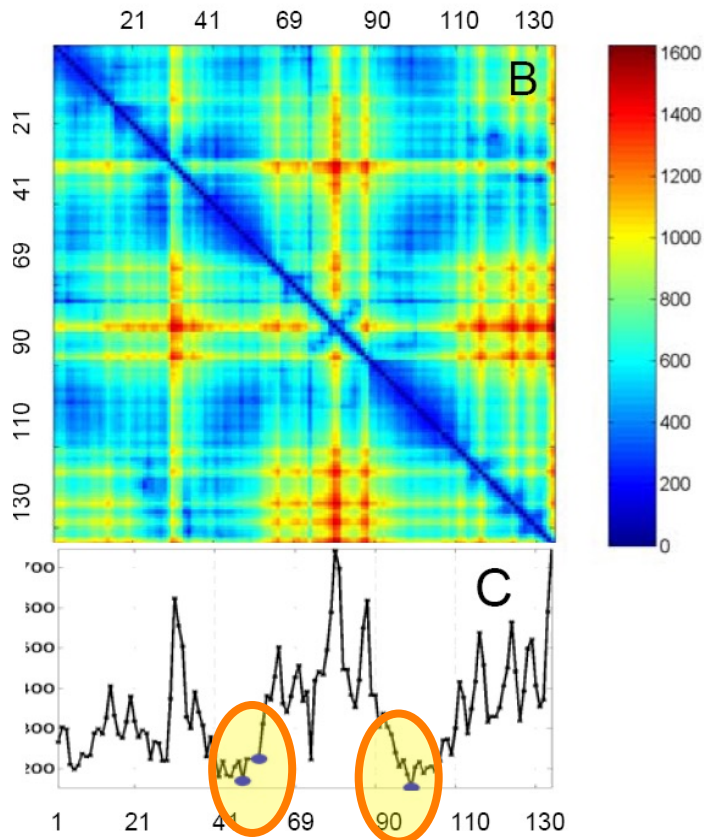# Fluctuations as determinant of communication



$$H(n, i) = 1 + \sum_{k=1}^{n-1} H(n, k) m_{ki}$$

$$H(j, i) = \sum_{k=1}^{n} \left[ \Gamma_{ki}^{-1} - \Gamma_{ji}^{-1} - \Gamma_{kj}^{-1} + \Gamma_{jj}^{-1} \right] d_k$$

$$C(i, j) = \left[ \Gamma_{ii}^{-1} + \Gamma_{jj}^{-1} - 2\Gamma_{ij}^{-1} \right] \sum_{k=1}^{n} d_k.$$

Commute distance ~ $<(\Delta R_{ij})^2>$

Chennubhotla & Bahar (2007) *PLoS Comp Bio*

# Communication times



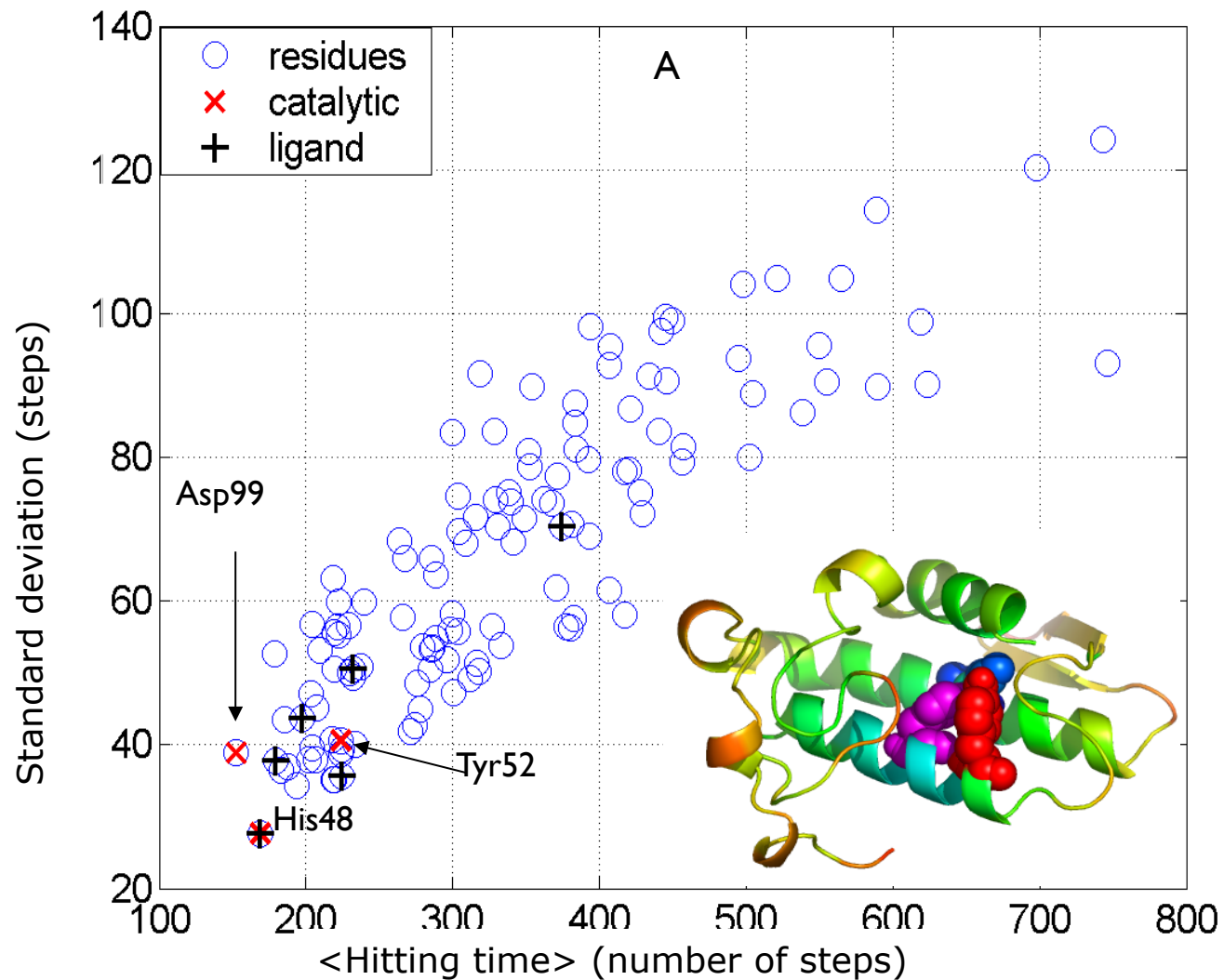Distribution of Commute Times for Phospholipase A2 (1bk9)
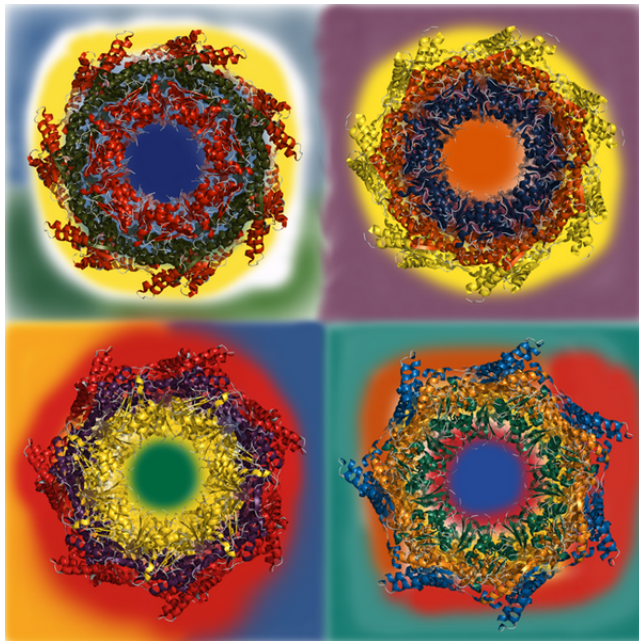
His48 ,Tyr52, Asp99 – catalytic residues

See also

Nadler, Lafon, Kevrekidis & Coifman (2005) Diffusion Maps, Spectral Clustering and Eigenfunctions of Fokker-Planck Operators, NIPS 18; Coifman et al (2005) PNAS 102, 7426.

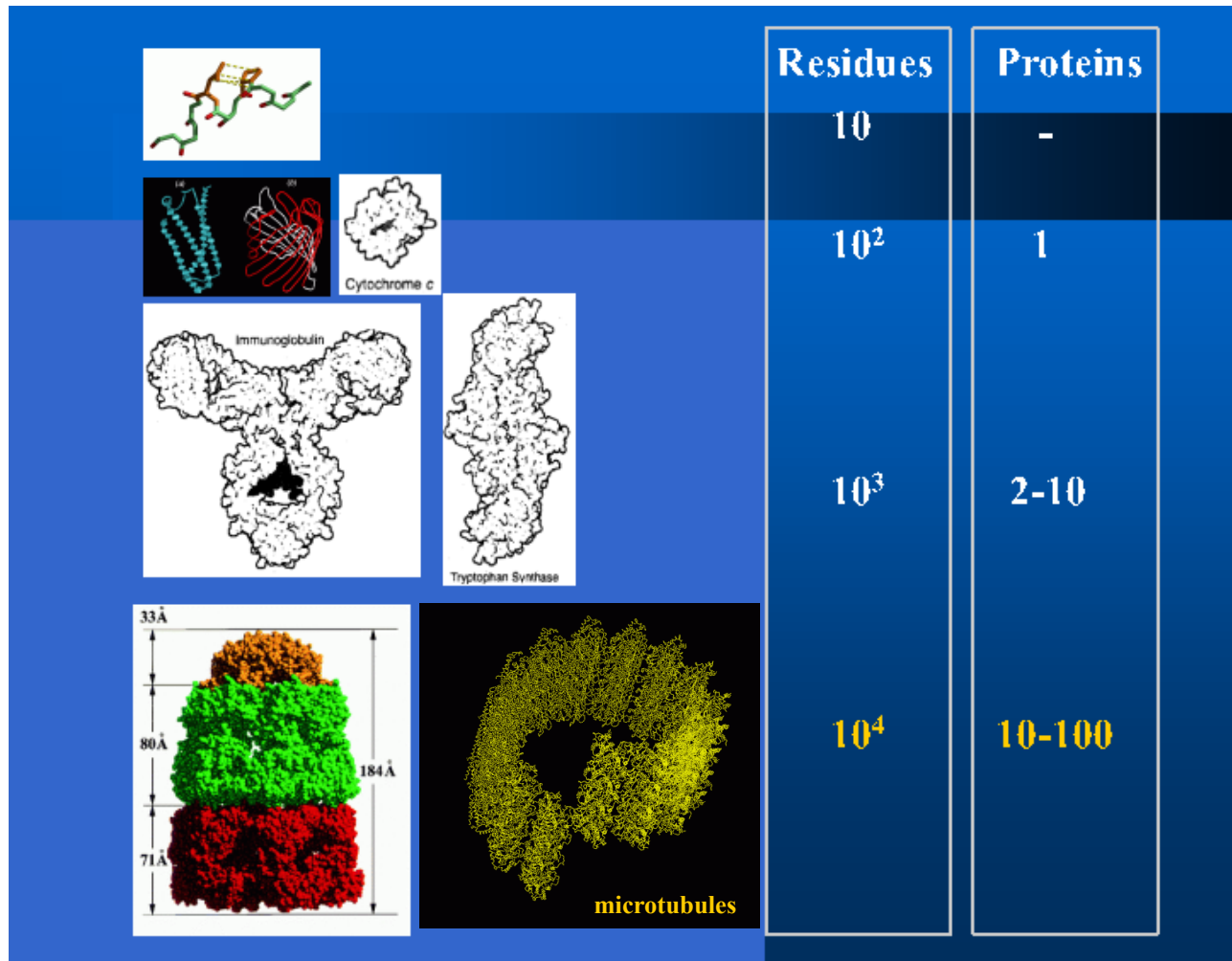# Active sites are distinguished by effective communication properties

# CONCLUSION



- Proteins are designed to favor functional changes in their structure. Pre-existing soft modes facilitate substrate binding.

- Collective mechanics/allosteric dynamics are mediated by conserved residues

- The intrinsic motions confer enhanced flexibility at substrate recognition sites

- Correlated mutations at recognition sites enable substrate specificity while conferring conformational adaptability

- Accurate modeling of protein dynamics is essential to assessing target druggability

# DISCUSSION

● Different tools for different time/length windows: MD cannot explore long-time processes for multimeric systems; ANM does not incorporate detailed atomic forces

● Not all evolutionarily correlated sites refer to structural or dynamic correlations

● Accurate modeling of protein dynamics is essential to computer-aided drug discovery, but not sufficient for quantitative evaluation of binding affinity

● Druggability simulations identify druggable sites, but not the type of drugs that optimally bind those sites

ProDy
Protein Dynamics Analysis in Python

NMWiz

Dr. Indira Shrivastava
Assist Prof, DCSB, Pitt

Markus Dittrich, PhD
NRBSC Group Leader
Pittt Supercomputing Center

Dr. Timothy R Lezon
Assistant Prof, DCSB, Pitt

Drs. Ahmet Bakan and Anindita Dutta

Dr. Chakra Chennubhotla
Assist Prof, DCSB, Pitt