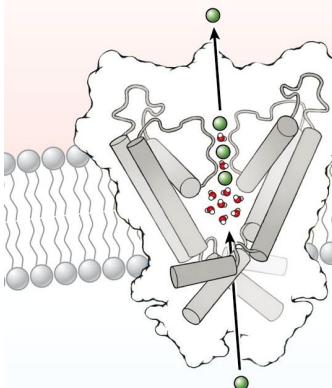


Lectures 21-22: Protein structure

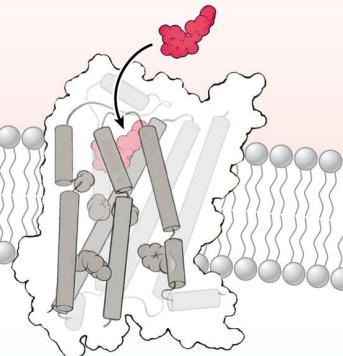
- Molecular dynamics
- Amino-acid coevolution
 - Mutual information
 - Maximum entropy modeling

Molecular dynamics (MD) simulations = Computational microscope

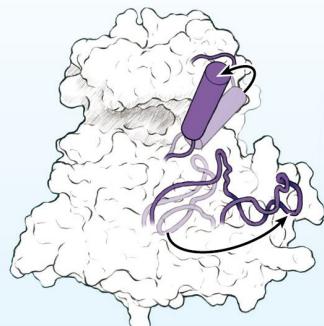
a Transport across membrane



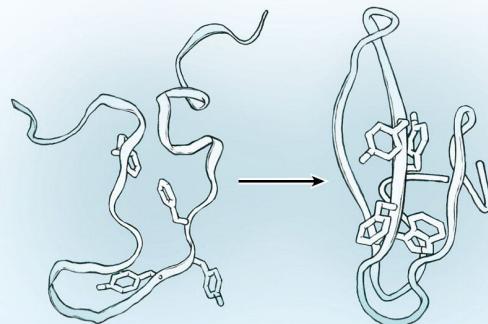
b Ligand binding



c Conformational change



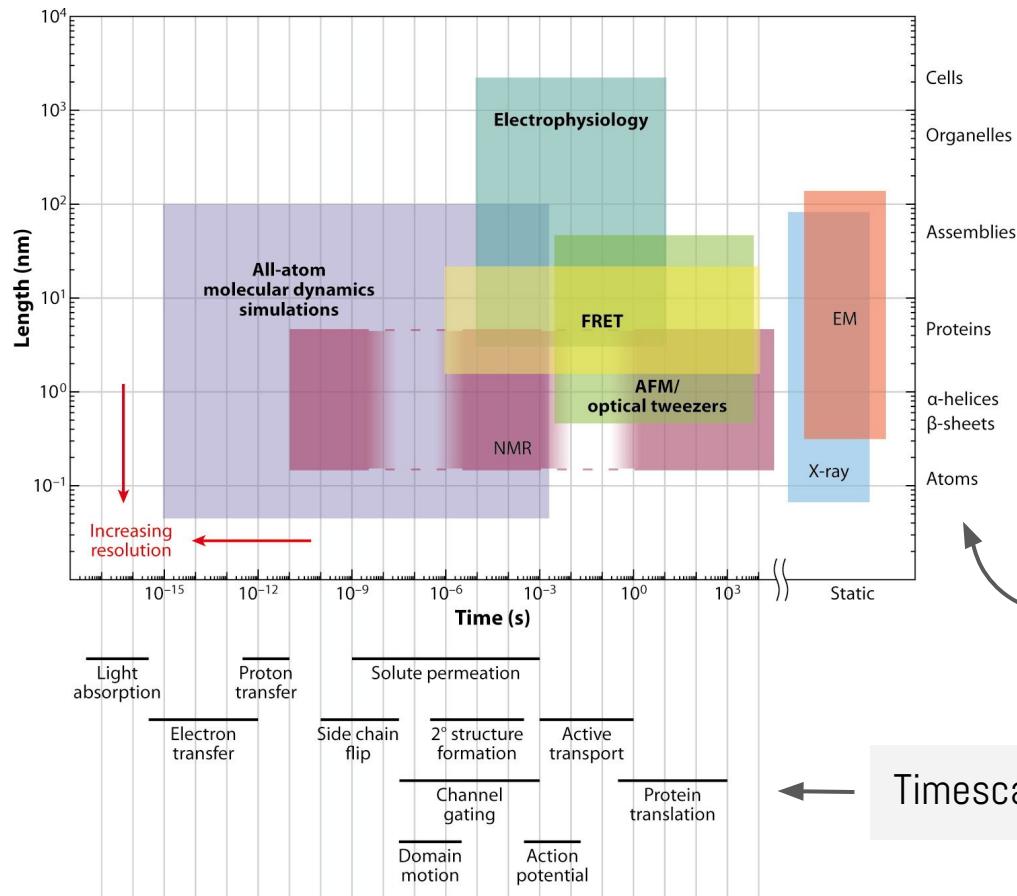
d Protein folding



MD simulations reveal the workings of biomolecular systems at a spatial and temporal resolution that is often difficult to access experimentally.

- Positions and velocities of atoms are computed using Newton's laws of motion.

Spatiotemporal resolution of various techniques



Data on single molecules (as opposed to only on ensembles) are in boldface.

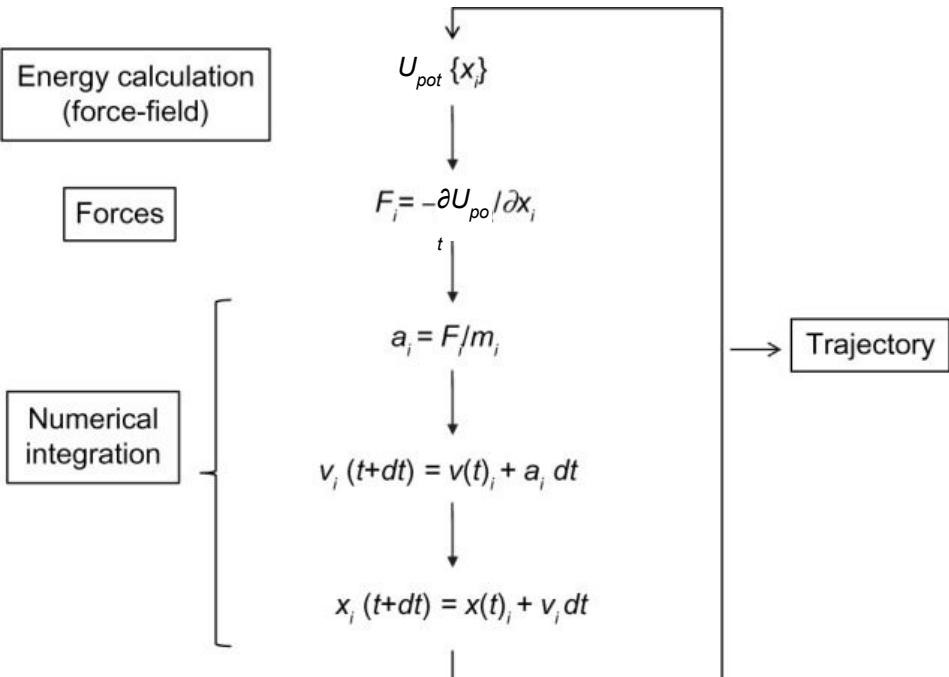
- AFM, atomic force microscopy
- EM, electron microscopy
- FRET, Forster resonance energy transfer
- NMR, nuclear magnetic resonance

Spatial resolution of biological features

Timescales of molecular processes

Molecular dynamics (MD) simulations = Computational microscope

The basic MD algorithm.



The simulation output – the trajectory – is an ordered list of $3N$ atom coordinates for each simulation time (or snapshot).

U_{pot} : potential energy

t : simulation time

dt : iteration time

For each spatial coordinate of the N simulated atoms (i):

- x : atom coordinate
- F : forces component
- a : acceleration
- m : atom mass
- v : velocity.

Force field and the energy function

The potential energy of N interacting atoms $U(\mathbf{r}_1, \dots, \mathbf{r}_N)$ is a function of their positions $\mathbf{r}_i = (x_i, y_i, z_i)$.

The force acting upon i th atom is determined by the gradient (vector of first derivatives) with respect to atomic displacements:

$$\mathbf{F}_i = -\nabla_{\mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_N) = -\left(\frac{\partial U}{\partial x_i}, \frac{\partial U}{\partial y_i}, \frac{\partial U}{\partial z_i}\right)$$

Find the positions $\mathbf{r}_i(t + \Delta t)$ at time $t + \Delta t$ in terms of the already known positions at time t .

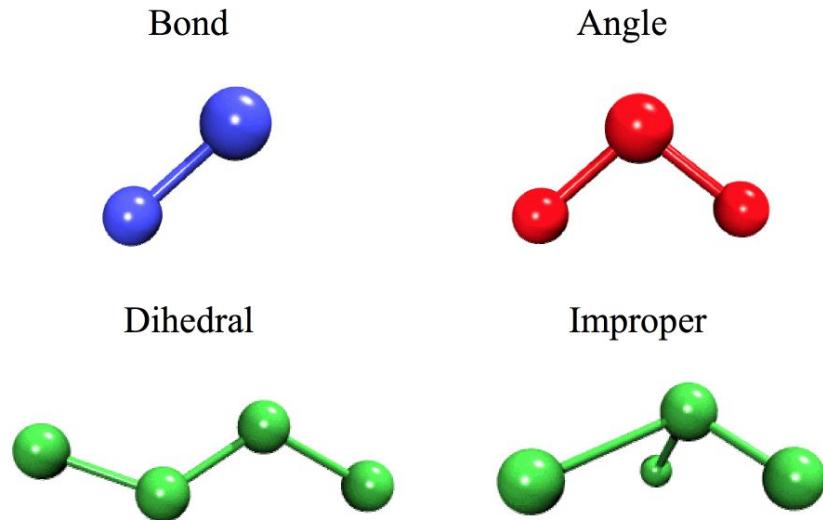
Verlet algorithm:

$$\mathbf{r}_i(t + \Delta t) \cong 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i(t)}{m_i} \Delta t^2$$

Force field: energy function used to compute the forces acting on atoms (due to interatomic interactions) during an MD simulation.

Force field and the energy function

$$U(\vec{R}) = \underbrace{\sum_{bonds} k_i^{bond} (r_i - r_0)^2}_{U_{bond}} + \underbrace{\sum_{angles} k_i^{angle} (\theta_i - \theta_0)^2}_{U_{angle}} + \underbrace{\sum_{dihedrals} k_i^{dih} [1 + \cos(n_i \phi_i + \delta_i)]}_{U_{dihedral}} + \underbrace{\sum_i \sum_{j \neq i} 4\epsilon_{ij} \left[\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] + \sum_i \sum_{j \neq i} \frac{q_i q_j}{\epsilon r_{ij}}}_{U_{nonbond}}$$



U_{bond} : oscillations about the equilibrium bond length

U_{angle} : oscillations of 3 atoms about an equilibrium angle

$U_{dihedral}$: torsional rotation of 4 atoms about a central bond

$U_{nonbond}$: non-bonded energy terms (electrostatics and Lenard-Jones)

Force field: energy function used to compute the forces acting on atoms (due to interatomic interactions) during an MD simulation.

Steps in a typical MD simulation

1. Prepare molecule: Read in pdb and psf file
2. Minimization: Reconcile observed structure with force field used ($T = 0$)
3. Heating: Raise temperature of the system
4. Equilibration: Ensure system is stable
5. Dynamics: Simulate under desired conditions (NVE, NpT, etc); Collect your data
6. Analysis: Collect your data; Evaluate observables (macroscopic level properties); Or relate to single molecule experiments.

Protein Data Bank (PDB)

www.rcsb.org: 3D shapes of proteins, nucleic acids, and complex assemblies.

RCSB PDB Deposit Search Visualize Analyze Download Learn More MyPDB

138878 Biological Macromolecular Structures Enabling Breakthroughs in Research and Education

PDB-101 EMDDataBank Worldwide Protein Data Bank Foundation

Search by PDB ID, author, macromolecule, sequence, or ligands Go Advanced Search | Browse by Annotations

Facebook Twitter YouTube

Welcome

Deposit

Search

Visualize

Analyze

Download

Learn

A Structural View of Biology

This resource is powered by the Protein Data Bank archive-information about the 3D shapes of proteins, nucleic acids, and complex assemblies that helps students and researchers understand all aspects of biomedicine and agriculture, from protein synthesis to health and disease.

As a member of the wwPDB, the RCSB PDB curates and annotates PDB data.

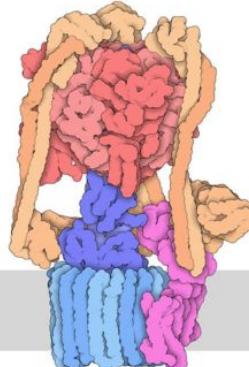
The RCSB PDB builds upon the data by creating tools and resources for research and education in molecular biology, structural biology, computational biology, and beyond.

New Video: What is a Protein?



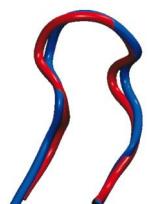
PDB-101 VIDEO WHAT IS A PROTEIN?

March Molecule of the Month

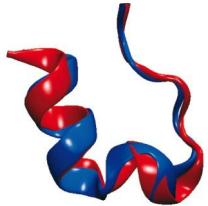


Vacuolar ATPase

Simulations of structurally diverse proteins



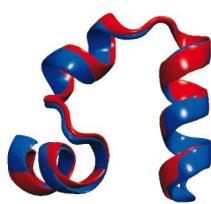
Chignolin



Trp-cage



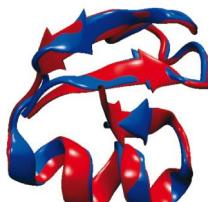
BBA



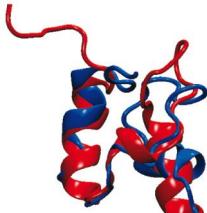
Villin



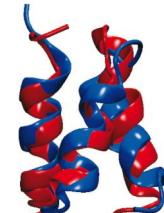
WW domain



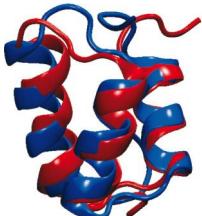
NTL9



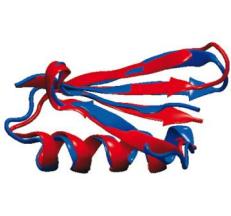
BBL



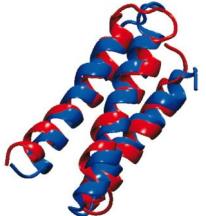
Protein B



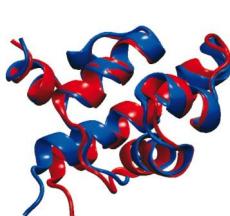
Homeodomain



Protein G



a3D



λ repressor

Simulations with a single force field.

- 12 structurally diverse proteins fold spontaneously to a structure (blue) closely resembling that determined experimentally (red).
- Simulation snapshots chosen automatically based on a clustering analysis that did not exploit knowledge of the experimental structure.
- Total simulation time per protein: 104 – 2,936 μ s – allowing observation of at least 10 folding & 10 unfolding events for each protein.

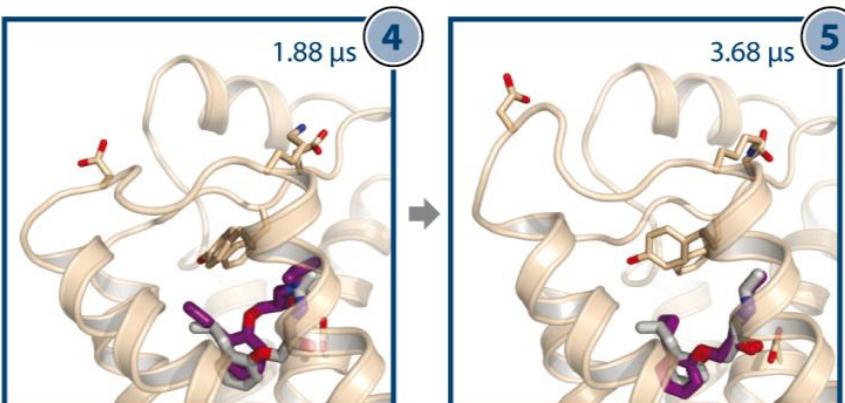
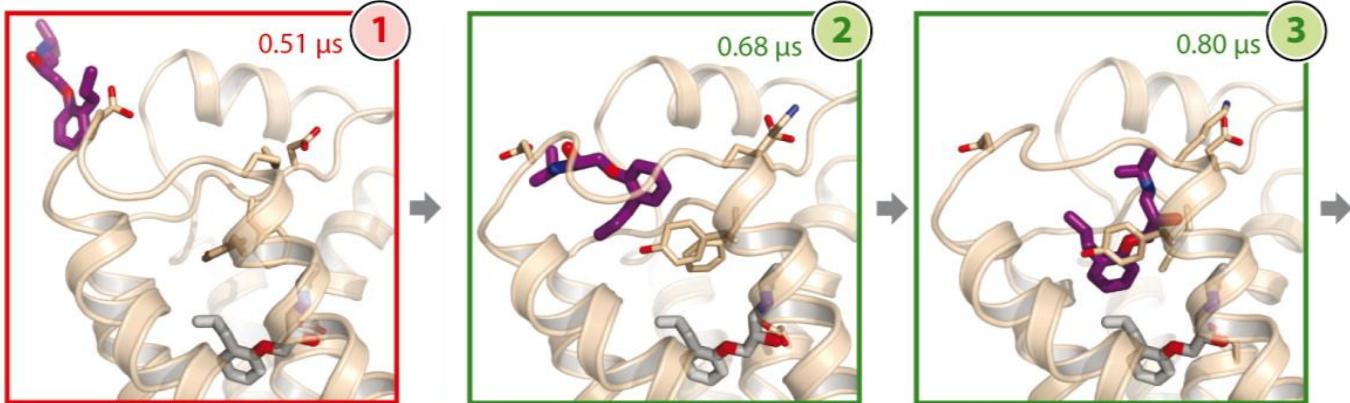
Beta-blockers binding spontaneously to the β_2 -adrenergic receptor

Metastable Intermediate stages of beta blocker binding.

1: Ligand moves from bulk solvent...

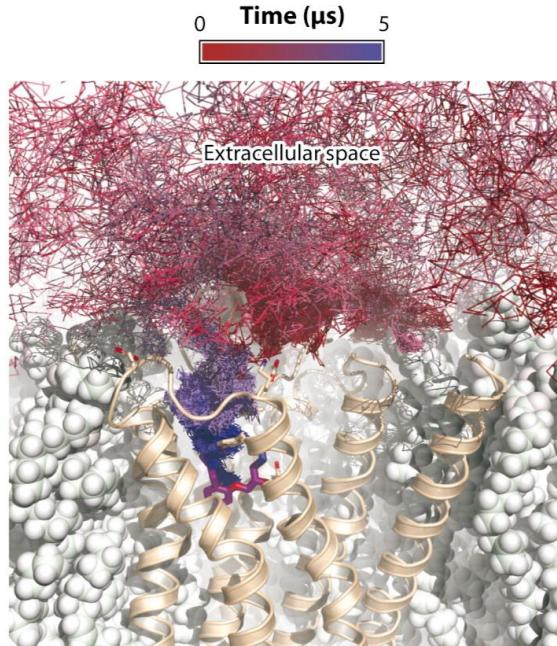
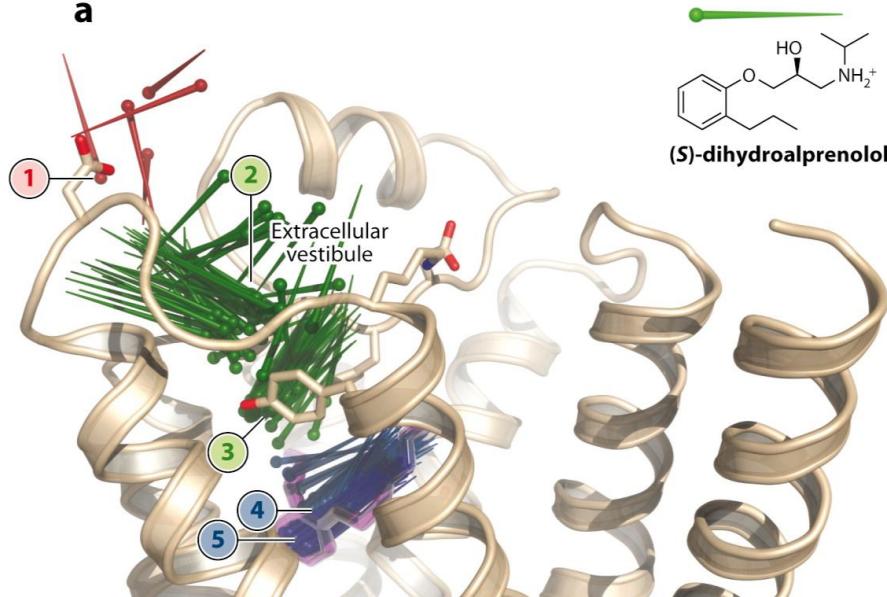
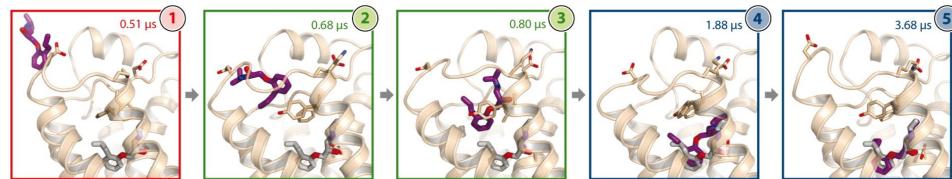
2, 3: ... into the extracellular vestibule, and finally...

4, 5: ... into the binding pocket.



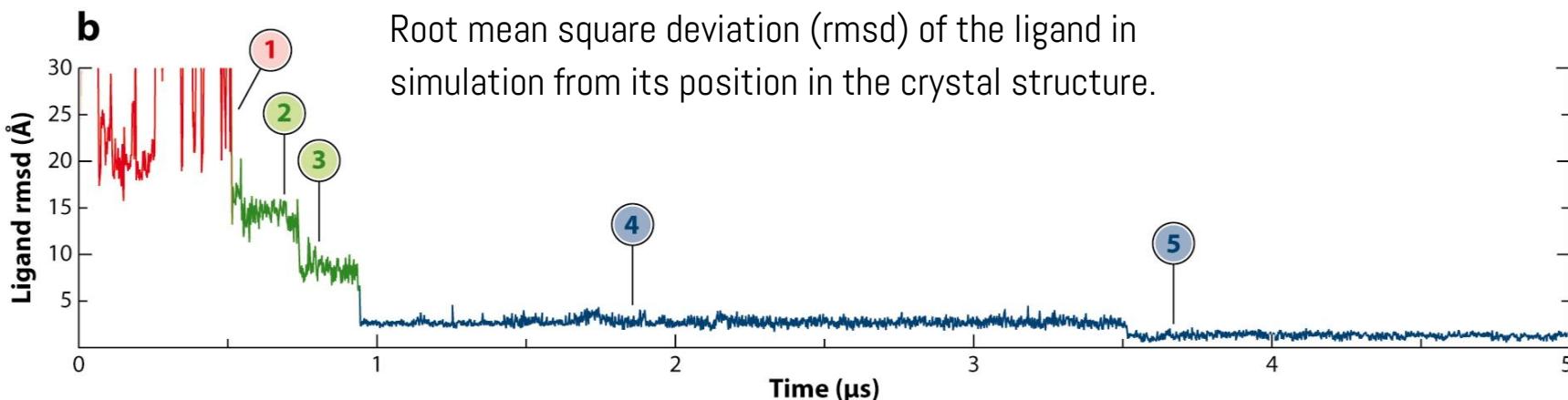
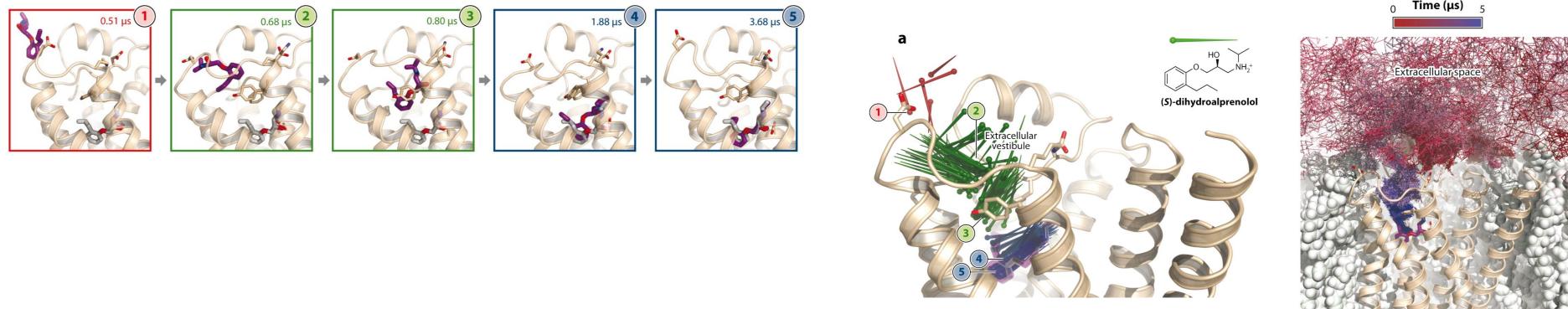
Beta blockers – aka beta-adrenergic blocking agents – reduce blood pressure by blocking the effects of epinephrine (adrenaline).

Beta-blockers binding spontaneously to the $\beta 2$ -adrenergic receptor



Pins:
successive
positions

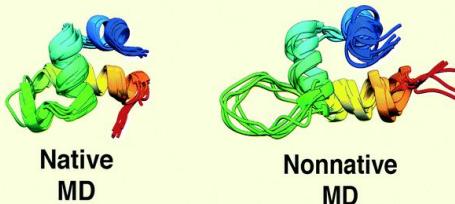
Beta-blockers binding spontaneously to the $\beta 2$ -adrenergic receptor



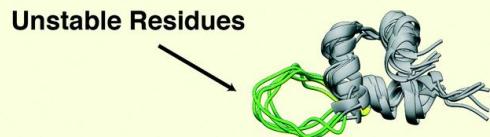
MD simulations for protein design

A. Modulating Protein Stability

- I. Perform simulations in native and nonnative environments



- II. Inform designs by analyzing the dynamics of unstable residues



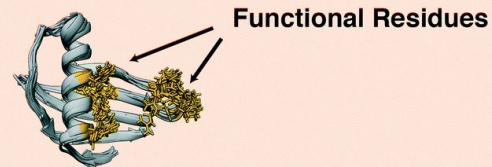
- III. Perform simulations of designs to determine the impact of mutations on protein stability

B. Engineering Functional Regions

- I. Perform simulations that capture functional dynamics



- II. Inform designs by analyzing the dynamics of functional residues

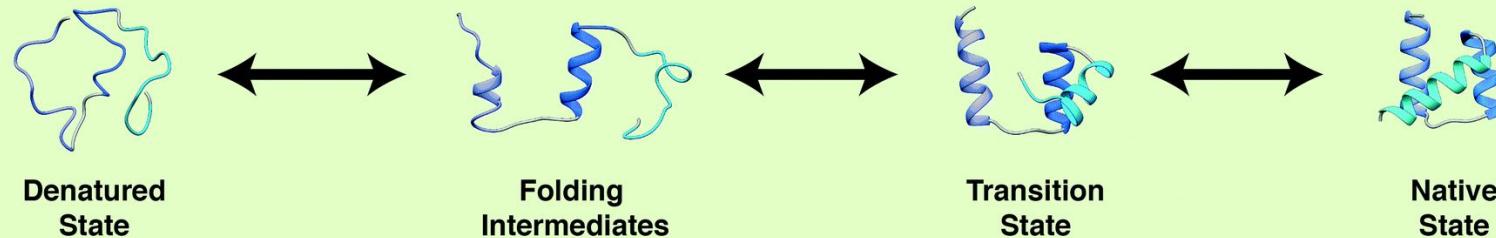


- III. Perform simulations of designs to assess the impact of mutations on function regions

MD simulations for protein design

C. Insights From Folding Pathways

I. Simulate the unfolding/folding pathway and partition the trajectory into conformational states



II. Inform designs with insights from specific conformations or transitions along the folding pathway

Destabilize the
Denatured State

Stabilize
Folding Intermediates

Design Fast
Folding Variants

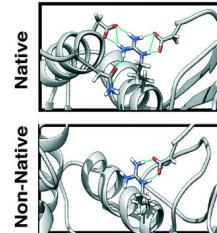
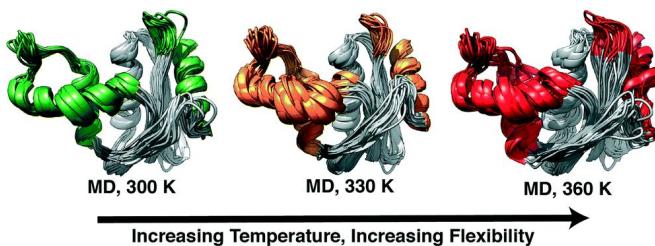
Probe / Alter the
Folding Pathway

III. Perform simulations of designs to assess the impact of mutations on the folding landscape

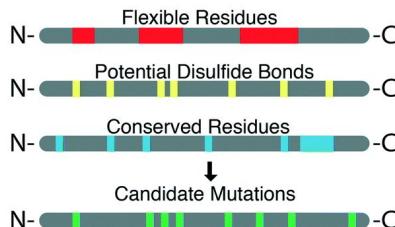
MD simulations for protein design

Direct

1. Simulate the design target over a range of temperatures; then analyze the simulations to resolve atomistic details of flexible sites.



2. Select mutations by combining insights from MD simulations with available data.



3. Assess design(s) computationally,



iteratively design,
simulate, and assess

4. Evaluate the design experimentally

Indirect

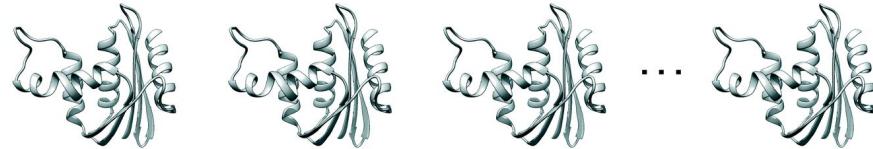
1. Construct multiple designs using available data

Design 1

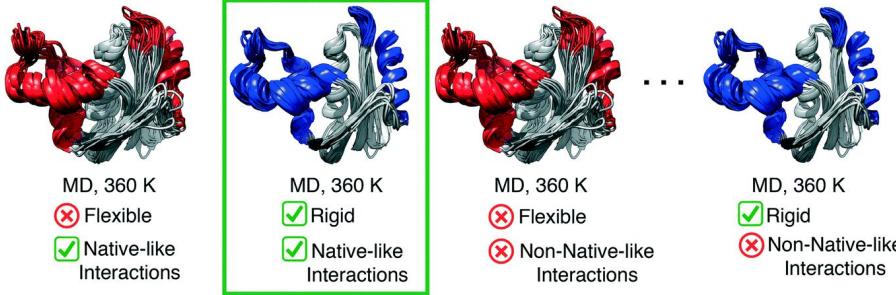
Design 2

Design 3

Design N



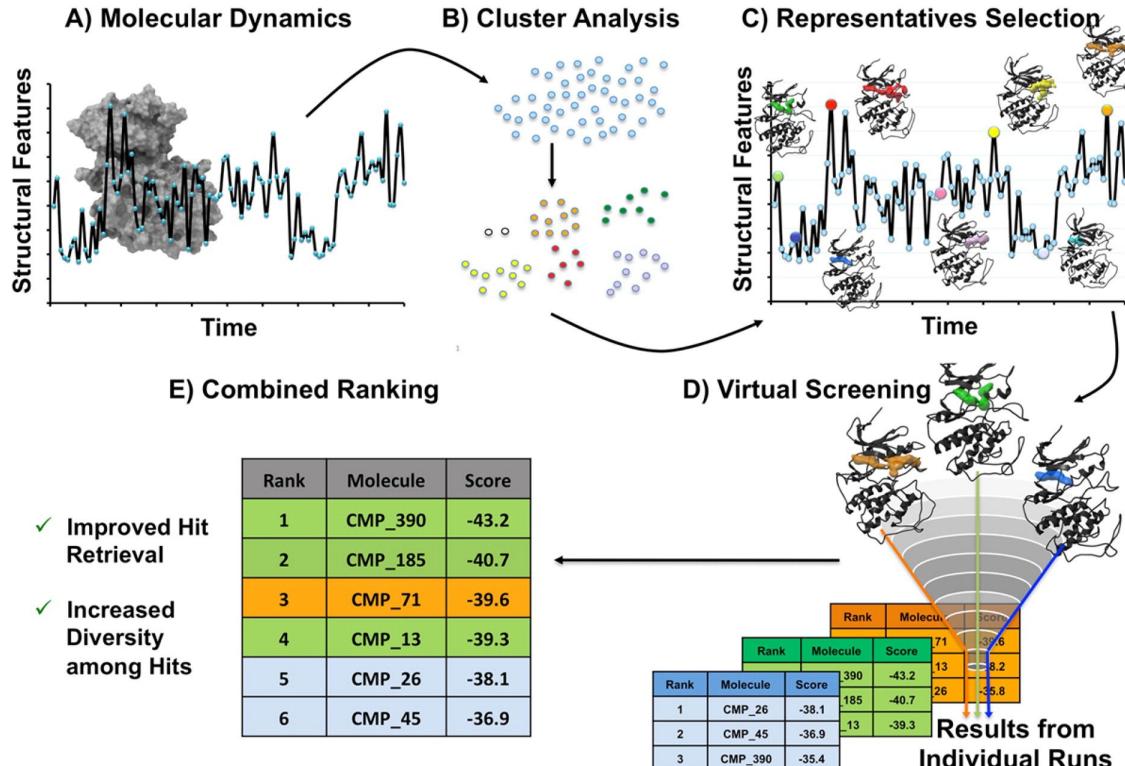
2. Perform simulations of designs at high temperature and rank designs according to their flexibility and other structural or dynamic criteria



3. Evaluate the top scoring design(s) experimentally

Virtual screening: docking & MD simulations

- (A) Use MD trajectory to explore the receptor conformational space.
- (B) Extract several snapshots from the trajectory; clustering to eliminate redundancy.
- (C) From each cluster, select a representative structure (e.g., medoid).
- (D) Carry out virtual ligand screening independently at each representative conformation.
- (E) Return activity predictions by independent runs and combine together in a global ranking.



Distributed computing & Crowdsourcing

Folding@home: folding.stanford.edu

- Distributed computing project for MD simulations (e.g., protein folding, computational drug design).
- Uses the idle resources of personal computers owned by volunteers from all over the world.

Foldit: fold.it

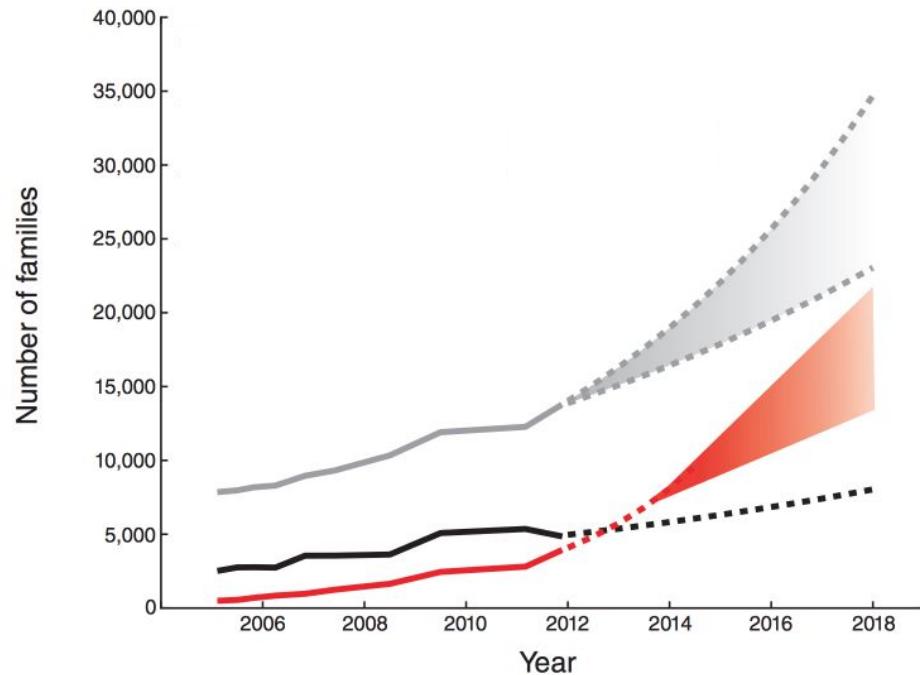
- An online game that poses complex puzzles about how proteins fold.
- Helped solve the structure of a protein-sniping enzyme critical for reproduction of the AIDS virus within 3 weeks; Identified targets for drugs to neutralize it.



Lectures 21-22: Protein structure

- Molecular dynamics
- Amino-acid coevolution
 - Mutual information
 - Maximum entropy modeling

Direct vs. indirect interactions



New protein families being discovered by high-throughput sequencing

Experimental structure-determination

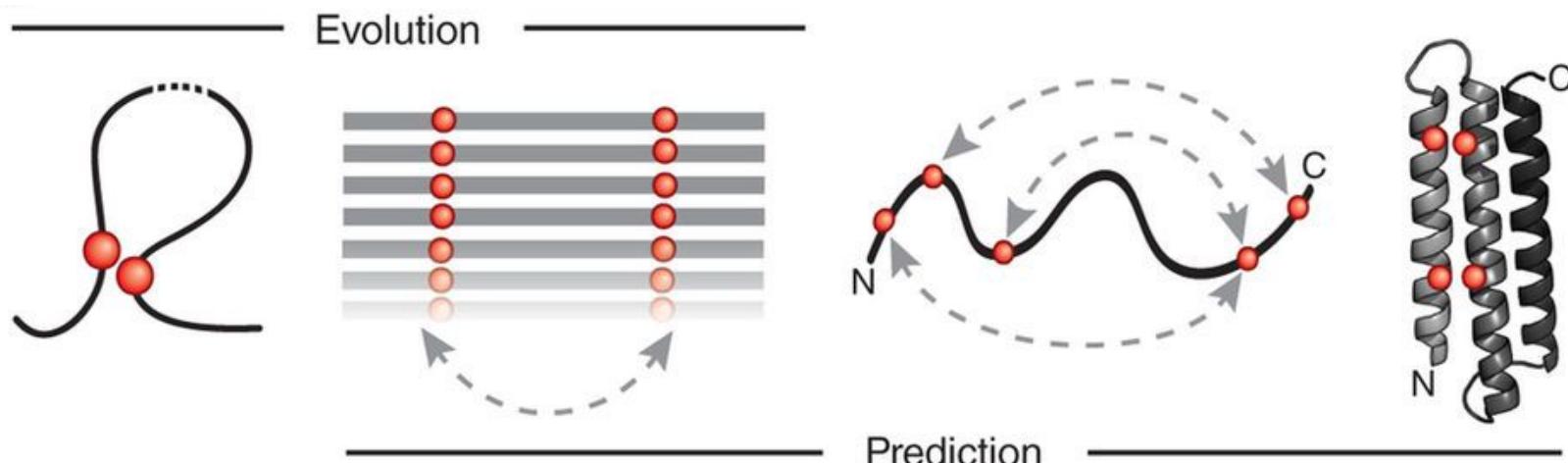
Predicting protein 3D structure from sequence

Evolutionary pressure to maintain favorable interactions b/w physically interacting AA residues in 3D.

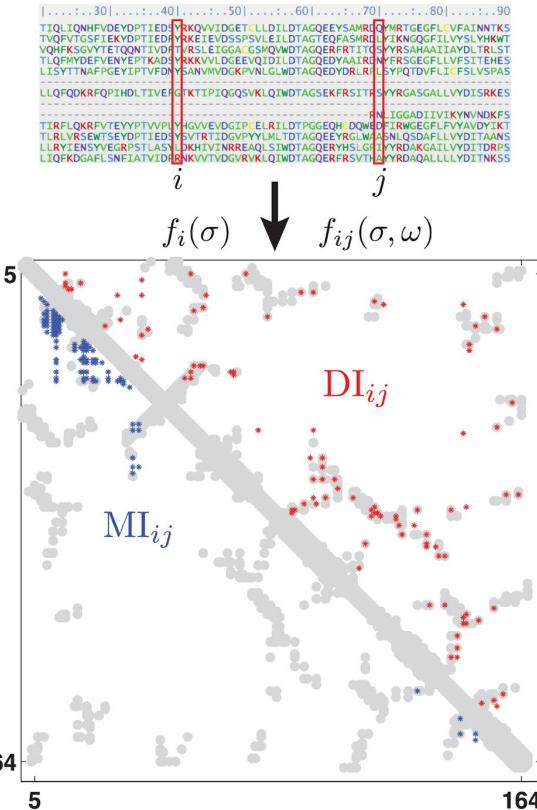
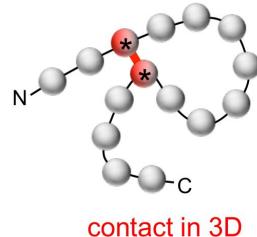
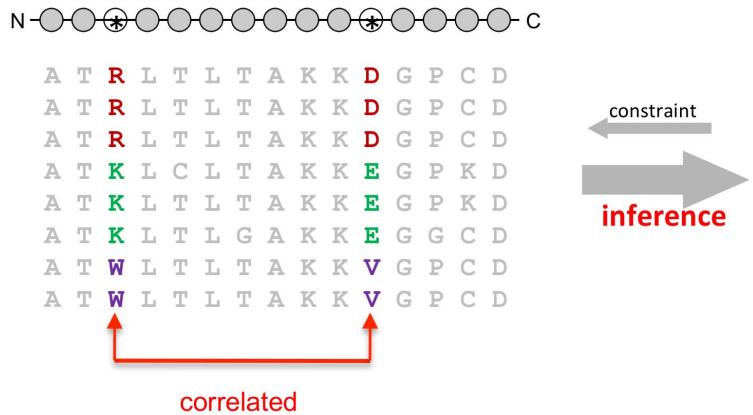
Visible record of residue covariation in related protein sequences.

Inverse problem – inferring directly causative residue couplings (evolutionary couplings) from the covariation record – challenging due to transitive correlations & other confounding effects.

ECs can be used to predict the unknown 3D structure of a protein from a set of sequences alone.

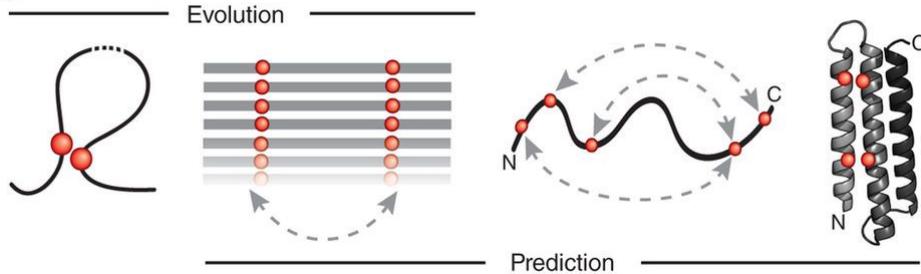


Predicting protein 3D structure from sequence



Marks (2011) PLoS One; Marks (2012) Nat. Biotech.
Stein (2015) PLoS Comp. Biol.

Predicting protein 3D structure from sequence



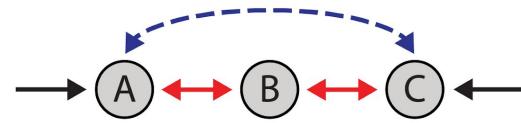
Growth in sequence databases from massively parallel sequencing.

- Availability of sufficient sequences of sufficient diversity.
- Known protein families are growing in size from a few sequences to many thousands of sequences (advances in DNA sequencing tech).

Reduction of conformational search space by cooperative probability models.

- Global probability models account for the fact that interactions along an entire protein chain are mutually interdependent in a way that is inherently cooperative.
- Pair interactions are modified by interactions with other parts of the system and cannot be factored (probabilities are not a simple product of independent terms).
- Compared with molecular dynamics simulations, statistical approaches are many orders of magnitude more efficient in reducing a huge conformational search space to manageable proportions.

Direct vs. indirect interactions

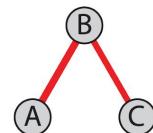
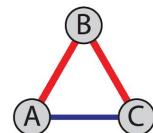


Correlation

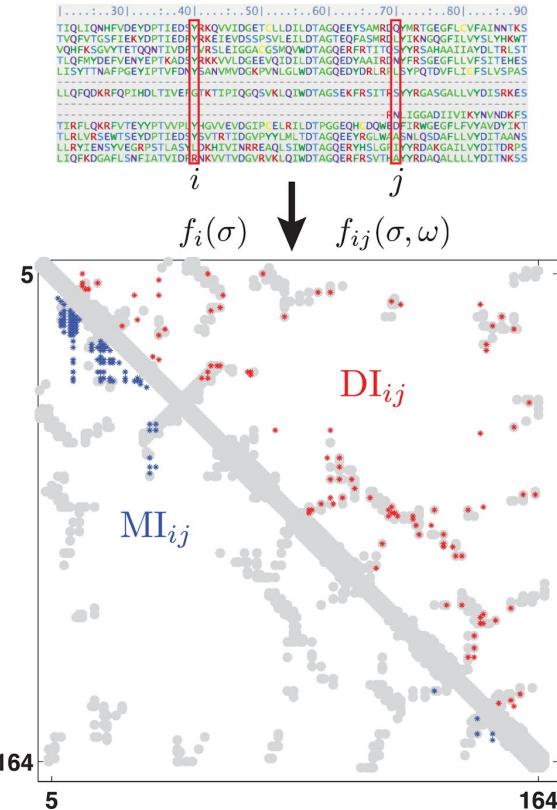
	A	B	C
A	white	red	blue
B	red	white	red
C	blue	red	white

Partial correlation

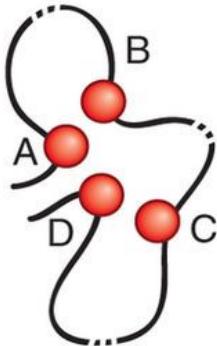
	A	B	C
A	white	red	white
B	red	white	red
C	white	red	white



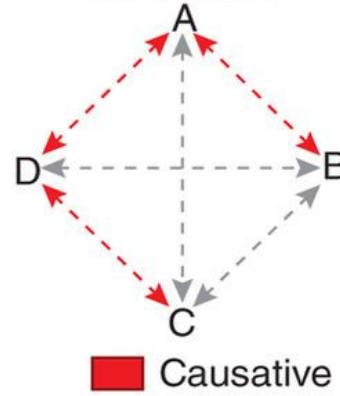
Direct vs. indirect interactions



Physical contacts



Observed correlations



Predicted contacts

A	B	C	D
A			
B	■		
C		■	
D	■		■

■ Causative ■ Transitive

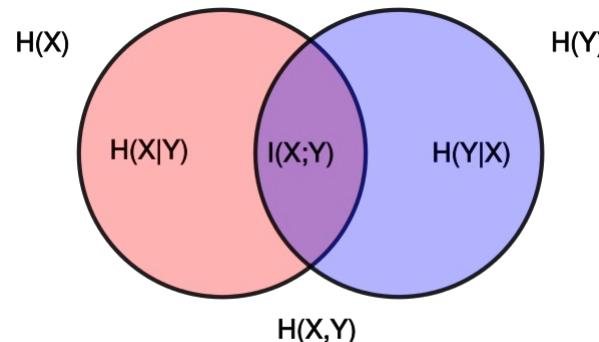
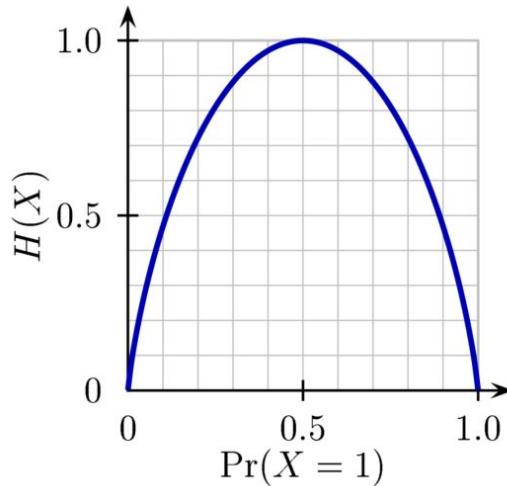
Information theory

Entropy (H): the average amount of information produced by a stochastic source of data.

Mutual information: MI two random variables $I(X, Y)$ quantifies the amount of information obtained about one random variable, through the other random variable.

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

$$I(X; Y) = \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} = H(Y) - H(Y|X)$$



Global probabilistic models of residue coupling (maximum-entropy)

$$\mathbf{x} = (x_1, \dots, x_L) \in \Omega^L$$



Pairwise maximum-entropy distribution

$$P(x_1, \dots, x_L) = \frac{1}{Z} \exp \left(\sum_i h_i(x_i) + \sum_{i < j} e_{ij}(x_i, x_j) \right)$$

Parameter inference

- mean-field (MF)

$$e_{ij}^{\text{MF}}(\sigma, \omega) = -\left(C^{-1}\right)_{ij}(\sigma, \omega)$$

- sparse maximum-likelihood (SML)

$$e_{ij}^{\text{SML}}(\sigma, \omega) = -\left(C_{1,\lambda}^{-1}\right)_{ij}(\sigma, \omega)$$

- pseudolikelihood maximization (PLM)

$$\left\{ \mathbf{h}^{\text{PLM}}(\sigma), \mathbf{e}^{\text{PLM}}(\sigma, \omega) \right\} = \arg \min_{\mathbf{h}(\sigma), \mathbf{e}(\sigma, \omega)} \left\{ -\ln l_{\text{PL}} + \lambda_h \|\mathbf{h}\|_2^2 + \lambda_e \|\mathbf{e}\|_2^2 \right\}$$



Pair scoring functions

- direct information

$$\text{DI}_{ij} = \sum_{\sigma, \omega} P_{ij}^{\text{dir}}(\sigma, \omega) \ln \left(\frac{P_{ij}^{\text{dir}}(\sigma, \omega)}{f_i(\sigma) f_j(\omega)} \right)$$

- Frobenius norm

$$\|e_{ij}\|_{\text{F}} = \left(\sum_{\sigma, \omega} e_{ij}(\sigma, \omega)^2 \right)^{1/2}$$

- average product-corrected Frobenius norm

$$\text{APC-FN}_{ij} = \|e_{ij}\|_{\text{F}} - \frac{\|e_{i\cdot}\|_{\text{F}} \|e_{\cdot j}\|_{\text{F}}}{\|e_{\cdot \cdot}\|_{\text{F}}}$$

Global probabilistic models of residue coupling (maximum-entropy)

$a = (a_1, a_2 \dots, a_N)$ A sequence made of monomers a_i taking values from a given alphabet

$$P(a|J, h) = \frac{1}{Z} \exp \left(\sum_{i=1}^{N-1} \sum_{j=i+1}^N J_{ij}(a_i, a_j) + \sum_{i=1}^N h_i(a_i) \right)$$

Probability of a sequence within the model.

$h(a_i)$: parameters that represent the propensity of symbol to be found at a certain position.

$J(a_i, a_j)$: represent an interaction, quantifying how compatible the symbols at both positions are with each other.

Global probabilistic models of residue coupling (maximum-entropy)

$$a = (a_1, a_2 \dots, a_N)$$

$$P(a|J, h) = \frac{1}{Z} \exp \left(\sum_{i=1}^{N-1} \sum_{j=i+1}^N J_{ij}(a_i, a_j) + \sum_{i=1}^N h_i(a_i) \right)$$

The idea of maximum-entropy: For a given set of sample covariances and frequencies, the model represents the **distribution with the maximal entropy** of all distributions reproducing those covariances and frequencies.

$$\begin{aligned} F[P] = & - \sum_a P(a) \log P(a) \\ & + \sum_{i < j} \sum_{x,y} \lambda_{ij}(x,y) \left(P_{ij}(x,y) - f_{ij}(x,y) \right) \\ & + \sum_i \sum_x \lambda_i(x) \left(P_i(x) - f_i(x) \right) \\ & + \Omega \left(1 - \sum_a P(a) \right). \end{aligned}$$

The unique distribution P that maximizes the functional to the *left*.

$f_i(a)$: frequency of finding symbol a at position i .

$f_{ij}(a, b)$: frequency of finding symbols a & b at positions i and j in the same sequence.

Global probabilistic models of residue coupling (maximum-entropy)

$$a = (a_1, a_2 \dots, a_N)$$

$$P(a|J, h) = \frac{1}{Z} \exp \left(\sum_{i=1}^{N-1} \sum_{j=i+1}^N J_{ij}(a_i, a_j) + \sum_{i=1}^N h_i(a_i) \right)$$

$$\begin{aligned} F[P] &= - \sum_a P(a) \log P(a) \\ &\quad + \sum_{i < j} \sum_{x,y} \lambda_{ij}(x,y) (P_{ij}(x,y) - f_{ij}(x,y)) \\ &\quad + \sum_i \sum_x \lambda_i(x) (P_i(x) - f_i(x)) \\ &\quad + \Omega \left(1 - \sum_a P(a) \right). \end{aligned}$$

$$\begin{aligned} F_i &= \frac{1}{N} \sum_{j \neq i}^N F_{ij} \\ F_{ij}^{APC} &= F_{ij} - \frac{F_i F_j}{F} \\ F &= \frac{1}{N^2 - N} \sum_{i,j,i \neq j}^N F_{ij} \end{aligned}$$

The idea of maximum-entropy: For a given set of sample covariances and frequencies, the model represents the **distribution with the maximal entropy** of all distributions reproducing those covariances and frequencies.

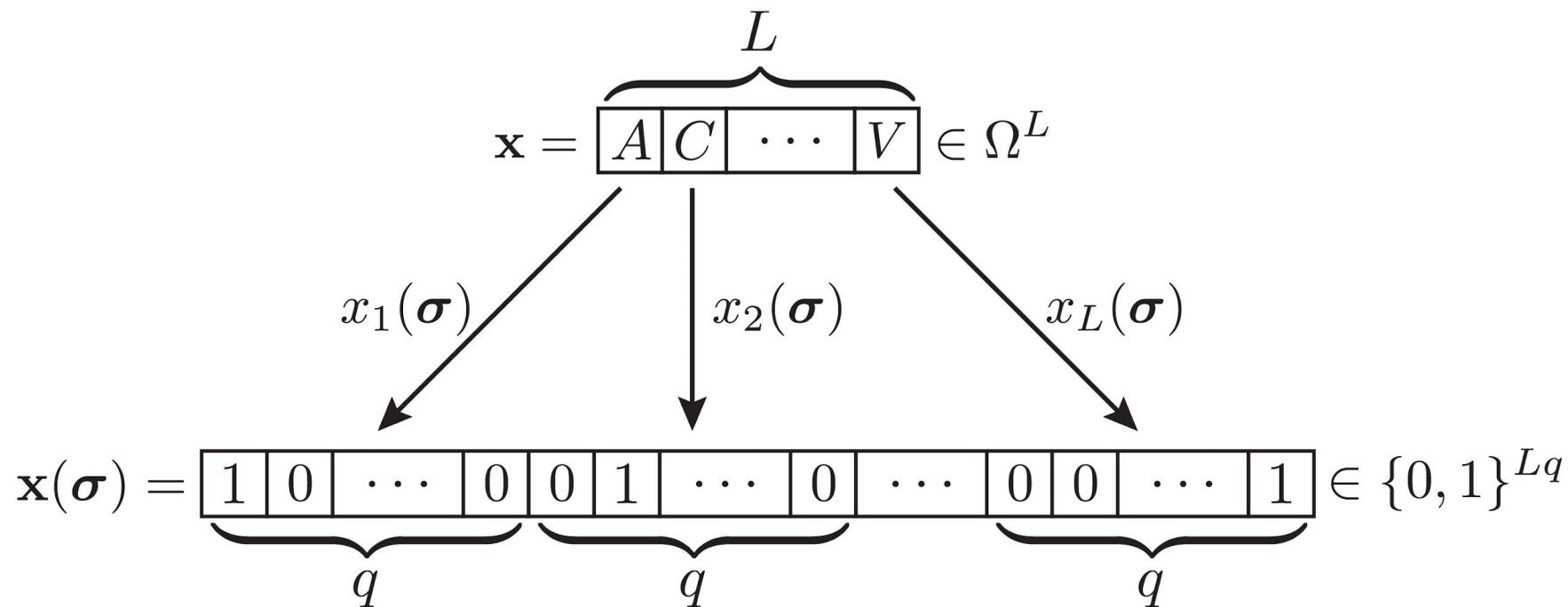
The unique distribution P that maximizes the functional to the *left*.

Final step:

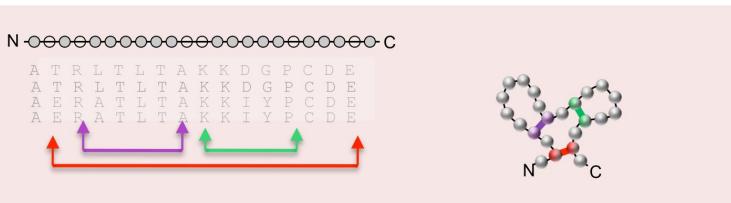
- average product correction (APC).

Global probabilistic models of residue coupling

Binary embedding of amino acid sequence



Global probabilistic models of residue coupling (maximum-entropy)



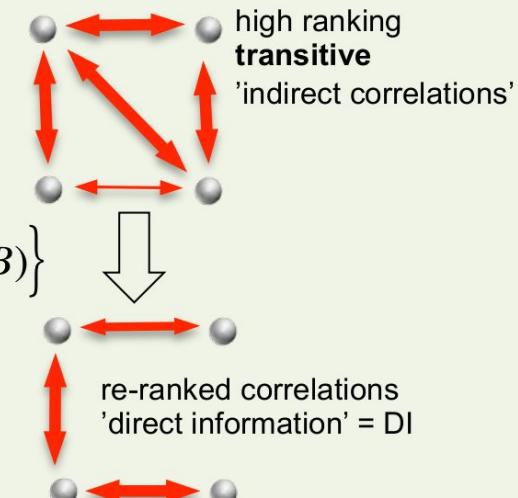
Calculate covariance matrix for each pair of sequence positions for all pairs of amino acids (A,B)

$$C_{ii}(A,B) = f_{ii}(A,B) - f_i(A)P_i(B)$$

$$C_{ij}^{-1}(A,B) = -e_{ij}(A,B)_{i \neq j}$$

$$P_{ij}^{Dir}(A,B) = \frac{1}{Z} \exp\left\{e_{ij}(A,B) + \tilde{h}_i(A) + \tilde{h}_j(B)\right\}$$

$$DI_{ij} = \sum_{A,B=1}^q P_{ij}^{Dir}(A,B) \ln \frac{P_{ij}^{Dir}(A,B)}{f_i(A)f_j(B)}$$

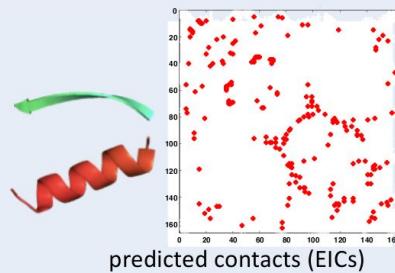


Identify maximally informative pair couplings using **statistical model** of entire protein to infer residue-residue co-evolution

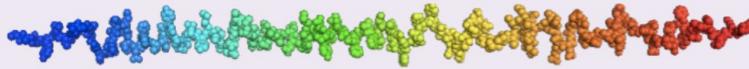
From contacts to structure

Analyze the highest scoring pairs to produce ranked list of residue pairs which we predict to be close in 3D space. Use these pairs as predicted close “evolutionary inferred contacts”, EICs, in folding calculations

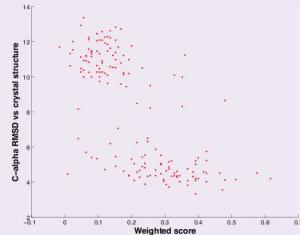
```
assign (resid 143 and name CA) (resid 123 and name CA) 4 4 3  
assign (resid 16 and name CA) (resid 10 and name CA) 4 4 3  
assign (resid 141 and name CA) (resid 82 and name CA) 4 4 3  
assign (resid 129 and name CA) (resid 87 and name CA) 4 4 3  
assign (resid 92 and name CA) (resid 11 and name CA) 4 4 3  
assign (resid 116 and name CA) (resid 81 and name CA) 4 4 3
```



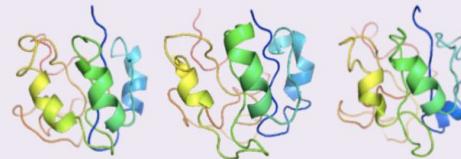
Start with extended structure
use **distance geometry** and **simulated annealing** with predicted constraints, EICs, to fold the chain



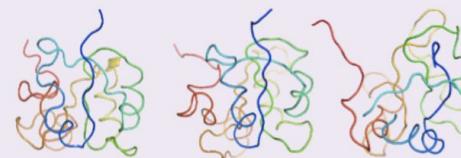
Rank predicted structures using quality measure of backbone alpha torsion and beta sheet twist



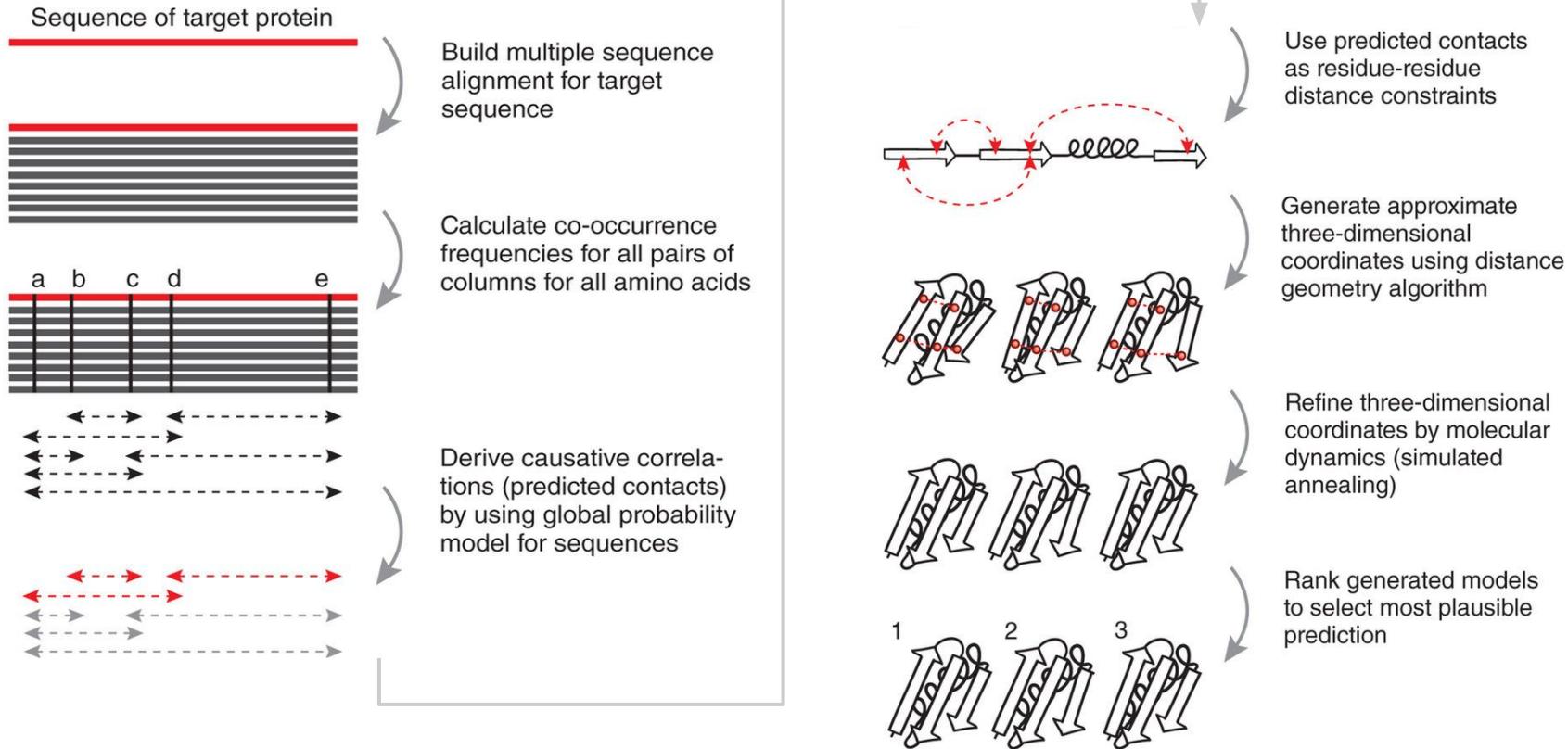
good scores



bad scores



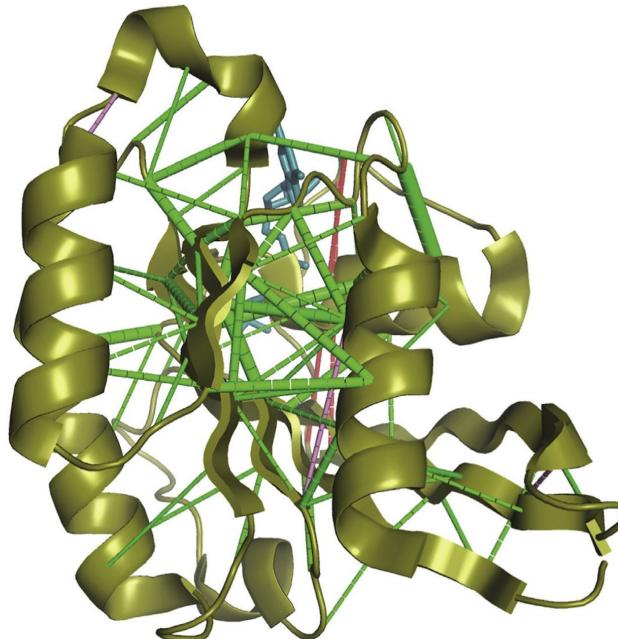
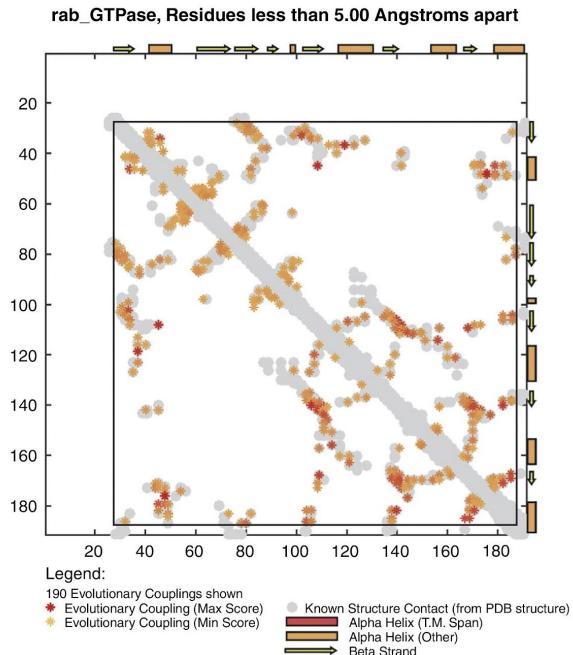
Predicting protein 3D structure from sequence



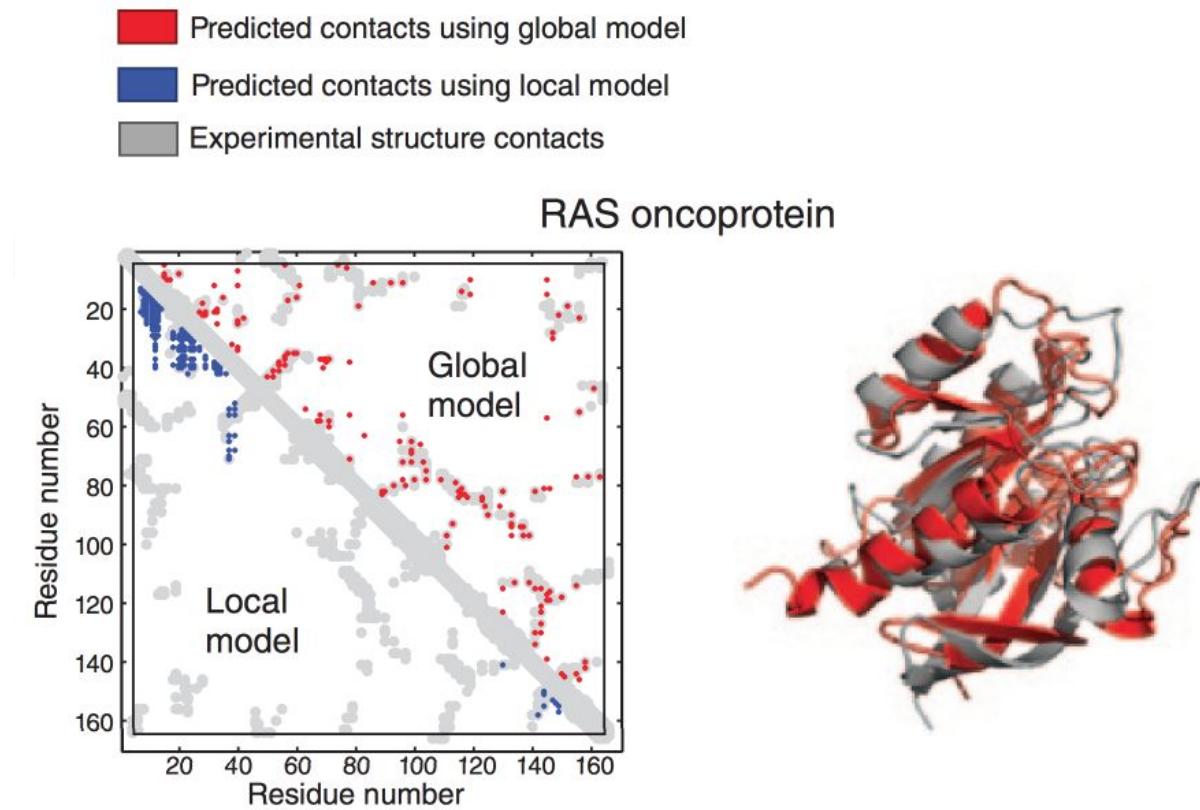
Three-dimensional structure of target protein

Marks (2012) Nat. Biotech.

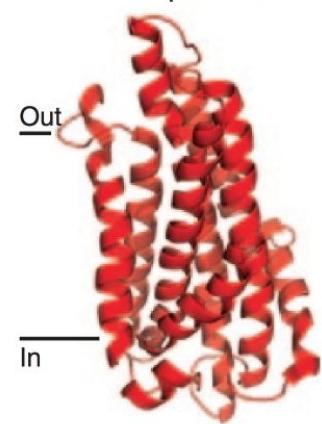
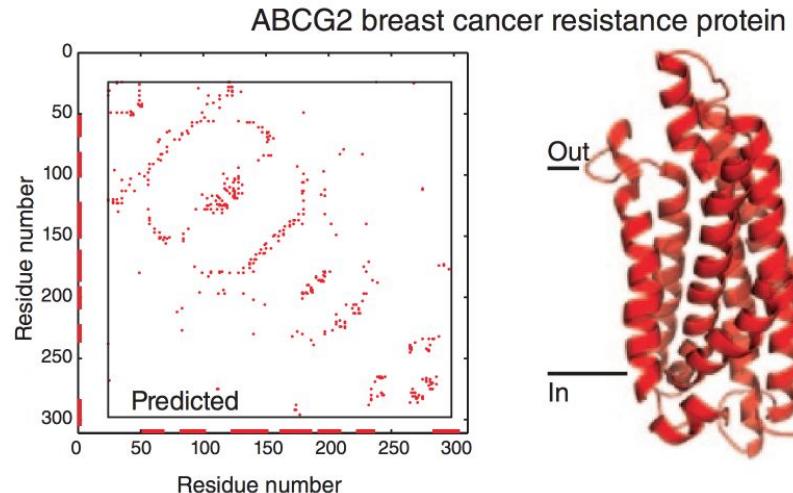
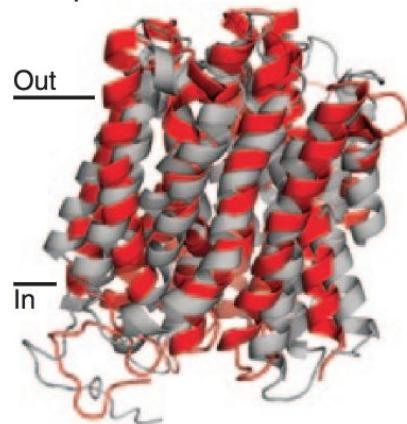
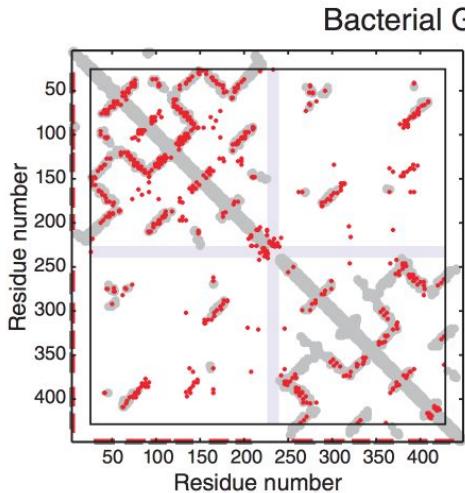
Predictions of 3D structures based on evolutionary coupling



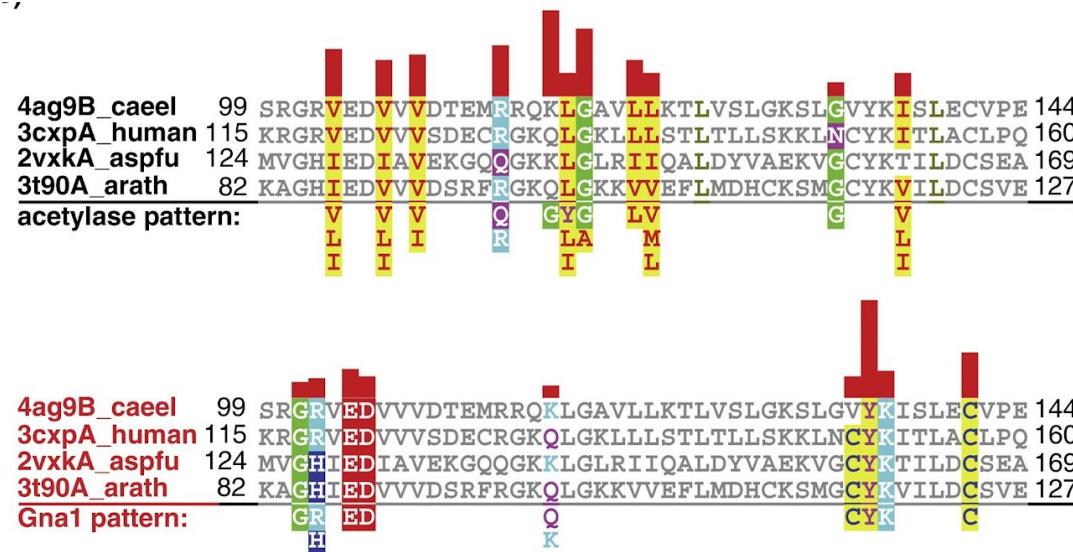
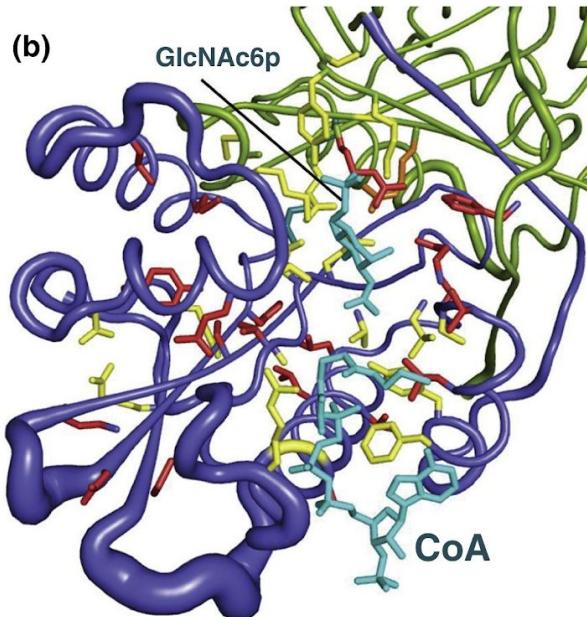
Predictions of 3D structures based on evolutionary coupling



Predictions of 3D structures based on evolutionary coupling



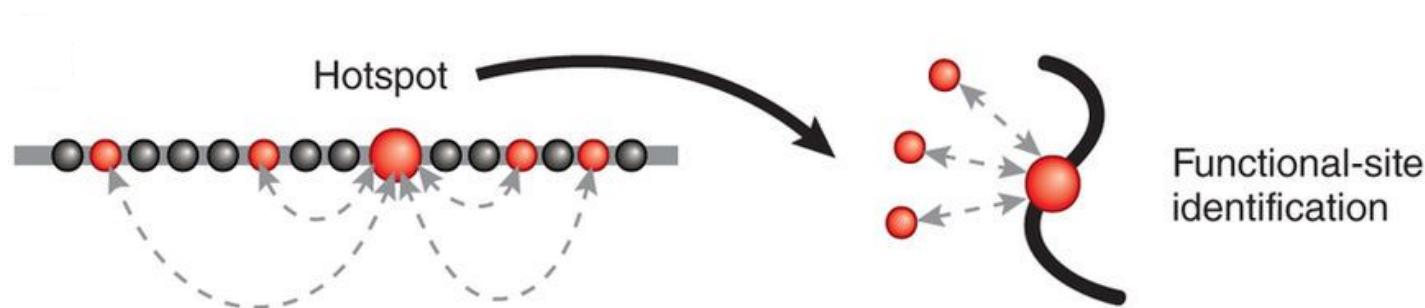
Predictions of 3D structures based on evolutionary coupling



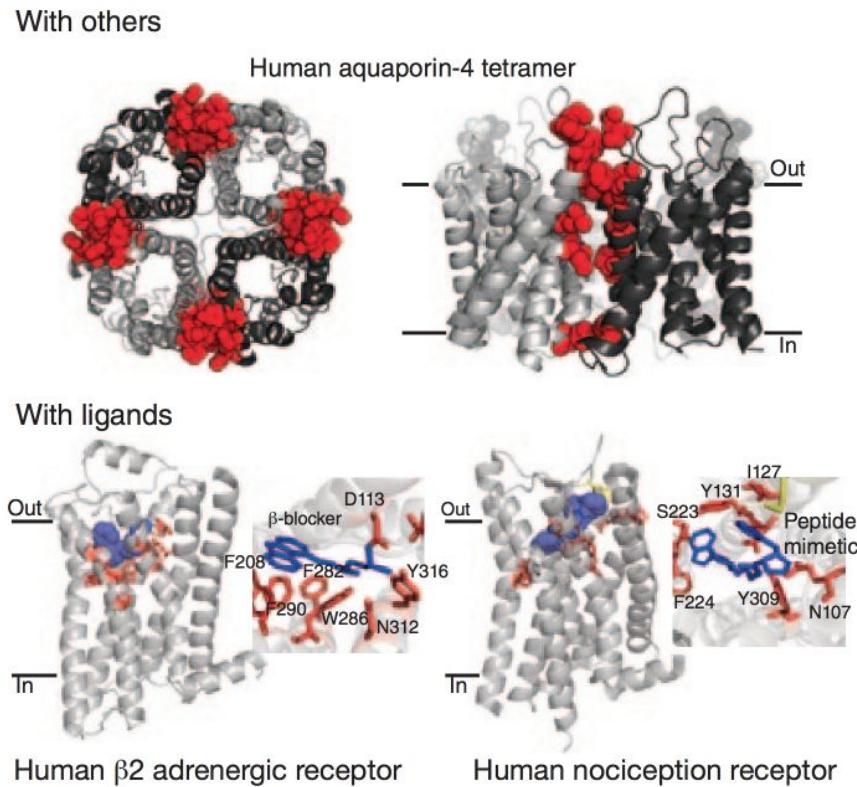
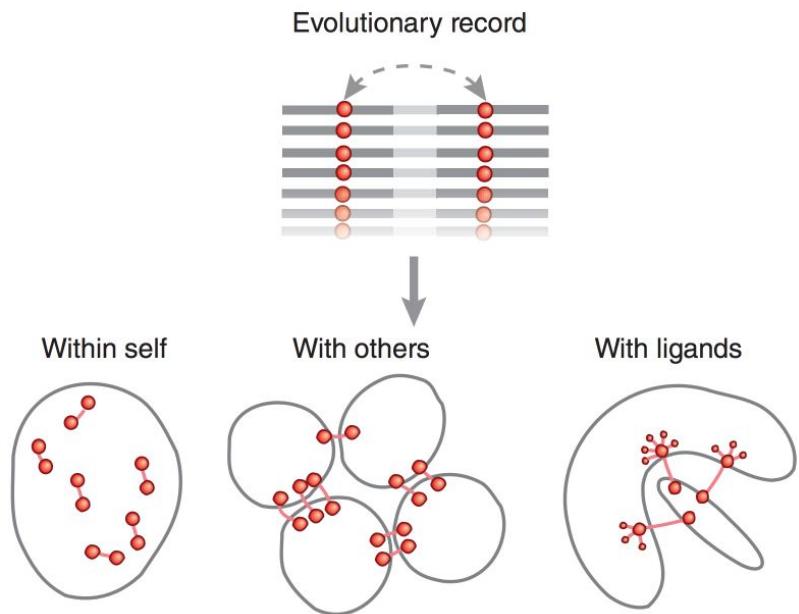
Detecting functional hotspots

Residues subject to a high number of evolutionary pair constraints represent likely functional hotspots.

- Such highly constrained residues include residues in functional sites (for e.g., interaction with external ligands).
- Not detectable by analysis of single-residue conservation.

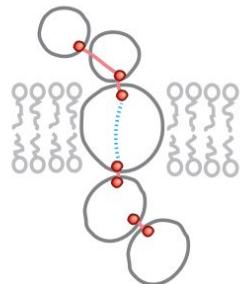


Predicting protein-protein & protein-ligand interactions

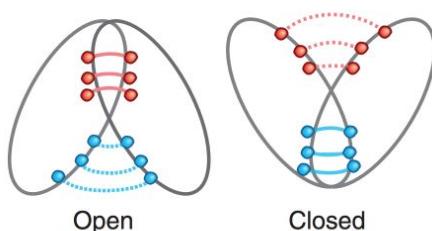


Predicting conformational changes

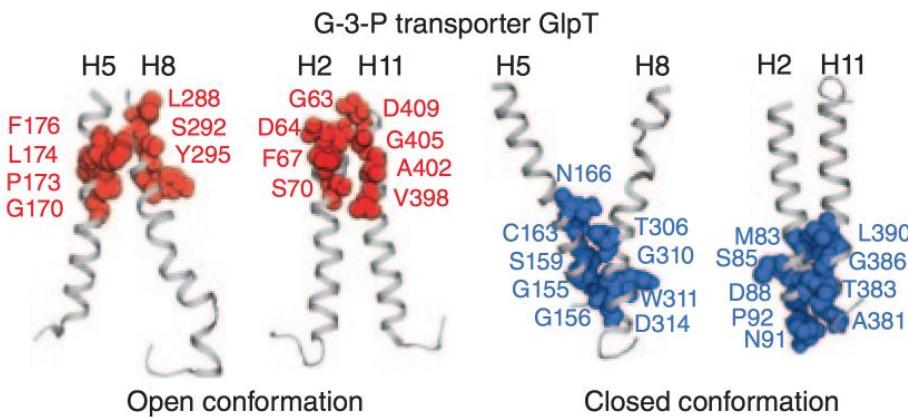
Information transmission



Conformational plasticity



Conformational plasticity



Hybrid approaches for determining protein 3D structure

