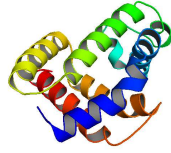


Protein Structure Determination II



Scott McAllister
Princeton University
CS597A, Fall 2005



Outline

- Notable Methods
- ASTRO-FOLD Framework
 - α -helix Prediction
 - β -sheet Prediction
 - Interhelical contact prediction
 - Derivation of Restraints
 - Generation of additional distance bounds
 - Tertiary Structure Prediction
- Results
- Discussion



Protein Structure Prediction

Amino acid sequence [PDB: 1q4sA]

MHRTSNGSHATGGNLPDVASHYPVAYEQQLDGTGTFVIDEMTPERATASVEVDTLRQRWGLVHGGAYCALAEMLA
TEATVAVVHEKGMMMAVGGSNHSTFRPVKEGHVRAEAVRIHAGSTTWFWDVSLRDDAGRLCAVSSMSIAVRPRRD

Helical structure

MHRTSNGSHATGGNLPDVASHYPVAYEQQLDGTGTFVIDEMTPERATASVEVDTLRQRWGLVHGGAYCALAEMLA
TEATVAVVHEKGMMMAVGGSNHSTFRPVKEGHVRAEAVRIHAGSTTWFWDVSLRDDAGRLCAVSSMSIAVRPRRD

Beta strand and sheet structure

MHRTSNGSHATGGNLPDVASHYPVAYEQQLDGTGTFVIDEMTPERATASVEVDTLRQRWGLVHGGAYCALAEMLA
TEATVAVVHEKGMMMAVGGSNHSTFRPVKEGHVRAEAVRIHAGSTTWFWDVSLRDDAGRLCAVSSMSIAVRPRRD



3D Protein Structure



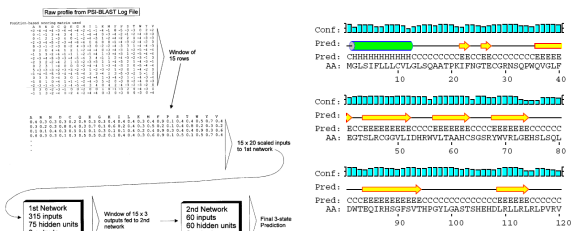
Notable Methods – PSI-PRED

- Prediction of Secondary Structure
- 3 stages:
 - Generation of a sequence profile
 - PSI-BLAST on non-redundant database
 - Creates PSSM for later input
 - Prediction of initial secondary structure
 - Neural network (Feed forward, back propagation)
 - Filtering of the predicted structure
- Web server:
<http://bioinf.cs.ucl.ac.uk/psipred>

Jones, DT. *J. Mol. Biol.* (1999)



Notable Methods – PSI-PRED



Jones, DT. *J. Mol. Biol.* (1999)



Notable Methods - RAPTOR

- Fold recognition method as optimization problem
- Scoring function that accounts for
 - mutation
 - environment
 - gaps
 - secondary structure
 - pairwise interactions

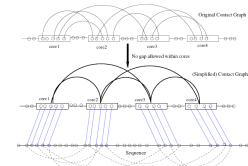


Fig. 1. Template contact graph and an example of alignment between a template and a sequence. A specifically representative structure. A solid line in the original contact graph represents an interaction between two non-sequential residues. A dashed line shows that if two query sequence residues are aligned to two interacting template residues, then the interaction score of those two query sequence residues must be included in the energy function. The interaction score between two residues of the query sequence is the sum of the interaction scores of two query sequence residues which are aligned by two interacting template residues.

$$\min W_m E_m + W_s E_s + W_p E_p + W_g E_g + W_{ss} E_{ss}$$

Xu, J and M. Li. *J. Bioinf Comput Biol.* (2003)



Notable Methods – Robetta

- Combines template-based and ab initio approach
 - Uses template if available, otherwise starts with ab initio
- Also includes methods to detect domains and model loops

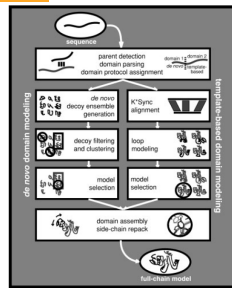


Fig. 1. Robetta process overview. The query sequence is scanned for homology with experimentally determined structures, domain boundaries are determined, and each domain is modeled separately using either the de novo or template-based protocols, assembled into a full-chain model, and side-chains are repacked to produce a full-chain all-atom complete model.

D. Chivian et al. *Proteins*. (2005)

Notable Methods – Skolnick

- Combine aspects of comparative modeling/fold recognition/ab initio
- Apply clustering algorithm to select among conformers

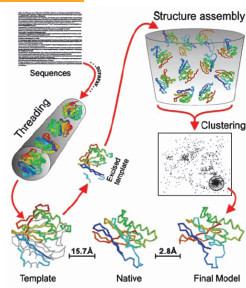
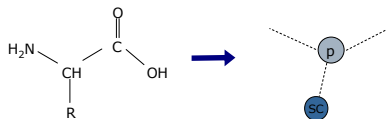


Fig. 1. Overview of the *in silico* structure prediction methodology that consists of template identification by the recursive threading algorithm (R3), CAS fragment assembly, and fold selection by noise clustering (NS). The entire process for 1ay0c is shown as an example.

Zhang, Y. and J. Skolnick. *PNAS*. (2004)

Notable Methods - UNRES

- United Residue approach
- Represent each protein as unified peptide group and unified side chain



- Minimize with this coarse force field, then refine a detailed atomistic force field

Liwo, et al. *J. Comput Chem*. (1997a,b)

Notable Methods – PREDICT

- Method for the structure prediction of membrane proteins, especially GPCRs
 - Does NOT rely on homology to rhodopsin
- Approach
 - Generation of a large number of protein "decoys"
 - Simultaneous optimization and scoring of decoys
 - Optimization includes helix orientations, helix vertical alignments, helix positions, π -stacking of aromatic residues, helical tilts

Shacham, S. *Proteins*. (2004)

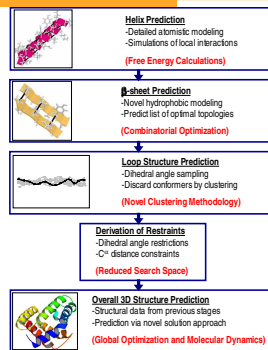
Notable Methods – Available Servers

- 3D Jury (consensus)
 - <http://bioinfo.pl>
- 3D-pssm/pyre
 - <http://www.sbg.bio.ic.ac.uk/~3dpssm/>
- 123D+
 - <http://123d.ncfcf.gov/run123D+.html>
- FFAS03
 - <http://ffas.ljcrf.edu/ffas-cgi/cgi/ffas.pl>
- FORESST
 - <http://abs.cit.nih.gov/foresst2/>

Notable Methods – Available Servers

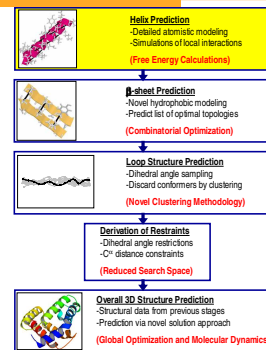
- 3GenesilicoD Jury (consensus)
 - <http://genesilico.pl/meta/>
- PredictProtein
 - <http://cubic.bioc.columbia.edu/predictprotein>
- mGenTHREADER/PSIPRED
 - <http://bioinf.cs.ucl.ac.uk/psipred/psiform.html>
- Robetta
 - <http://rosetta.bakerlab.org>

ASTRO-FOLD



Klepeis, JL and Floudas, CA. *Biophys J.* (2003)

ASTRO-FOLD



Klepeis, JL and Floudas, CA. *Biophys J.* (2003)

Helix Formation

- o Physical characteristics
 - n Well-defined backbone and hydrogen bonding patterns
- o Physical understanding
 - n Local forces: Hierarchical folding
 - n Non-local forces: Hydrophobic collapse
- o Experimental Evidence
 - n Helix formation occurs rapidly
 - n Sequence is sufficient to identify initiation/termination



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2002)

Helix Prediction – Key Ideas

Overlapping oligopeptides

Decompose polypeptide to identify local sites of helix formation and termination

Ensemble of low energy states

Calculate properties of proteins using data from many low energy states rather than a single state

Free energy calculations Klepeis & Floudas 2000

Model proteins using detailed energy calculations including entropic and solvation contributions

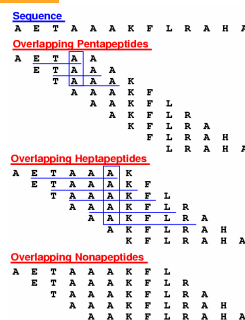
Deterministic global optimization Floudas 2000

Predict low energy states using powerful global optimization approaches such as aBB

Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2002)

Dividing the problem

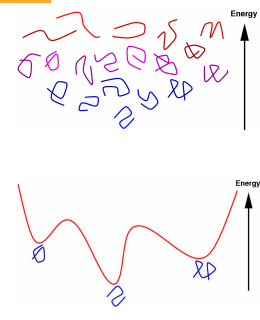
- o Decompose sequence into smaller, overlapping oligopeptides
- o Captures local interactions
- o Free energy calculations on oligopeptides combined to yield overall prediction



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2002)

Generating an ensemble

- o Create low energy states and the global minimum state
- o Formulated as a nonconvex optimization problem
 - n requires global optimization techniques



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2002)

Overall Gibbs Free Energy

- Potential**
 Scheraga & coworkers

$$\sum_{i,j \in \text{DB}} v_{ij} \left[\left(\frac{r_{ij}}{r_{ij}^0} \right)^{12} - \left(\frac{r_{ij}}{r_{ij}^0} \right)^6 \right] + \sum_{i,j \in \text{DB}} v_{ij} \left[\left(\frac{r_{ij}}{r_{ij}^0} \right)^{12} - \left(\frac{r_{ij}}{r_{ij}^0} \right)^6 \right] + \sum_{i,j \in \text{DB}} \frac{332 q_i q_j}{r_{ij}^2} + \sum_{k \in \text{TOB}} \frac{A_k}{r_{ij}^2} \quad (1 \leq k \leq 100; n_k = 0, 1)$$
 $F_{\text{vac}} -$
- Entropic**

$$-\frac{k_B}{2} \ln [\text{Det}(\mathbf{H}_{\text{vac}, \gamma})]$$
 $TS_{\text{vac}} +$
- Cavity**
 Honig & coworkers
 1988, 1993, 1995
 $F_{\text{cavity}} = \gamma(\text{SA}) + b$
 $F_{\text{cavity}} +$
- Polarization**
 Honig & coworkers
 1988, 1993, 1995
 $F_{\text{solvation}} = F_{\text{vac}}(r=80) - F_{\text{polar}}(r=1)$
 $F_{\text{solvation}} +$
- Ionization**
 Honig & coworkers
 1988, 1993, 1995
 $F_{\text{ionization}}(q_i) = kT \ln(Z)$
 $F_{\text{ionization}} +$

Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2002)

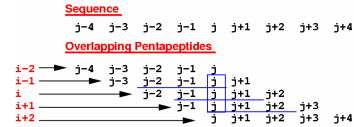
Helix Formation Probability

- Calculate probability of conformer i from free energy

$$p_i = \frac{\exp[-\beta(F_o - F_i)]}{\sum_j \exp[-\beta(F_o - F_j)]}$$

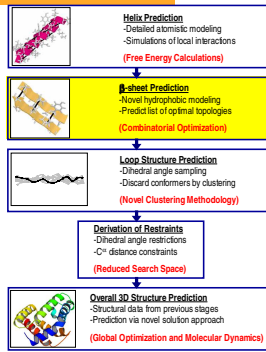
- Use to calculate probability of cluster formation for particular oligopeptide

$$P_{\text{AAA}} = \sum_{i \in \text{AAA}} p_i$$



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2002)

ASTRO-FOLD



Klepeis, JL and Floudas, CA. *Biophys J.* (2003)

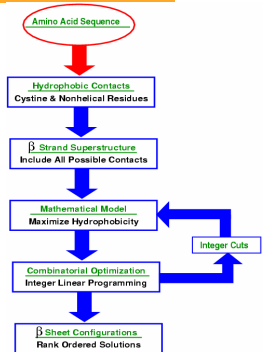
β-sheet Formation

- Major challenge for protein structure prediction
 - β-strand location not accurate
 - Topology prediction not reliable
- Physical understanding
 - Local forces not as important
 - Non-local forces (hydrophobic collapse) dominate
- Experimental evidence
 - Hydrophobic collapse proceeds rapidly



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2003)

β-sheet Prediction Flowchart



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2003)

β-sheet Prediction

- Residue-based aspect
 - Identify set of residues i , and the associated hydrophobicity, H_i
 - Binary variables introduced for residue-residue contact



- Strand-based aspect
 - Identify set of strands s_i , and the associated weight S_{s_i}
 - Binary variables introduced for strand-strand contacts



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2003)

B-sheet Formulation: Key Concepts

Binary variables

0-1 variables are used to characterize residue-to-residue and strand-to-strand contacts

Linear objective function

Objective is to maximize the hydrophobic potential as controlled by the binary variables

Linear constraints

Constraints account for different combinations of residue and strand contacts (e.g., parallel/antiparallel)

Integer cuts

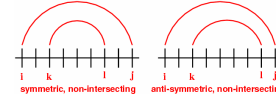
Iterative addition of these constraints allow for the generation of a ranked list of optimal solutions

Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2003)

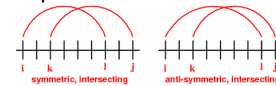


β-sheet Constraints

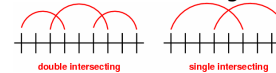
Allowable antiparallel forms



Allowable parallel forms



Disallow double intersecting loops



Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2003)



β-sheet Objective Function

Maximizing hydrophobic potential

$$\max \sum_i \sum_{j, P(i)+2 < P(j)} (H_i + H_j + H_{ij}^{\text{add}}) y_{ij} + \sum_{s_i} \sum_{s_j, Q(s_i) < Q(s_j)} (S_{s_i} + S_{s_j}) w_{s_i, s_j}$$

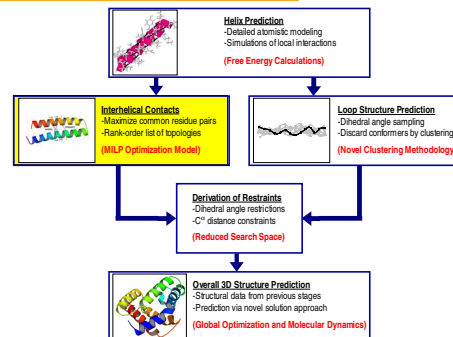
$$y_{ij} = \begin{cases} 1 & \text{if } i, j \text{ form contact} \\ 0 & \text{if } i, j \text{ do not form contact} \end{cases} \quad \forall i < j$$

$$w_{s_i, s_j} = \begin{cases} 1 & \text{if } s_i, s_j \text{ form contact} \\ 0 & \text{if } s_i, s_j \text{ do not form contact} \end{cases} \quad \forall s_i < s_j$$

Klepeis, JL and Floudas, CA. *J. Comput. Chem.* (2003)



ASTRO-FOLD for α-helical Bundles



McAllister, SR and Floudas, CA. *Proceedings, BIOMAT Conference* (2005).



Dataset Selection

- o Protein Sources
 - n 229 PDBSelect25¹ database
 - n 62 CATH² database
 - n 20 Zhang et al.³
 - n 7 Huang et al.⁴
- o Restrictions
 - n No β-sheets, at least 2 α-helices
 - n No highly similar sequences
- o Dataset
 - n 318 proteins in the database set

¹Hobohm, U. and C.Sander. *Prot Sci* 3 (1994) 522

²Orengo, C.A. et al. *Structure* 5 (1997) 1093.

³Zhang, C. et al. *PNAS* 99 (2002) 3581.

⁴Huang, E.S. et al. *J Mol Biol* 290 (1999) 267.

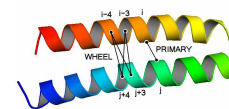
McAllister, SR, et al. (submitted 2005)



Probability Development

Contact Types

- n PRIMARY contact
 - o Minimum distance hydrophobic contact between 4.0 Å and 10.0 Å
- n WHEEL contact
 - o Only WHEEL position hydrophobic contacts between 4.0 Å and 12.0 Å
- o Classified as parallel or antiparallel contacts



McAllister, SR, et al. (submitted 2005)



Model Overview

- o Formulation: **Maximize interhelical residue-residue contact probabilities**
 - n Binary variable $y_{m,n}^a$ indicates antiparallel helical contact
 - n Binary variable $w_{i,j}^{m,n}$ indicates residue contact
- o Goal: Produce a **rank-ordered list** of the most likely helical contacts
 - n Contacts used to **restrict conformational space** explored during protein tertiary structure prediction

McAllister, SR, et al. (submitted 2005)



Pairwise Model Objective

- o Level 1 Objective
 - n Maximize probability of pairwise residue-residue contacts

$$\begin{aligned} \max \quad & \sum_m \sum_n y_{mn}^a \cdot \sum_i \sum_j w_{ij}^{mn} \cdot p_{ij;mn}^a \\ & + \sum_m \sum_n y_{mn}^p \cdot \sum_i \sum_j w_{ij}^{mn} \cdot p_{ij;mn}^p \\ & y_{mn}^a, y_{mn}^p, w_{ij}^{mn} = \{0, 1\} \end{aligned}$$

McAllister, SR, et al. (submitted 2005)



Pairwise Model Constraints

- o Level 1 Constraints
 - n At most one contact per position

$$\sum_{j>i} w_{ij} + \sum_{j<i} w_{ij} \leq 1$$

- n Helix-helix interaction direction

$$y_{mn}^a + y_{mn}^p \leq 1 \quad \forall(m, n)$$

- n Linking interaction variables

$$\begin{aligned} w_{ij}^{mn} &\leq y_{mn}^a + y_{mn}^p \\ y_{mn}^a + y_{mn}^p - \sum_i \sum_j w_{ij}^{mn} &\leq 0 \end{aligned}$$

McAllister, SR, et al. (submitted 2005)



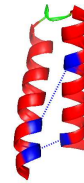
Pairwise Model Constraints

- o Level 1 Constraints
 - n Limit helical kinks

$$w_{ij}^{mn} + w_{j'j'}^{mn} \leq 1$$

$$\forall(i, i', j, j') : |\text{diff}(i, i')| - |\text{diff}(j, j')| \leq 2$$

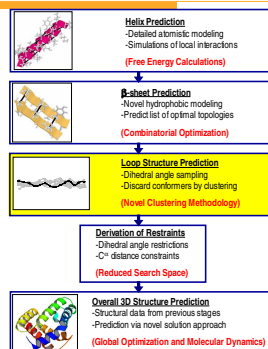
$$\text{or either } |\text{diff}(i, i')| \leq 5 \text{ or } |\text{diff}(j, j')| \leq 5$$



McAllister, SR, et al. (submitted 2005)



ASTRO-FOLD



Klepeis, JL and Floudas, CA. *Biophys J.* (2003)

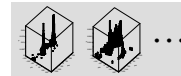


Loop Prediction - Methodology

Create ensemble by **dihedral angle sampling**

- n extracted $p(\phi, \psi)$ from ~2500 loops
- n sampled $p(\phi, \psi)$ at $5^\circ \times 5^\circ$ resolution
- n created ensembles of 2000 conformers for each loop

ARG PRO ...



Structure optimization with first principles force field

- n Dunbrack rotamer library
- n ECEPP/3 force field for structure optimization

Clustering to identify conformers that are close to native

Mönnigmann, M. and Floudas, CA. *Proteins.* (2005)



Loop Prediction – New Use of Clustering

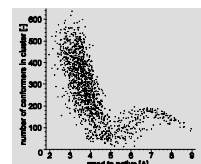
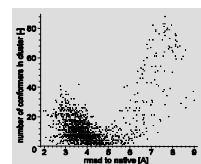
- o Clustering has been used **before** to
 - n **Group** conformers
 - n **Select** conformers that represent groups
- o **New use** of clustering
 - n **Discard** conformers that are far from native
- o First steps of approach
 - n Choose RMSD threshold t
 - n Calculate pairwise RMSD values for the ensemble
 - n For each conformer, record number of conformers with $\text{RMSD} \leq t$

Mönnigmann, M. and Floudas, CA. *Proteins*. (2005)



Loop Prediction – Clustering Example

- threshold $\approx 3.0\text{\AA}$
 - large clusters for small RMSDs unfortunately also for large RMSDs
 - not always advisable to consider centroid of largest cluster
-
- threshold $\approx 3.5\text{\AA}$
 - increasing threshold shows that clusters with large RMSDs are small basins only
 - large clusters with small RMSDs survive

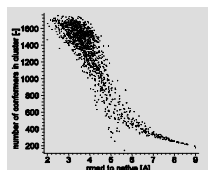
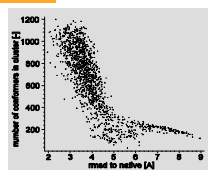


Mönnigmann, M. and Floudas, CA. *Proteins*. (2005)



Loop Prediction – Clustering Example

- threshold $\approx 4.0\text{\AA}$
 - for sufficiently large threshold distribution is monotonous
 - tail with large RMSDs becomes apparent
-
- threshold $\approx 4.5\text{\AA}$
 - distribution more conservative the larger threshold
 - for sufficiently large threshold clusters of conformers with large RMSDs can be discarded



Mönnigmann, M. and Floudas, CA. *Proteins*. (2005)



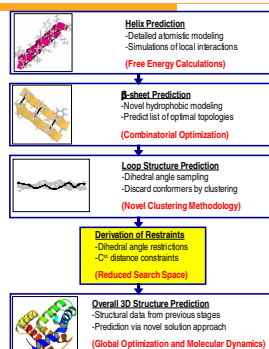
Loop Prediction – Clustering Algorithm

1. Choose threshold t , choose critical cluster size N_{crit}
2. Calculate cluster sizes N_i for all conformers in ensemble
3. If $N_i > N_{crit}$ for all i , stop
4. Discard conformers that generate clusters of size $N_i < N_{crit}$
5. Go back to step 2

Mönnigmann, M. and Floudas, CA. *Proteins*. (2005)



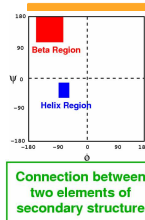
ASTRO-FOLD



Klepeis, JL and Floudas, CA. *Biophys J*. (2003)



Derivation of Restraints



Dihedral angle restraints

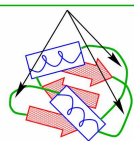
- Backbone **dihedral angles** restrained according to classification of residue as either helix or strand

Distance restraints

- $C^\alpha-C^\alpha$ distance restraints for **hydrogen bond network of helix** (residues i and $i+4$)
- $C^\alpha-C^\alpha$ distance restraints for **predicted interhelical contacts**

Bounds on loop residues

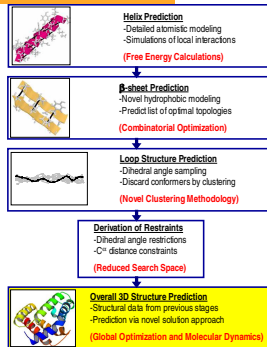
- Based on **dihedral angle deviation** of best identified conformer from loop clustering analysis



Klepeis, JL and Floudas, CA. *Biophys J*. (2003)



ASTRO-FOLD



Klepeis, JL and Floudas, CA. *Biophys J.* (2003)



Tertiary Structure Prediction: Key Ideas

Utilization of helix predictions

Enforce prediction of secondary structure and interhelical distances through rigorous constraint modeling

Mathematical formulation

Formulate tertiary structure prediction problem as a constrained global optimization problem

Energy modeling

Model proteins using detailed atomistic level force field with physically based terms

Global optimization approach

Predict overall tertiary structure using combination of global optimization and torsion angle dynamics

Klepeis, JL and Floudas, CA. *Biophys J.* (2003)



Tertiary Structure Prediction Formulation

$$\min_{\theta} E(\theta)$$

s.t. $E_{l,distance}(\theta) \leq E_{l,ref} \quad l = 1, \dots, N_{con}$
 $\theta_l^L \leq \theta_l \leq \theta_l^U \quad i = 1, \dots, N_{\theta}$

Objective: Nonconvex atomistic level forcefield

$$E = \sum_{i,j \in NB} \epsilon_{ij} \left[\left(\frac{r_{ij}^0}{r_{ij}} \right)^{12} - \left(\frac{r_{ij}^0}{r_{ij}} \right)^6 \right] + \sum_{i,j \in HB} \epsilon_{ij} \left[\left(\frac{r_{ij}^0}{r_{ij}} \right)^{12} - \left(\frac{r_{ij}^0}{r_{ij}} \right)^{10} \right]$$

$$+ \sum_{i,j \in BS} \frac{332 \, a_i a_j}{D r_{ij}} + \sum_{k \in TOH} \frac{A_k}{2} (1 \pm \cos \eta_k \phi_k)$$

Constraints

- Enforce bounds on backbone variables
- Enforce upper / lower distances through square well constraints

$$E_{distance} = \sum_{j \in upper} \begin{cases} A_j (d_j - d_j^U)^2 & \text{if } d_j > d_j^U \\ 0 & \text{otherwise} \end{cases}$$

$$+ \sum_{j \in lower} \begin{cases} A_j (d_j - d_j^L)^2 & \text{if } d_j < d_j^L \\ 0 & \text{otherwise} \end{cases}$$

Klepeis, JL and Floudas, CA. *Biophys J.* (2003)



Torsion Angle Dynamics

Initialization

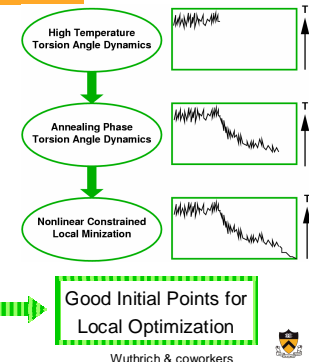
- Difficult to identify low energy feasible structures

Torsion Angle Dynamics

- Identify feasible low energy structures (satisfy constraints)
- Fast evaluation of simplified force field (steric based)
- Unconstrained formulation using penalty functions

Implementation

- Solve equations of motion as preprocessing for each constrained minimization



Klepeis, JL and Floudas, CA. *Biophys J.* (2003)

Wuthrich & coworkers



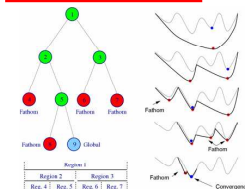
α BB Framework

$$\min_{\mathbf{x}} f(\mathbf{x})$$

s.t. $\mathbf{h}(\mathbf{x}) = \mathbf{0}$
 $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$
 $\mathbf{x} \in \mathbf{X} \subset \mathcal{R}^n$

$f, \mathbf{h}, \mathbf{g}$ twice continuously differentiable

- Based on a branch-and-bound framework
- Upper bound on the global solution is obtained by solving the full nonconvex problem to local optimality
- Lower bound is determined by solving a valid convex underestimation of the original problem
- Convergence is obtained by successive subdivision of the region at each level in the branch & bound tree
- Guaranteed ϵ -convergence for C^2 NLPs



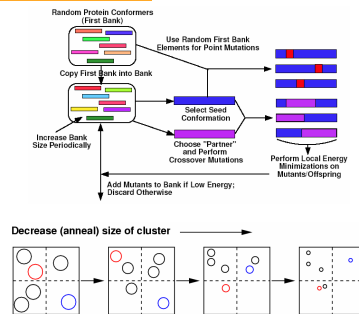
Klepeis, JL and Floudas, CA. *Biophys J.* (2003)



Conformational Space Annealing

Variation types

- Single random
- Backbone random
- Group crossover
- Connected crossover
- Annealing
 - Gradual reduction in space
 - Cluster members in low energy regions



Klepeis, JL, et al. *Biophys J.* (2003)



Hybrid Algorithm Motivation

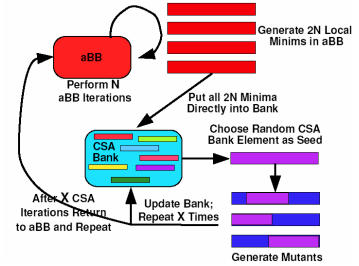
- α BB Features
 - Global minimum guarantee
 - Rigorous upper and lower bounds
 - Rigorous termination
 - Slow performance
- CSA Features
 - Stochastic global minimum search
 - Provides upper bound on solution
 - Heuristic termination
 - Faster performance

Klepeis, J.L. et al. *Biophys J.* (2003)



Alternating Hybrids

- Approach
 - Use output of a series of α BB runs to fill the CSA Bank
 - A number of CSA iterations are performed before refilling the bank with additional α BB solutions

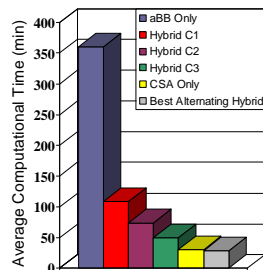


Klepeis, J.L. et al. *Biophys J.* (2003)



Hybrid Performance

- Runs on met-enkephalin
 - Hybrids identified global minimum structure (rigorous)
 - Best alternating hybrid faster than CSA only!

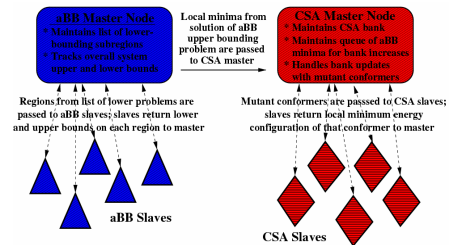


Klepeis, J.L. et al. *Biophys J.* (2003)



Parallelization

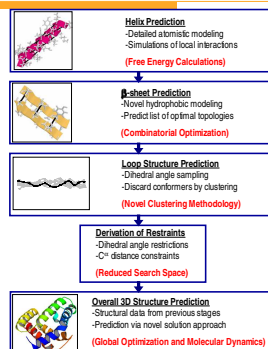
- Protein runs scale as $(N_{res})^{2-4}$
- Implement as distributed method
 - Two "master" nodes: α BB and CSA



Klepeis, J.L. et al. *Biophys J.* (2003)



ASTRO-FOLD



Klepeis, J.L. and Floudas, CA. *Biophys J.* (2003)



Secondary Structure Prediction Results

- Applied to a number of CASP5 targets

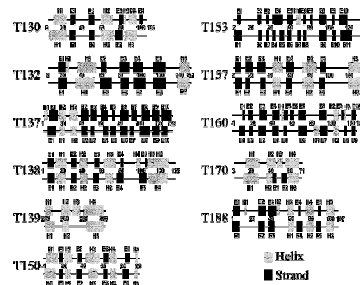


FIGURE 11 Comparison of predictions for helix and β -strand locations with respect to the experimental observations. For each target the top line represents the secondary structure content of the experimentally determined structure, whereas the second line identifies the subsequent prediction results.

Klepeis, J.L. and Floudas, CA. *Biophys J.* (2003)



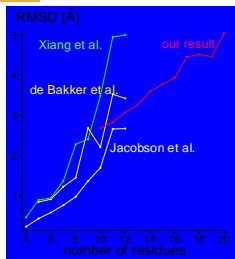
Loop Prediction – Results Comparison

Comparison difficult

- flexible stem residues
- fixed stems in all previous results
- we solve harder problem
- number of residues includes 3+3 stem residues in our case
- stem residues have tighter probability distributions

Results of comparison

- Jacobson et al. result with fixed stems better
- use information on stem geometry, if available
- new method results in very favorable slope
- new method is better than or only slightly worse than methods for fixed stems



Mönnigmann, M. and Floudas, CA. *Proteins*. (2005)

Tertiary Structure Prediction Results

- Successful on a number of small protein systems
- Able to address difficult structures from CASP5 targets

Protein	# of AA	RMSD
1gb1	56	4.2
bp1t	58	4.1
3ci2	63	5.4
r69	68	6.2
t59	75	5.4
t114	87	4.5
t105	95	5.8
t52	101	6.9



FIGURE 6. Comparison of predicted lowest energy tertiary structure (left) of 1gb1 and experimentally determined structure (right). All images generated with the RASPRO molecular visualization package (Sgall and Wilson-White, 1993).

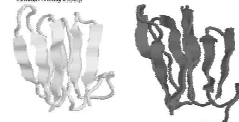
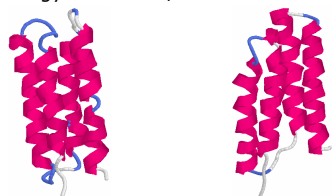


FIGURE 9. Comparison of predicted lowest energy tertiary structure (left) of t114 and experimentally determined structure (right). All images generated with the RASPRO molecular visualization package (Sgall and Wilson-White, 1993).

Klepeis, JL and Floudas, CA. *Biophys J*. (2003)

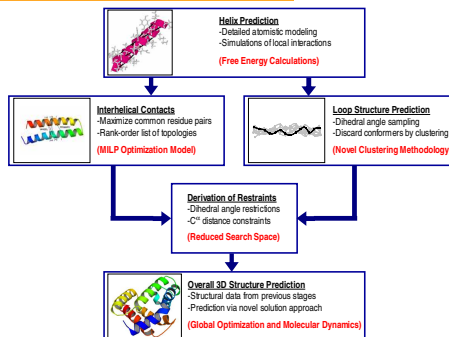
Results – Blind Structure Prediction

- PDB:1p68 (Prof. M. Hecht, Princeton Univ.)
- No information about secondary/tertiary structure
- α -helix: 5-21, 30-49, 56-75, 80-100
- Distance restraints: 63 intrahelical
- Best energy: -846 kcal/mol RMSD: 5.1Å



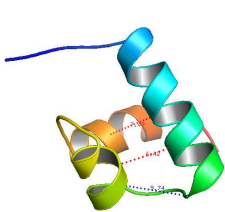
Klepeis, JL, et al. *Proteins*. (2005)

ASTRO-FOLD for α -helical Bundles

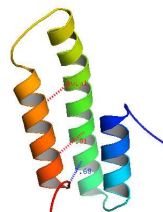


McAllister SR and Floudas, CA. *Proceedings, BIOMAT Conference* (2005).

Results – 2-3 helix bundles



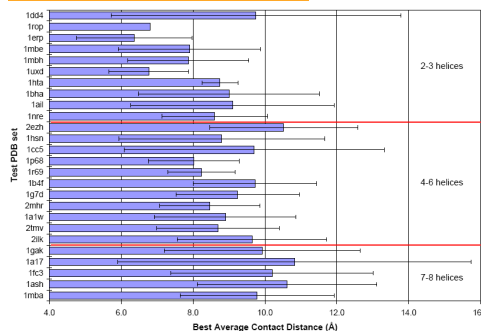
PDB:1mbh in PyMol



PDB:1nre in PyMol

McAllister et al. (submitted 2005)

Results – Contact Prediction Summary

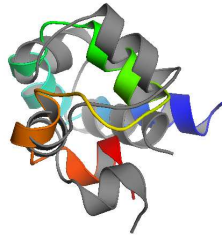


McAllister et al. (submitted 2005)

Results – Tertiary Structure Prediction

- PDB:1r69
 - 63 amino acid protein
 - Top 5 conformers

Conformer	Energy (kcal/mol)	RMSD (Å)
1	-358.77	6.05
2	-351.92	7.95
3	-347.71	6.72
4	-337.80	5.88
5	-337.52	8.04
157	-210.13	4.68



PDB:1r69 in PyMol

- +/- 20° on dihedral angles of loop predictions

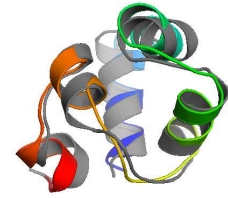
McAllister SR and Floudas, CA. *Proceedings, BIOMAT Conference* (2005).



Results – Tertiary Structure Prediction

- PDB:1r69
 - 63 amino acid protein
 - Top 5 conformers

Conformer	Energy (kcal/mol)	RMSD (Å)
1	-381.54	3.72
2	-376.10	3.40
3	-375.17	4.45
4	-363.77	2.03
5	-356.85	4.91



PDB:1r69 in PyMol

- +/- 10° on dihedral angles of experimental loops

McAllister SR and Floudas, CA. *Proceedings, BIOMAT Conference* (2005).



Discussion

?

