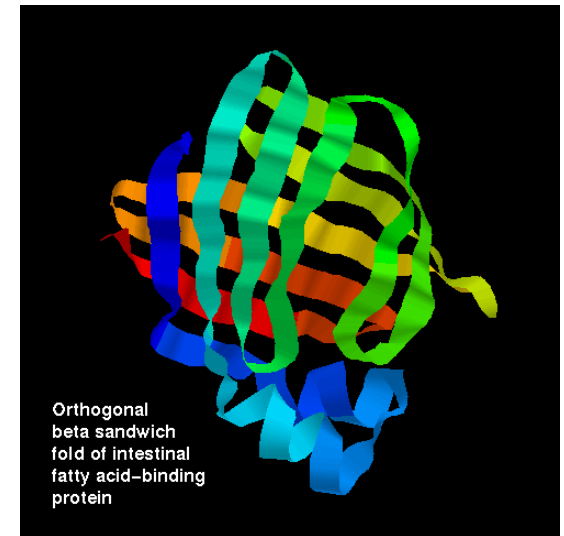
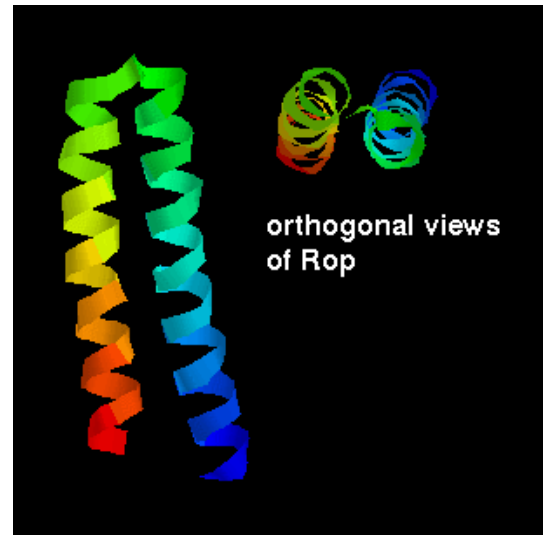


# 3D Structure

## *Prediction & Assessment Pt. 2*



**David Wishart**  
**3-41 Athabasca Hall**  
**david.wishart@ualberta.ca**

# Objectives

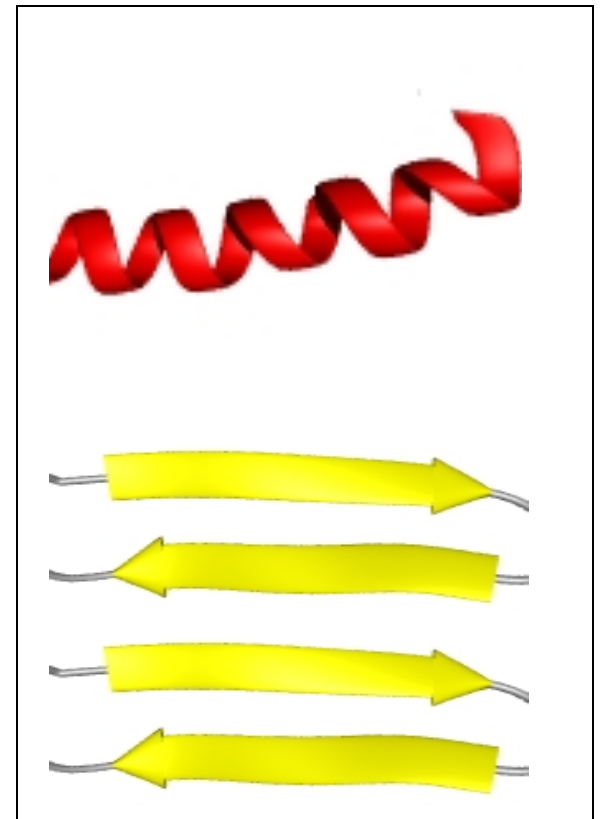
- **Become familiar with methods and algorithms for secondary Structure Prediction**
- **Become familiar with protein Threading (2D and 3D threading)**
- **Become acquainted with Ab initio protein structure prediction**

# 3D Structure Generation\*

- **X-ray Crystallography**
- **NMR Spectroscopy**
- **Homology or Comparative Modelling**
- **Secondary Structure Prediction**
- **Threading (2D and 3D threading)**
- **Ab initio Structure Prediction**

# Secondary (2°) Structure

Phi & Psi angles for Regular Secondary Structure Conformations		
Structure	Phi ( $\Phi$ )	Psi ( $\Psi$ )
Antiparallel $\beta$ -sheet	-139	+135
Parallel $\beta$ -Sheet	-119	+113
Right-handed $\alpha$ -helix	- 64	- 40
$3_{10}$ helix	-49	-26
$\pi$ helix	-57	-70
Polyproline I	-83	+158
Polyproline II	-78	+149
Polyglycine II	-80	+150



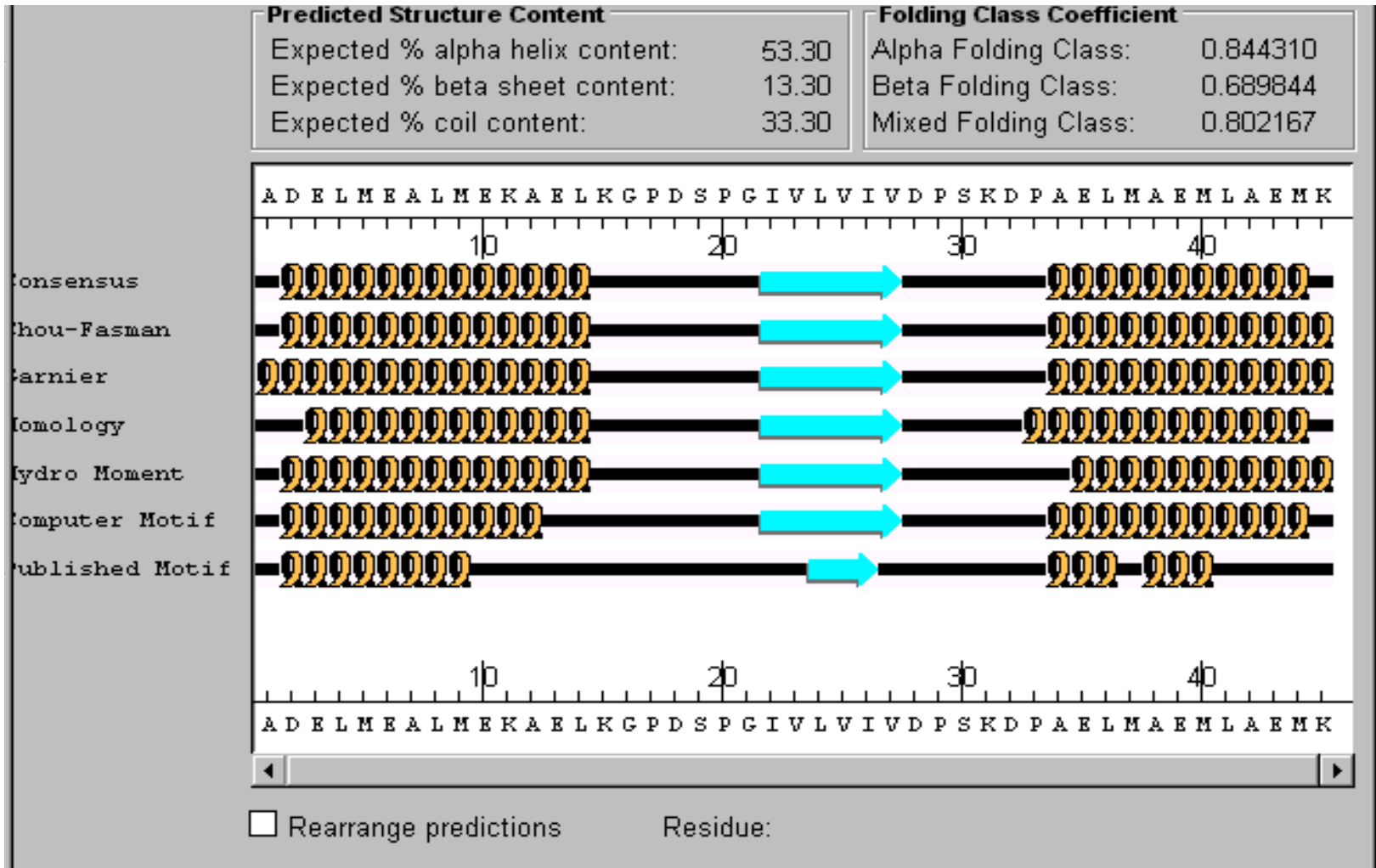
# Secondary Structure Prediction\*

- **One of the first fields to emerge in bioinformatics (~1967)**
- **Grew from a simple observation that certain amino acids or combinations of amino acids seemed to prefer to be in certain secondary structures**
- **Subject of hundreds of papers and dozens of books, many methods...**

# 2° Structure Prediction\*

- **Statistical** (Chou-Fasman, GOR)
- **Homology or Nearest Neighbor** (Levin)
- **Physico-Chemical** (Lim, Eisenberg)
- **Pattern Matching** (Cohen, Rooman)
- **Neural Nets** (Qian & Sejnowski, Karplus)
- **Evolutionary Methods** (Barton, Niemann)
- **Combined Approaches** (Rost, Levin, Argos)

# Secondary Structure Prediction



# Chou-Fasman Statistics\*

**Table 8**

Chou & Fasman Secondary Structure Propensity of the Amino Acids

	$P_{\alpha}$	$P_{\beta}$	$P_c$		$P_{\alpha}$	$P_{\beta}$	$P_c$
A	1.42	0.83	0.75	M	1.45	1.05	0.5
C	0.7	1.19	1.11	N	0.67	0.89	1.44
D	1.01	0.54	1.45	P	0.57	0.55	1.88
E	1.51	0.37	1.12	Q	1.11	1.1	0.79
F	1.13	1.38	0.49	R	0.98	0.93	1.09
G	0.57	0.75	1.68	S	0.77	0.75	1.48
H	1	0.87	1.13	T	0.83	1.19	0.98
I	1.08	1.6	0.32	V	1.06	1.7	0.24
K	1.16	0.74	1.1	W	1.08	1.37	0.45
L	1.21	1.3	0.49	Y	0.69	1.47	0.84

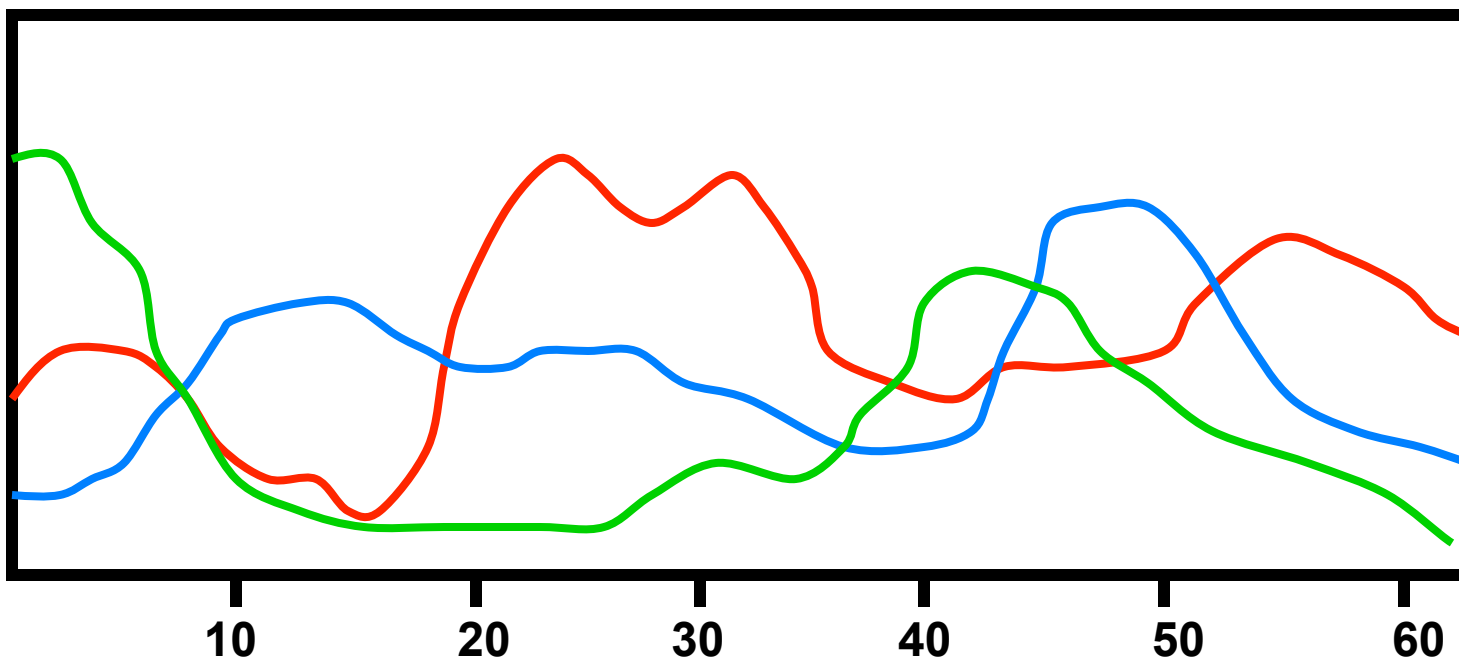


# Simplified C-F Algorithm\*

- **Select a window of 7 residues**
- **Calculate average  $P_{\alpha}$  over this window and assign that value to the central residue**
- **Repeat the calculation for  $P_{\beta}$  and  $P_{\gamma}$**
- **Slide the window down one residue and repeat until sequence is complete**
- **Analyze resulting “plot” and assign secondary structure (H, B, C) for each residue to highest value**

# Simplified C-F Algorithm

helix      beta      coil



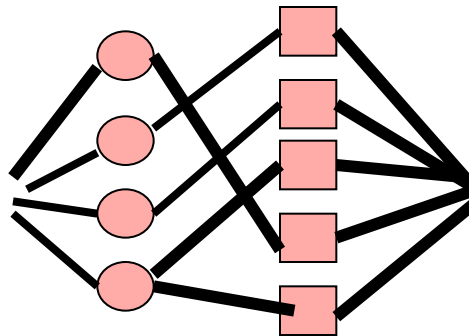
# Limitations of Chou-Fasman

- **Does not take into account long range information (>3 residues away)**
- **Does not take into account sequence content or probable structure class**
- **Assumes simple additive probability (not true in nature)**
- **Does not include related sequences or alignments in prediction process**
- **Only about 55% accurate (on good days)**

# The PhD Approach

Consensus:	CSNLSTCVLGKLSQDLHKLQTFPRT--GAG-P
1: sockeye	CSNLSTCVLGKLSQDLHKLQTFPRTNTGAGVP
2: chum	CSNLSTCVLGKLSQDLHKLQTFPRTNTGAGVP
3: pink	CSNLSTCVLGKLSQDLHKLQTFPRTNTGAGVP
4: coho	CSNLSTCMLGKLSQDLHKLQTFPRTNTGAGVP
5: pig	CSNLSTCVLSAYWRNLNMFHRSFGMGFGPETP
6: bovine	CSNLSTCVLSAYWKDLNMYHRSFGMGFGPETP
7: eel	CSNLSTCVLGKLSQELHKLQTYPRTDVGAGTP

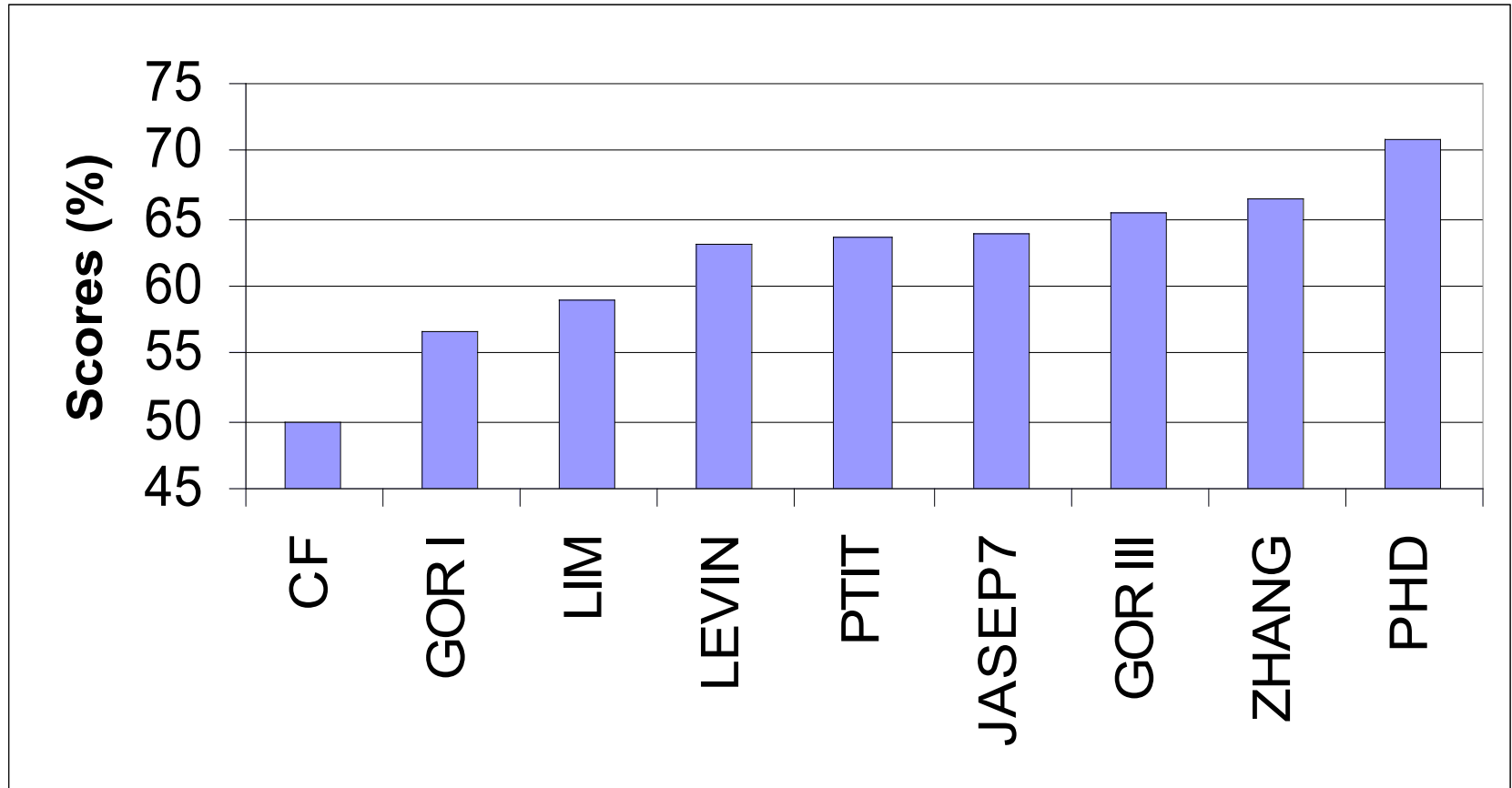
**PROFILE...**



# The PhD Algorithm\*

- *Search the SWISS-PROT database and select high scoring homologues*
- *Create a sequence “profile” from the resulting multiple alignment*
- *Include global sequence info in the profile*
- *Input the profile into a trained two-layer neural network to predict the structure and to “clean-up” the prediction*

# Prediction Performance



# Evaluating Structure 2° Predictions\*

- **Historically problematic due to tester bias (developer trains and tests their own predictions)**
- **Some predictions were up to 10% off**
- **Move to make testing independent and test sets as large as possible**
- **EVA – evaluation of protein secondary structure prediction**

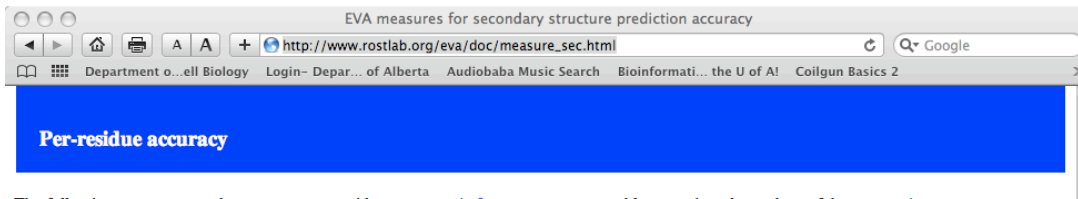
# EVA

Method	Number of proteins	Average ModZscore for all proteins	Comments
<a href="#">APSSP</a>	1312	76.5	
<a href="#">APSSP2</a>	672	82.9	
<a href="#">JNet</a>	?	?	no longer tested
<a href="#">JPred</a>	1218	73.8	
<a href="#">PHD</a>	1599	71.0	
<a href="#">PHDpsi</a>	1598	74.4	
<a href="#">PROF king</a>	1276	74.6	
<a href="#">PROFsec</a>	1554	76.7	
<a href="#">Prospect</a>	103	71.7	
<a href="#">PSIpred</a>	1461	77.9	
<a href="#">PSSP</a>	?	?	no longer tested
<a href="#">SAM-T99sec</a>	543	75.6	
<a href="#">SSpro1</a>	?	?	no longer tested
<a href="#">SSpro2</a>	1348	76.9	

- ~10 different methods evaluated in real time as new structures arrive at PDB
- Results posted on the web and updated weekly
- <http://www.pdg.cnb.uam.es/eva/>



# EVA- <http://www.pdg.cnb.uam.es/eva/>



The following scores are used to measure per-residue accuracy ([reference](#), note: most tables contain only a subset of these scores):

## 1. Prediction accuracy matrix:

$M_{ij}$  = number of residues observed in state  $i$  and predicted in state  $j$ , with  $i$  and  $j \in \{H, E, L\}$   
 note: the total number of residues observed in state  $i$  is:

$$obs_i = \sum_{j=1}^3 M_{ij}, \text{ with } j \in \{H, E, L\}$$

note: the total number of residues predicted in state  $i$  is (helix, strand, other)

$$prd_i = \sum_{j=1}^3 M_{ji}, \text{ with } i, j \in \{H, E, L\}$$

and the total number of residues is simply:

$$N_{res} = \sum_i obs_i = \sum_i prd_i = \sum_{i,j} M_{ij}$$

## 2. Three-state prediction accuracy: $Q_3$

Thus, the three-state per residue accuracy  $Q_3$  becomes:

$$Q_3 = 100 \cdot \frac{1}{N_{res}} \cdot \sum_{i=1}^3 M_{ii}$$

## 3. Per-state percentages:

To define accuracy for a particular state (helix, strand, other), there are two possible vari

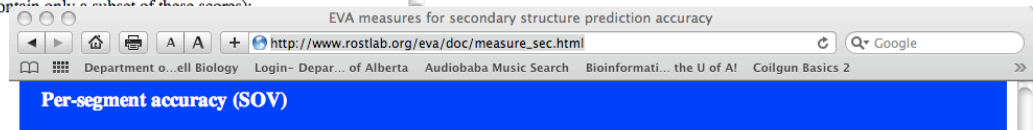
- How many of the observed helix residues (strand/other) were correctly predicted?  
 Given are the correctly predicted residues as percentage of all residues OBSERVE

$$Q_i^{obs} = 100 \cdot \frac{M_{ii}}{obs_i}$$

- How many of the predicted helix (strand/other) residues were correctly predicted?  
 Given are the correctly predicted residues as percentage of all residues PREDICTI

$$Q_i^{prd} = 100 \cdot \frac{M_{ii}}{prd_i}$$

## 4. Information index:



The Segment OVerlap measure for prediction accuracy is defined by (for [more details](#)):

### • Per-stage segment overlap:

$$SOV_i = \frac{1}{N_i} \sum_{S1, S2} \frac{MINOV(S1;S2) + DELTA(S1; S2)}{MAXOV(S1;S2)}$$

with the following definitions:

$S1$  and  $S2$  are the observed and predicted secondary structure segments (in state  $i$ , which can be either H, E or L)  
 $LEN(S1)$  is the number of residues in the segments  $S1$   
 $MINOV(S1;S2)$  is the length of actual overlap of  $S1$  and  $S2$ , i.e. the extent for which both segments have residues  
 $MAXOV(S1;S2)$  is the length of the total extent for which either of the segments  $S1$  or  $S2$  has a residue in state  $i$   
 $DELTA(S1;S2)$  is the integer value defined as being equal to the following

$$DELTA(S1;S2) = \min \left\{ \begin{array}{l} MAXOV(S1;S2) - MINOV(S1;S2) \\ MINOV(S1;S2) \\ INT(0.5 \cdot LEN(S1)) \\ INT(0.5 \cdot LEN(S2)) \end{array} \right\}$$

THE SUM is taken over  $S$ , all the pairs of segments  $\{S1;S2\}$ , where  $S1$  and  $S2$  have at least one residue in state  $i$   
 $N(i)$  is the number of residues in state  $i$  defined as follows:

$$N_i = \sum_{S(i)} LEN(S1) + \sum_{S'(i)} LEN(S1)$$

The two sums are taken over  $S$  and  $S'$ :  
 $S(i)$  is the number of all the pairs of segments  $\{S1;S2\}$ , where  $S1$  and  $S2$  have at least one residue in state  $i$   
 $S'(i)$  is the number of segments  $S1$  that do not produce any segment pair

### • Segment OVerlap quantity measure for all three states:

$$SOV = SOV_3 = \frac{1}{N} \cdot \sum_{S(i)} \frac{MINOV(S1;S2) + DELTA(S1;S2)}{MAXOV(S1;S2)} \cdot LEN(S1)$$

with:

$$N = \sum N_i$$

## **2° Structure Evaluation\***

- **Q3 score – standard method in evaluating performance, 3 states (H,C,B) evaluated like a multiple choice exam with 3 choices. Same as % correct**
- **SOV (segment overlap score) – more useful measure of how segments overlap and how much overlap exists**

# Best of the Best

- **PredictProtein-PHD (74%)**
  - <http://www.predictprotein.org/meta.php>
- **Jpred (73-75%)**
  - <http://www.compbio.dundee.ac.uk/www-jpred/>
- **PSIpred (77%)**
  - <http://bioinf.cs.ucl.ac.uk/psipred/>
- **Proteus and Proteus2 (88%)**
  - <http://wks80920.ccis.ualberta.ca/proteus/>
  - <http://www.proteus2.ca/proteus2/>

# Meta PredictProtein

- About
- Submission
- Help
- Downloads
- Register
- MetaPP

Welcome David [My Queries](#) [Edit Account](#) [Logout](#)

Description of field (click on description for help)

Type the required information into the fields (and select one or more services)

Your email address (watch for typos)

One-line name of protein

Paste, or type your sequence

- amino acids in one-letter code (any number of spaces allowed)
- other possible formats

For retrieving protein sequences from databases we recommend the Sequence Retrieval System [SRS6](#)

Ads by Google

[Protein Substrates](#)

High capacity NC slides Your reliable source for substrates [www.thermo.com/advance](http://www.thermo.com/advance)

[Proteomics data analysis](#)

Compare multiple protein lists ProteinCenter FastTrack publication [www.proxeon.com](http://www.proxeon.com)

[Custom Peptide Antibodies](#)

Design peptides and develop better antibodies using Antigen Profiler [www.OpenBiosystems.com](http://www.OpenBiosystems.com)

[Protein Evolution](#)

Superior to Directed Evolution Next Generation Technologies

## Available Services

Choose (at least one) checkbox(es) to request respective services for your protein

Homology-based prediction of 3D structure (not always possible)

<input type="checkbox"/>	Homology Modelling	<b>3D-JIGSAW</b>	<a href="#">Go There</a>	<a href="#">About</a>
--------------------------	--------------------	------------------	--------------------------	-----------------------

Jpred - A Secondary Structure Prediction Server

http://www.compbio.dundee.ac.uk/www-jpred/

Most Visited ▾ Getting Started Latest Headlines ↗

Jpred - A Secondary Structure Pre... +



# Jpred 3

Incorporating Jnet

## A Secondary Structure Prediction Server

Sequence:

[Help](#)  
[Advanced](#)

Make Prediction

Clear

[The Barton Group - The University of Dundee](#)

Citation: Cole C, Barber JD & Barton GJ. *Nucleic Acids Res.* 2008. [\[Advanced Access\]](#)

[More citations](#)



## Site Navigation

## Server Navigation

- PSIPRED Server
- PSIPRED help
- Server Overview
- Server Citation
- News
- History
- Build A TDB File
- Software Download
- Login

## The PSIPRED Protein Structure Prediction Server

The PSIPRED Protein Structure Prediction Server aggregates several of our structure prediction methods into one location. Users can submit a protein sequence, perform the prediction of their choice and receive the results of the prediction via e-mail. You may select one of three prediction methods to apply to your sequence:

PSIPRED - a highly accurate method for protein secondary structure prediction  
MEMSAT and MEMSAT-SVM - our widely used transmembrane topology prediction method  
and one of GenTHREADER, pGenTHREADER and pDomTHREADER - sequence profile based fold recognition methods. [More...](#)

**For queries regarding PSIPRED:** [psipred@cs.ucl.ac.uk](mailto:psipred@cs.ucl.ac.uk)

### Choose Prediction Method

- Predict Secondary Structure (PSIPRED v3.0)
- Predict Transmembrane Topology (MEMSAT3 & MEMSAT-SVM)
- SVM Prediction of TM Topology and Helix Packing (MEMPACK) - **NEW!**
- Fold Recognition (GenTHREADER - quick)
- Fold Recognition (pGenTHREADER - with profiles and predicted secondary structure)
- Fold Recognition (pDomTHREADER - annotates multiple domain on chains)

[Help...](#)

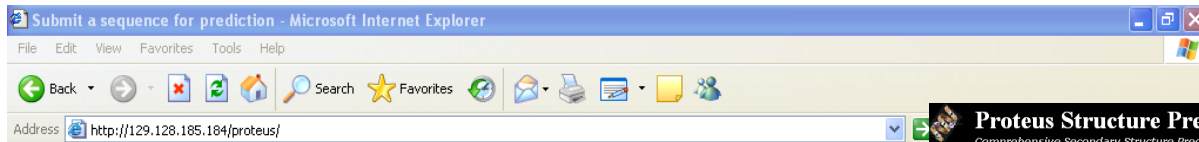
### Input Sequence (single letter amino acid code)

[Help...](#)

If you wish to test these services follow this link to retrieve a [test fasta sequence](#).

### Filtering Options

# Proteus



**Proteus Structure Prediction Server**  
 Comprehensive Secondary Structure Predictions

HOME DOCUMENTATION SAMPLE OUTPUT CONTACT & DOWNLOAD

## Welcome to Proteus

Proteus is a high-performing integrated web server and a prediction methods (PSIPRED, JNET and TRANSSEC) and a robust PDB-based structure alignment process to generate protein. Proteus is able to achieve a very high level of accuracy. Proteus query protein shows no similarity whatsoever to any known. Proteus is not restricted to generating accurate secondary structure predictions for integral membrane proteins (both helix-containing proteins) homologues or a portion of a homologue in the PDB. The 2005 updated PDB secondary structure database available for download.

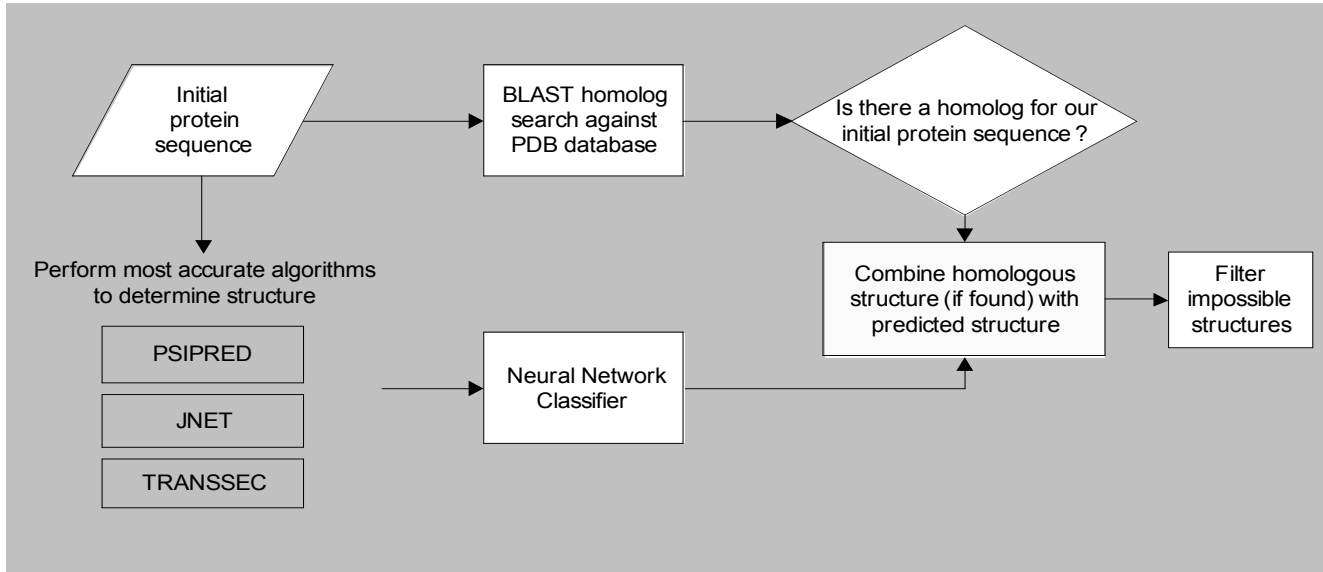
### Submitting a sequence

You can submit your sequence in FASTA Format by pasting it to the server. A prediction will be returned to you.

### Paste Single Sequence (FASTA format)

```
<div data-bbox="555 265 955 315" data-label="Complex-Block" style="background-color: black; color: white; padding: 5px;">
    Proteus Structure Prediction Server
    Comprehensive Secondary Structure Predictions
    HOME DOCUMENTATION SAMPLE OUTPUT CONTACT & DOWNLOAD
  </div>
  <div data-bbox="555 315 955 565" data-label="Complex-Block" style="background-color: #f0f0f0; padding: 5px;">
    Proteus prediction (ID=5689539) complete
    Summary:
    <ul>
      <li>Time of Submission: 21:34:16 Apr 01, 2006</li>
      <li>Sequence Name:</li>
      <li>Number of residues read in: 844</li>
      <li>Number of useable PDB homologs found: 12
          <ul>
            <li>1F3ZA, e-value = 1.0E-105 NMR STRUCTURE OF THE Y174 AUTOINHIBITED DBL HOMOLOGY DOMAIN</li>
            <li>2C9HA, e-value = 5.0E-65 SOLUTION STRUCTURE OF THE SH3 DOMAIN OF HUMAN PROTO...</li>
            <li>1GCPA, e-value = 8.0E-34 CRYSTAL STRUCTURE OF VAV SH3 DOMAIN</li>
            <li>1YXDA, e-value = 6.0E-17 CRYSTAL STRUCTURE OF THE DR/PH DOMAINS OF LEUKEMIA-...</li>
            <li>1XC6A, e-value = 2.0E-13 CRYSTAL STRUCTURE OF HUMAN SHOA IN COMPLEX WITH DR/PH</li>
            <li>1KZ7A, e-value = 7.0E-12 CRYSTAL STRUCTURE OF THE DR/PH FRAGMENT OF MURINE DBS IN</li>
            <li>1FOEA, e-value = 5.0E-11 CRYSTAL STRUCTURE OF RAC1 IN COMPLEX WITH THE GUANINE</li>
            <li>1KI1B, e-value = 1.0E-9 GUANINE NUCLEOTIDE EXCHANGE REGION OF INTERSECTIN IN</li>
            <li>1UJYA, e-value = 7.0E-9 SOLUTION STRUCTURE OF THE SH3 DOMAIN IN RAC/CDC42 GUANINE</li>
            <li>1AIZEA, e-value = 1.0E-8 NMR STRUCTURE OF THE COMPLEX BETWEEN THE C32S-Y7V MUTANT OF</li>
            <li>1UFFA, e-value = 5.0E-8 SOLUTION STRUCTURE OF THE FIRST SH3 DOMAIN OF HUMAN</li>
            <li>100TA, e-value = 6.0E-8 CRYSTAL STRUCTURE OF THE SH3 DOMAIN FROM A S. CEREVISIAE</li>
          </ul>
        </li>
      <li>Number of sequence alignments used for ab-initio predictions: 49</li>
      <li>Overall confidence value: 83.4%</li>
      <li>Predicted % Helix content: 37 % (313 residues)</li>
      <li>Predicted % Beta sheet content: 14 % (115 residues)</li>
      <li>Predicted % Coil content: 49 % (416 residues)</li>
    </ul>
  </div>
  <div data-bbox="375 430 415 445" data-label="Section-Header" style="background-color: #f0f0f0; padding: 5px;">
    Legend:
  </div>
  <div data-bbox="390 445 555 540" data-label="Text" style="background-color: #f0f0f0; padding: 5px;">
    H = Helix
    E = Beta Strand
    C = Coil
    Line 1 = sequence (single letter IUPAC code),
    Line 2 = secondary structure (H, E or C)
    Line 3 = confidence score (0-9, 0 = low, 9 = high)
    A '*' character above the overall prediction residue.
  </div>
  <div data-bbox="375 550 505 565" data-label="Section-Header" style="background-color: #f0f0f0; padding: 5px;">
    Predicted Secondary Structure:
  </div>
  <div data-bbox="395 565 560 885" data-label="Text" style="background-color: #f0f0f0; padding: 5px;">
    1
    MELWRQCTHVLIQKRVLPSPHRVTWDGAQVCELAAQALRE
    CCHHHHHHHHHHHCCCCCCCCCCCCCHHHHHHHHHHCC
    98589999998763567777788777658999999984
    61
    NLRPQNSQFLCLKNIRITFLSTCCCKFLKRSLEFAFDL
    CCCCCCHHHHHHHHHHHHHHHHHHHHCCCCCHHCCCCCH
    788775799999999999999999985776565478757E
    121
    AQNRGINPFPTEESVGDIEDIYSGLSDQIDDTVEEDEDL
    CCCCCCCCCCCCCCCCCHHCCCCCCCCCCCCCCCC
    66777777667777777555578777766667886
    181
    ***** 240
    EPIVSMPPKMTEDYKRCCLREIQQTEERYDTLGSIQQHLKPLQFLKFPDIEIIFINI
    CCCCCCCCCCHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
    9889999999999999999999999999999999999999999999999999999999
    241
    ***** 300
    EDLLRVHTFLKMKKALGTPGAANLYQVFIKYERFLVYGRYCSQVESASKHLDRVAAA
    HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
    9999999999999999999999999999999999999999999999999999999999
    301
    ***** 360
    REDVQMKLECSQRANRGRFTLRDLLMVPQRVLRKYLHLLQELVKHTQAMEKKNLRLL
    HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
    9999999999999999999999999999999999999999999999999999999999
  </div>
  <div data-bbox="575 565 715 580" data-label="Section-Header" style="background-color: #f0f0f0; padding: 5px;">
    Graphical Alignment of PDB Homologs:
  </div>
  <div data-bbox="575 580 900 705" data-label="Figure" style="background-color: #f0f0f0; padding: 5px;">
    QUERY
    1F3ZA
    2C9HA
    1GCPA
    1YXDA
    1KCGA
    1KZ7A
    1FOEA
    1KI1B
    1UJYA
    1AIZEA
    1UFFA
    100TA
    <img alt="Graphical alignment of PDB homologs showing sequence identity bars for various proteins like 1F3ZA, 2C9HA, 1GCPA, etc." data-bbox="575 580 900 705"/>
  </div>
  <div data-bbox="76 810 355 855" data-label="Image" style="background-color: #f0f0f0; padding: 5px;">
    start Snagit Submit a sequen... Local Dis
  </div>
  </div>
```

# Proteus Methods\*



Query sequence:

DLQTTGADHSATVNPDQQLIMTKHSATVTPENKCVFFPNYRGYRYDCTRTRDSFYRWCSLTGTYSGSWKYCAATDYAKC

Predicted structure (consensus method):

DLQTTGADHSATVNPDQQLIMTKHSATVTPENKCVFFPNYRGYRYDCTRTRDSFYRWCSLTGTYSGSWKYCAATDYAKC  
 --BBBBBB--BBB-----BBBB--BB--BBBB-----BBBBBB--BBB-----BBBB-----

Structure of homolog:

KCVFFPNYRGYRYDCTRTRDSFYRWCSLTGTYSGSWKYCAATDYAKCAFPFVYRGQTYD  
 --BBBBBB--BBB-----BBBB--BB--BBBB-----BBBBBB--BBB-----

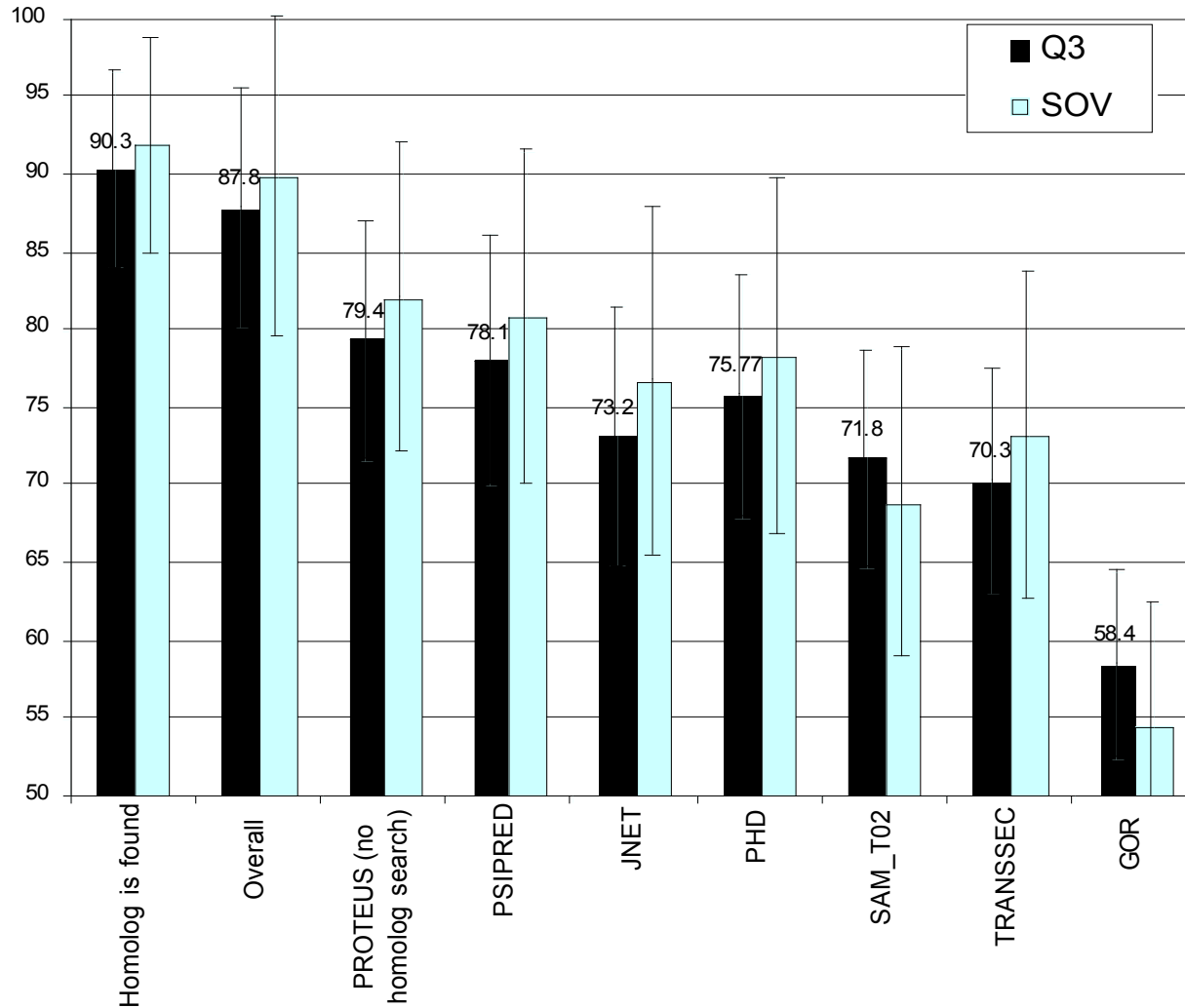
Predicted structure retained

Homologous structure overlaid on predicted structure

DLQTTGADHSATVNPDQQLIMTKHSATVTPENKCVFFPNYRGYRYDCTRTRDSFYRWCSLTGTYSGSWKYCAATDYAKC  
 --BBBBBB--BBB-----BBBB--BB--BBBB-----BBBBBB--BBB-----BBBB-----

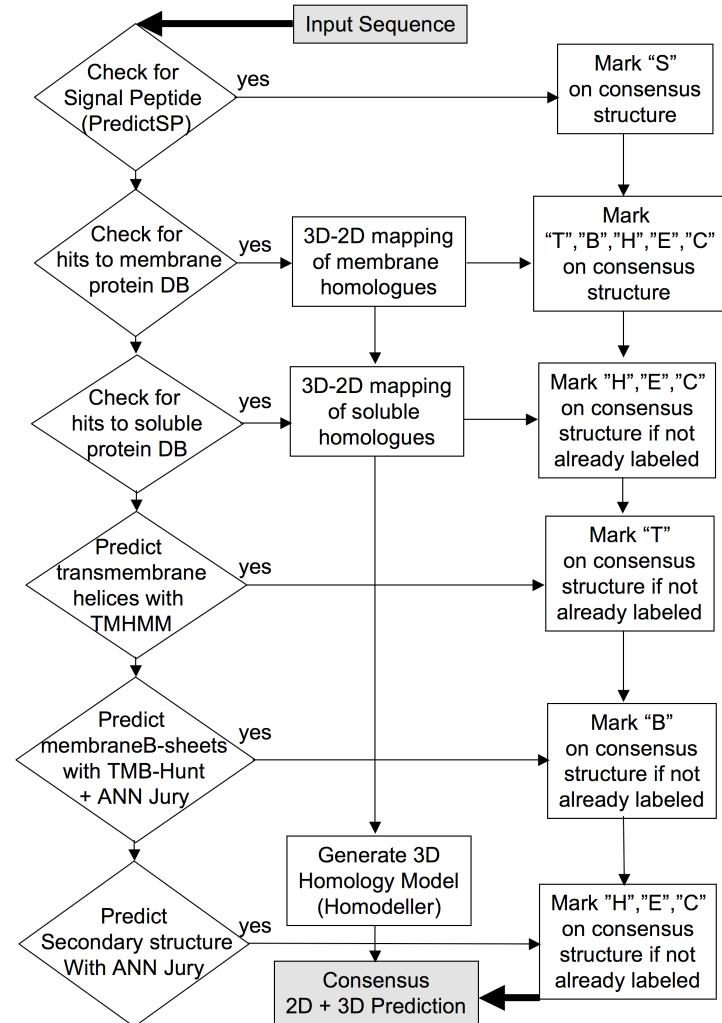


# Performance Comparison



# Proteus2\*

The screenshot displays the Proteus2 web interface, which is used for protein structure prediction. It includes a submission form, a 'Proteus2 prediction (ID=2464967) complete' summary, and a 'Homology Modelling' section. The summary lists various metrics such as sequence name, number of residues, and confidence values. The homology modeling section shows a template ID (2ILPA) and a percent homology of 96%. The interface also features a 3D visualization of the protein structure with various controls for viewing and interacting with the model.



# Proteus2 Performance\*

<i>Transmembrane Helix Prediction Performance (TMH Benchmark test set)</i>		
Program or Server	Q2	# False positives
PROTEUS2	91%	0
TMHMM	80%	1
HMMTOP	80%	6
DAS	72%	16
<i>Transmembrane Helix Prediction Performance (PPT-DB-TMH test set)</i>		
Program or Server	Q2	# False neg. (Missed Prots)
PROTEUS2	87%	0
TMHMM	82%	8
<i>Transmembrane Beta Barrel Detection Performance (PPT-DB "All" protein data set)</i>		
Program or Server	Q2	Accuracy (TMB vs glob)
PROTEUS2	100%	100%
TMB-Hunt	78%	99%
<i>Transmembrane Beta Strand Prediction Performance (PPT-DB -TMB test set)</i>		
Program or Server	Q2	
PROTEUS2	86%	
Pred-TMBB	73%	
<i>Non-membrane Secondary Structure Prediction Performance (EVA Test Set)</i>		
Program or Server	Q3	SOV
PROTEUS2	81	82
Porter	77	76
JNET	72	73
PSIPred	77	78

# Definition\*

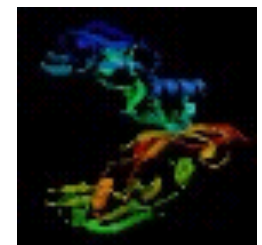
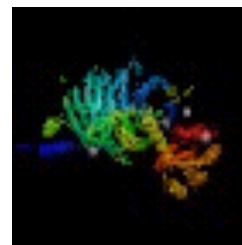
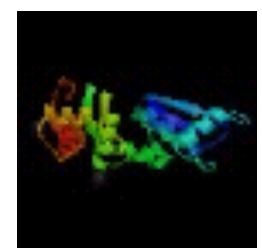
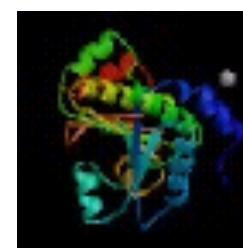
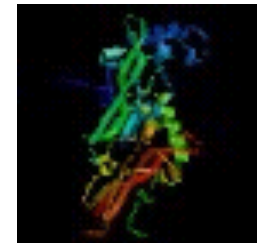
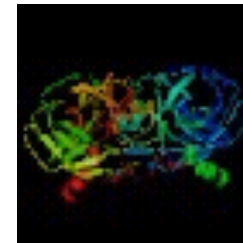
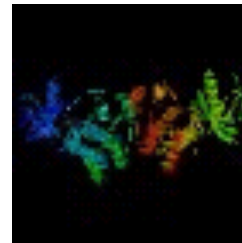
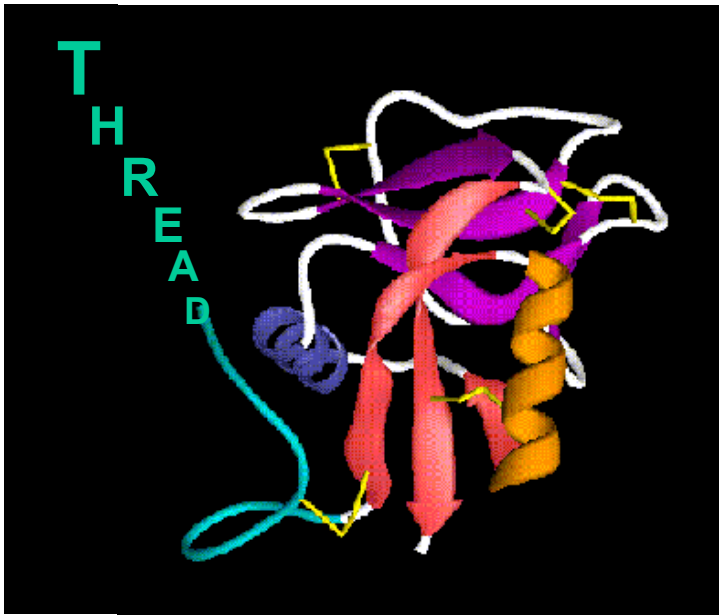
- **Threading** - A protein fold recognition technique that involves incrementally replacing the sequence of a known protein structure with a query sequence of unknown structure. The new “model” structure is evaluated using a simple heuristic measure of protein fold quality. The process is repeated against all known 3D structures until an optimal fit is found.

# Why Threading?\*

- **Secondary structure is more conserved than primary structure**
- **Tertiary structure is more conserved than secondary structure**
- **Therefore very remote relationships can be better detected through 2° or 3° structural homology instead of sequence homology**

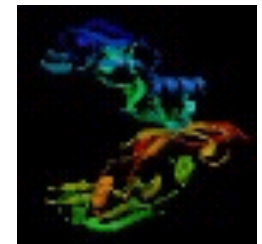
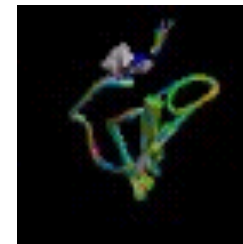
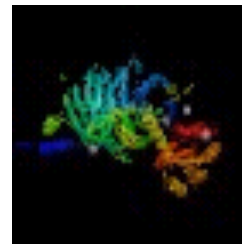
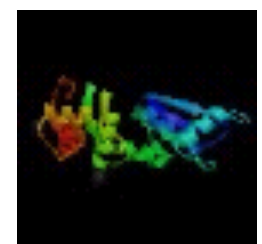
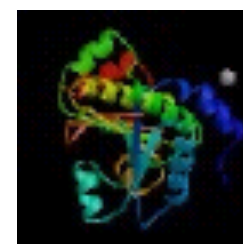
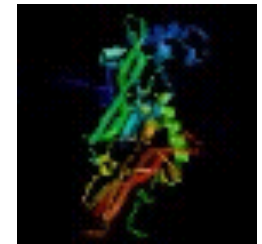
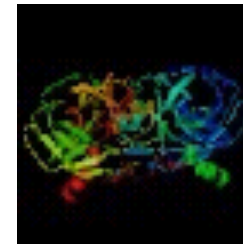
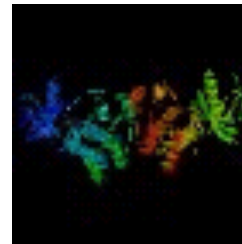
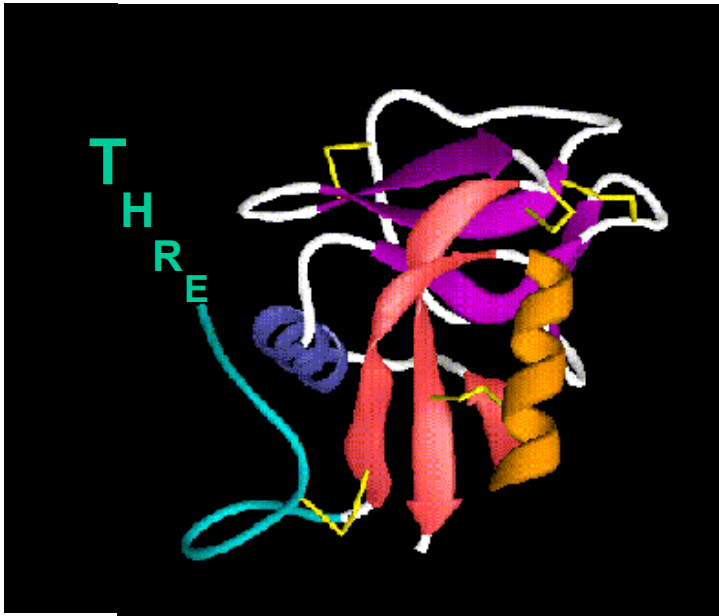
# Visualizing Threading

THREADINGSEQNCEECNQESGNI  
ERHTHREADINGSEQNCETHREAD  
GSEQNCEQCQESGIDAERTHR...



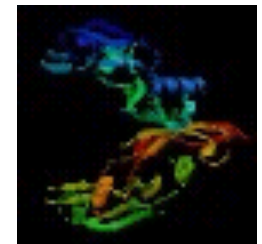
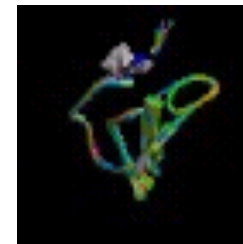
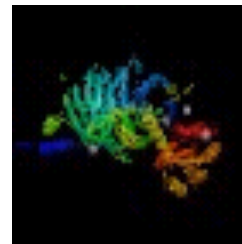
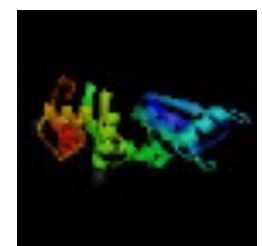
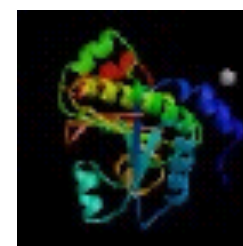
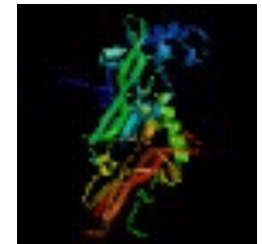
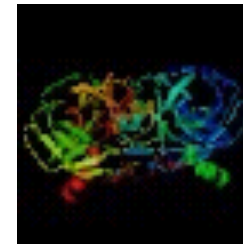
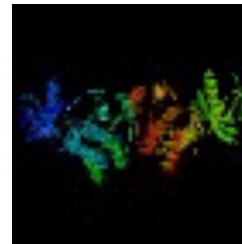
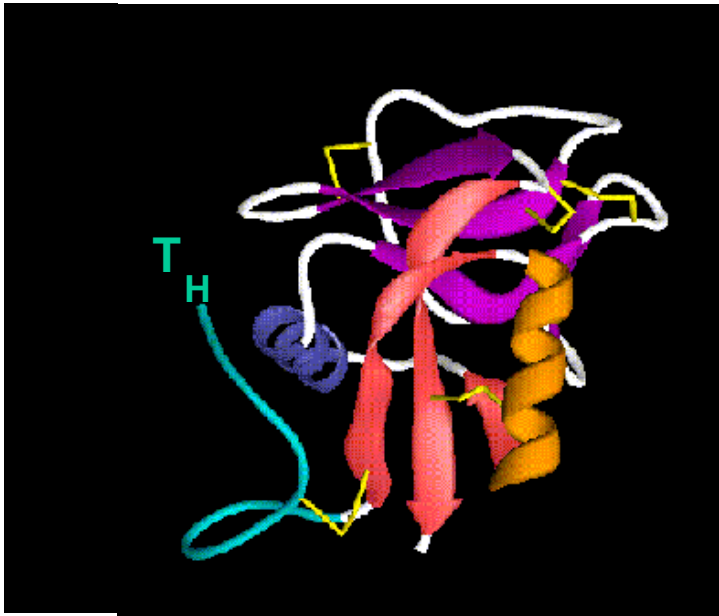
# Visualizing Threading

THREADINGSEQNCEECNQESGNI  
ERHTHREADINGSEQNCETHREAD  
GSEQNCEQCQESGIDAERTHR...



# Visualizing Threading

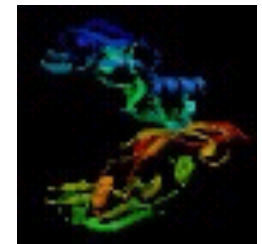
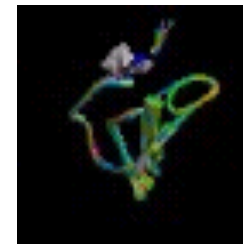
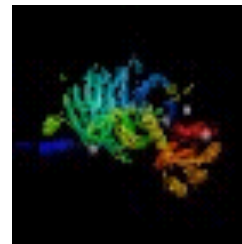
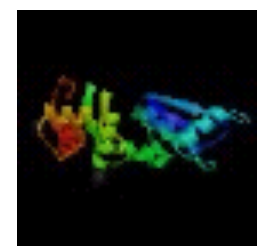
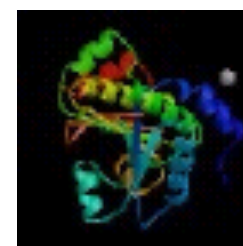
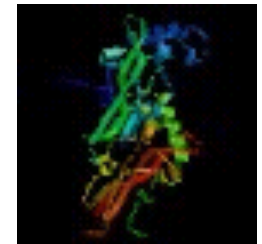
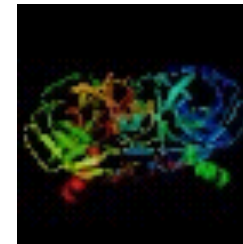
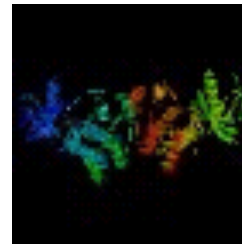
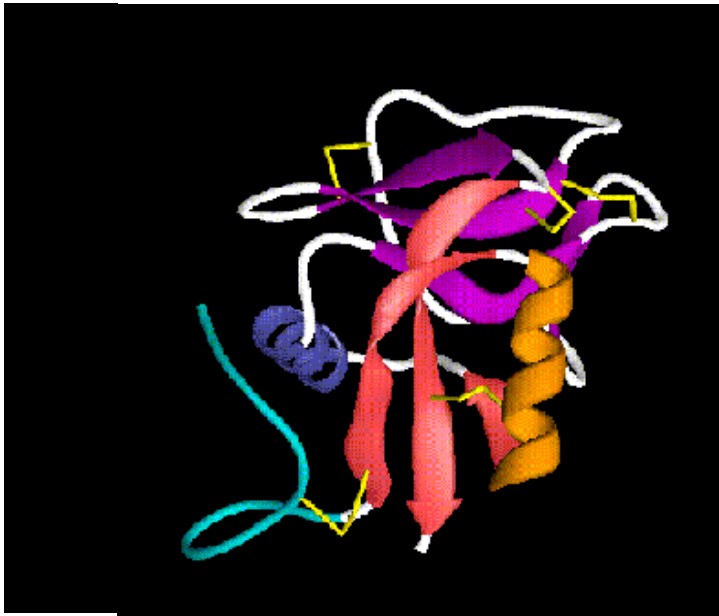
THREADINGSEQNCEECSNQGSGNI  
ERHTHREADINGSEQNCETHREAD  
GSEQNCEQCQESGIDAERTHR...



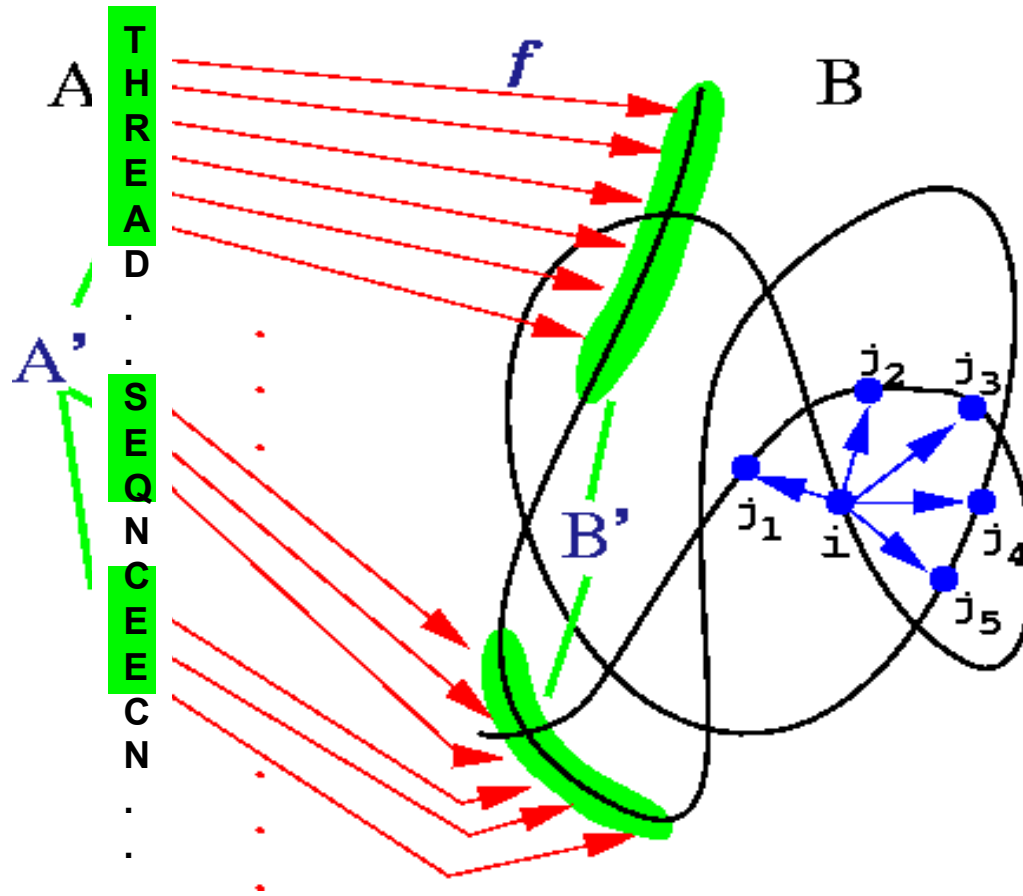


# Visualizing Threading

THREADINGSEQNCEECSNQESGNI  
ERHTHREADINGSEQNCETHREAD  
GSEQNCEQCQESGIDAERTHR...



# Visualizing Threading



# Threading\*

- **Database of 3D structures and sequences**
  - Protein Data Bank (or non-redundant subset)
- **Query sequence**
  - Sequence < 25% identity to known structures
- **Alignment protocol**
  - Dynamic programming
- **Evaluation protocol**
  - Distance-based potential or secondary structure
- **Ranking protocol**

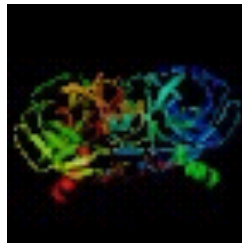
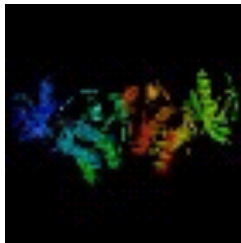
# 2 Kinds of Threading\*

- **2D Threading or Prediction Based Methods (PBM)**
  - Predict secondary structure (SS) or ASA of query
  - Evaluate on basis of SS and/or ASA matches
- **3D Threading or Distance Based Methods (DBM)**
  - Create a 3D model of the structure
  - Evaluate using a distance-based “hydrophobicity” or pseudo-thermodynamic potential

# 2D Threading Algorithm\*

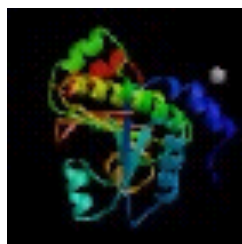
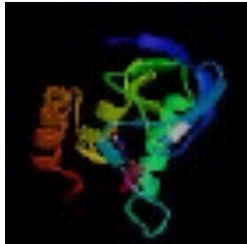
- *Convert PDB to a database containing sequence, SS and ASA information*
- *Predict the SS and ASA for the query sequence using a “high-end” algorithm*
- *Perform a dynamic programming alignment using the query against the database (include sequence, SS & ASA)*
- *Rank the alignments and select the most probable fold*

# Database Conversion



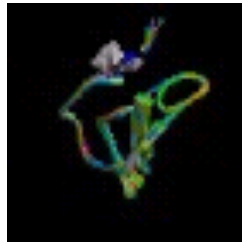
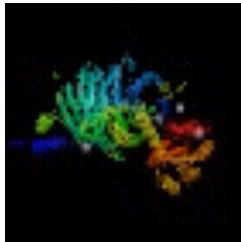
>Protein1

THREADINGSEQNCEECSGNI  
HHHHHHCCCCEEEECCHHHHHH  
ERHTHREADINGSEQNCETHREAD  
HHCCEEEECCECCCHHHHHHHHHH



>Protein2

QWETRYEWQEDFSHAECNQESGNI  
EEEEECCHHHHHHHHHHHHHHHH  
YTREWQHGFDSASQWETRA  
CCCCEEEECCEEEECCECC

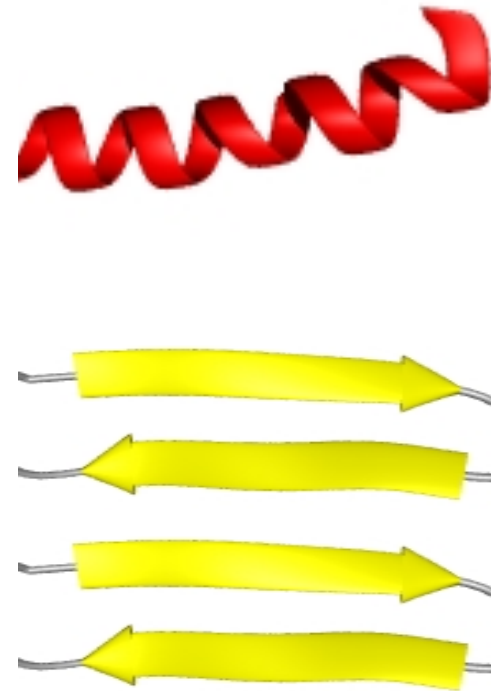


>Protein3

LKHGMNSNWEDFSHAECNQESG  
EEECCEEEECCEEEECCECC

# Secondary Structure

Phi & Psi angles for Regular Secondary Structure Conformations		
Structure	Phi ( $\Phi$ )	Psi ( $\Psi$ )
Antiparallel $\beta$ -sheet	-139	+135
Parallel $\beta$ -Sheet	-119	+113
Right-handed $\alpha$ -helix	+64	+40
$3_{10}$ helix	-49	-26
$\pi$ helix	-57	-70
Polyproline I	-83	+158
Polyproline II	-78	+149
Polyglycine II	-80	+150



# 2° Structure Identification\*

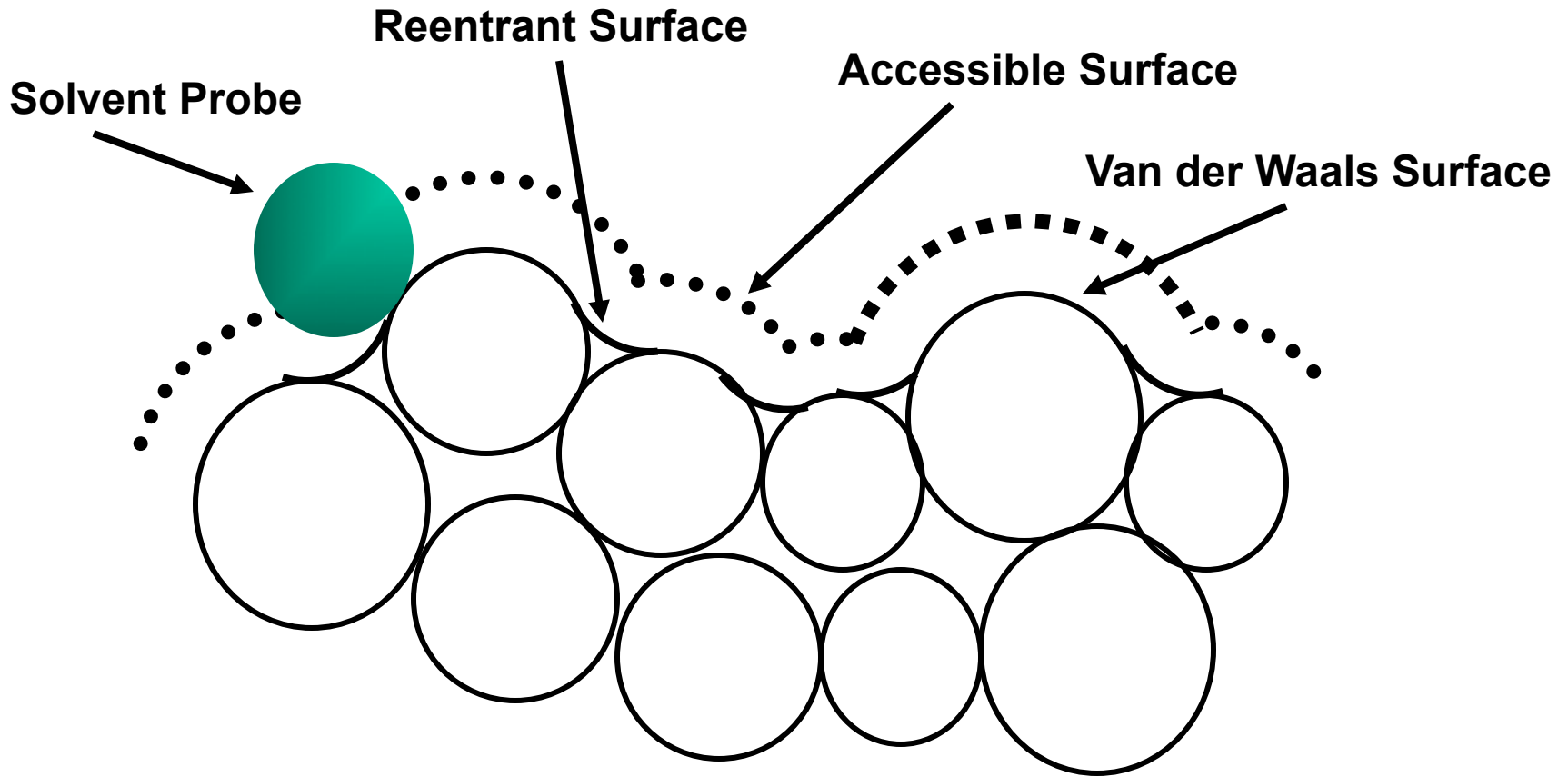
- **DSSP** - Database of Secondary Structures for Proteins (<http://swift.cmbi.ru.nl/gv/start/index.html>)
- **VADAR** - Volume Area Dihedral Angle Reporter (<http://vadar.wishartlab.com/>)
- **PDB** - Protein Data Bank ([www.rcsb.org](http://www.rcsb.org))
- **STRIDE** (<http://webclu.bio.wzw.tum.de/cgi-bin/stride/stridecgi.py>)



QHTAWCLTSEQHTAAVIWDCETPGKQNGAYQEDCA  
HHHHHCCEEEEEEEEEEEECCHHHHHHCCCCC



# Accessible Surface Area



# ASA Calculation\*

- **DSSP** - Database of Secondary Structures for Proteins (<http://swift.cmbi.ru.nl/gv/start/index.html>)
- **VADAR** - Volume Area Dihedral Angle Reporter (<http://vadar.wishartlab.com/>)
- **GetArea** - <http://curie.utmb.edu/getarea.html>



QHTAWCLTSEQHTAAVIWDCETPGKQNGAYQEDCAMD  
**BBPPBEEEEEPBPBPBPBBPEEEPBPEPEEEEEEEEEEE**  
10562987994152515104789414969899999999

# Other ASA sites

- **Connolly Molecular Surface Home Page**
  - <http://www.biohedron.com/>
- **Naccess Home Page**
  - <http://www.bioinf.manchester.ac.uk/naccess/>
- **MSMS**
  - [http://www.scripps.edu/~sanner/html/msms\\_home.html](http://www.scripps.edu/~sanner/html/msms_home.html)
- **Surface Racer**
  - <http://apps.phar.umich.edu/tsodikovlab/>

# 2D Threading Algorithm

- *Convert PDB to a database containing sequence, SS and ASA information*
- *Predict the SS and ASA for the query sequence using a “high-end” algorithm*
- *Perform a dynamic programming alignment using the query against the database (include sequence, SS & ASA)*
- *Rank the alignments and select the most probable fold*

# ASA Prediction\*

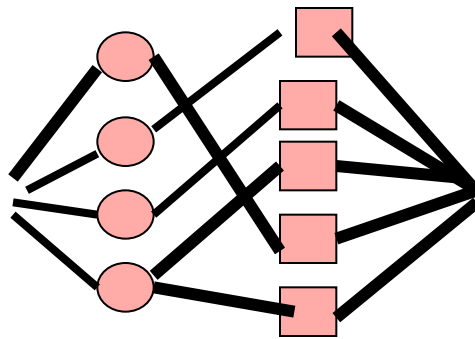
- **NetSurfP (70%)**

- <http://www.cbs.dtu.dk/services/NetSurfP/>

- **PredAcc (70%?)**

- <http://mobyli.rpbs.univ-paris-diderot.fr/cgi-bin/portal.py?form=PredAcc>

**QHTAW...**



**QHTAWCLTSEQHTAAVIW**  
**BBPPBEEEEEPBPBPBPB**

# 2D Threading Algorithm

- ***Convert PDB to a database containing sequence, SS and ASA information***
- ***Predict the SS and ASA for the query sequence using a “high-end” algorithm***
- ***Perform a dynamic programming alignment using the query against the database (include sequence, SS & ASA)***
- ***Rank the alignments and select the most probable fold***

# Dynamic Programming

	G	E	N	E	T	I	C	S
G	10	0	0	0	0	0	0	0
E	0	10	0	10	0	0	0	0
N	0	0	10	0	0	0	0	0
E	0	0	0	10	0	10	0	0
S	0	0	0	0	0	0	0	10
I	0	0	0	0	0	10	0	0
S	0	0	0	0	0	0	0	10

	G	E	N	E	T	I	C	S
G	60	40	30	20	20	0	10	0
E	40	50	30	30	20	0	10	0
N	30	30	40	20	20	0	10	0
E	20	20	20	30	20	10	10	0
S	20	20	20	20	20	0	10	10
I	10	10	10	10	10	20	10	0
S	0	0	0	0	0	0	0	10

G E N E T I C S  
 | | | | \* | |  
 G E N E S I S





# A Simple Example...\*

	A	T	V	D
A	1			
V				
V				
D				

	A	T	V	D
A	1	1		
V				
V				
D				

	A	T	V	D
A	1	1	0	0
V				
V				
D				

	A	T	V	D
A	1	1	0	0
V	0			
V				
D				

	A	T	V	D
A	1	1	0	0
V	0	1	1	
V				
D				

	A	T	V	D
A	1	1	0	0
V	0	1	1	2
V				
D				

# A Simple Example..\*

A A T V D  
 A 1 1 0 0 0  
 V 0 1 1 2 1  
 V  
 D

A A T V D  
 A 1 1 0 0 0  
 V 0 1 1 2 1  
 V 0 1 1 2 2  
 D 0 1 1 1 3

A A T V D  
 A 1 1 0 0 0  
 V 0 1 1 2 1  
 V 0 1 1 2 2  
 D 0 1 1 1 3

A A T V D  
 | | | |  
 A - V V D

A A T V D  
 | | | |  
 A V V D

A A T V D  
 | | | |  
 A V - V D

# Let's Include 2° info & ASA\*

	H	E	C		E	P	B		
$S_{ij}^{\text{strc}}$	H	1	0	0	$S_{ij}^{\text{asa}}$	E	1	0	0
	E	0	1	0		P	0	1	0
	C	0	0	1		B	0	0	1

$$S_{ij}^{\text{total}} = k_1 S_{ij}^{\text{seq}} + k_2 S_{ij}^{\text{strc}} + k_3 S_{ij}^{\text{asa}}$$

# A Simple Example...\*

**E E E C C**  
**A A T V D**  
EA **2**  
EV  
CV  
CD

**E E E C C**  
**A A T V D**  
EA **2 2**  
EV  
CV  
CD

**E E E C C**  
**A A T V D**  
EA **2 2 1 0 0**  
EV  
CV  
CD

**E E E C C**  
**A A T V D**  
EA **2 2 1 0 0**  
EV **1**  
CV  
CD

**E E E C C**  
**A A T V D**  
EA **2 2 1 0 0**  
EV **1 3 3**  
CV  
CD

**E E E C C**  
**A A T V D**  
EA **2 2 1 0 0**  
EV **1 3 3 3**  
CV  
CD

# A Simple Example...

	<b>E E E C C</b>				
	<b>A A T V D</b>				
<b>E A</b>	<b>2 2 1 0 0</b>				
<b>E V</b>	<b>1 3 3 3 2</b>				
<b>C V</b>					
<b>C D</b>					

	<b>E E E C C</b>				
	<b>A A T V D</b>				
<b>E A</b>	<b>2 2 1 0 0</b>				
<b>E V</b>	<b>1 3 3 3 2</b>				
<b>C V</b>	<b>0 2 3 5 4</b>				
<b>C D</b>	<b>0 2 3 4 7</b>				

	<b>E E E C C</b>				
	<b>A A T V D</b>				
<b>E A</b>	<b>2 2 1 0 0</b>				
<b>E V</b>	<b>1 3 3 3 2</b>				
<b>C V</b>	<b>0 2 3 5 4</b>				
<b>C D</b>	<b>0 2 3 4 7</b>				



# 2D Threading Performance

- In test sets 2D threading methods can identify 30-40% of proteins having very remote homologues (i.e. not detected by BLAST) using “minimal” non-redundant databases (<700 proteins)
- If the database is expanded ~4x the performance jumps to 70-75%
- Performs best on true homologues as opposed to postulated analogues

# 2D Threading Advantages\*

- **Algorithm is easy to implement**
- **Algorithm is very fast (10x faster than 3D threading approaches)**
- **The 2D database is small (<500 kbytes) compared to 3D database (>1.5 Gbytes)**
- **Appears to be just as accurate as DBM or other 3D threading approaches**
- **Very amenable to web servers**

Secondary structure element alignment

http://protein.cribi.unipd.it/ssea/

Most Visited Getting Started Latest Headlines

Secondary structure element align... +



## Secondary Structure Element Alignment (SSEA)

Database alignment mode

One vs. One Alignment

Quick Help and References

Scoring and Benchmarking

Download the Software

Please see our [acceptable usage policy](#).

Name of sequence (optional)

Fold Library

PDB 95 fold library

Secondary Structure

Options:

Predict secondary structure from sequence

Local Alignment

Global Alignment

Z-Score Cutoff:

Submit Query

Reset

(c) [Silvio Tosatto](#) for [BioComputing GRUP](#) 08 / 2004

Done

<http://protein.cribi.unipd.it/ssea/>



PDBeFold

http://www.ebi.ac.uk/msd-srv/ssm/

EMBL-EBI

Enter Text Here Find Help | Feedback

Databases Tools Research Training Industry About Us Help Site Index

**PDB**  
HOME

**Structure Similarity**  
pdbe.org/fold

Bringing Structure to Biology  
FEEDBACK

### Documentation and links

#### PDBeFold links

- o Tips
- o Visualisation
- o Performance
- o Privacy
- o FAQs
- o Version log
- o Fold Links
- o Comparisons
- o Publications
- o PDBeFOLD tutorial

#### Other links

- o PDBePISA
- o CCP4 CoordLib
- o Rasmol
- o Rastop
- o Jmol
- o PDB
- o SCOP
- o PDBeMotif
- o GeneCensus
- o FSSP
- o CATH
- o PDBSum
- o UniProt

### PDBeFold (Structure Similarity)

PDBeFold (SSM) is an interactive service for comparing protein structures in 3D.

#### PDBeFold functionality:

- o pairwise comparison and 3D alignment of protein structures
- o multiple comparison and 3D alignment of protein structures
- o examination of a protein structure for similarity with the whole PDB archive or SCOP archive
- o best Ca-alignment of compared structures
- o download and visualisation of best-superposed structures using Rasmol (Unix/Linux platforms), Rastop (Windows machines) and Jmol(platform-independent server-side java viewer)
- o linking the results to other services - PDBeMotif, SCOP, GeneCensus, FSSP, CATH, PDBSum, UniProt

**Launch PDBeFold**

PDBeFold: A comparison with other protein matching services.  
PDBeFold is used as a structure search engine in PDBePISA.  
PDBeFold queries may be launched from any web site (instructions).  
PDBeFold is based on the CCP4 Coordinate Library.

We are having hardware issues that may occasionally affect the PISA and SSM services.  
If you are experiencing problems with the service please wait for 20 minutes and try again.  
We apologise for any inconvenience.

**We welcome your feedback! Please send any questions, comments, suggestions and bug reports using the FEEDBACK button on the top of the page.**

Terms of Use | EBI Funding | Contact EBI | © European Bioinformatics Institute 2011. EBI is an Outstation of the European Molecular Biology Laboratory.

W3C XHTML 1.0

Done

http://www.ebi.ac.uk/msd-srv/ssm/

# Servers - HHPred

The screenshot shows a web browser window with the address bar displaying `http://toolkit.tuebingen.mpg.de/hhpred`. The page title is "HHpred - Homology detection & structure prediction by HMM-HMM comparison". The browser's address bar also shows a Google search engine icon.

The website header includes navigation links: HOME, Login, PDBAlert, Personal Databases, Contact, Imprint, Disclaimer, and Help. The main heading is "Bioinformatics Toolkit" from the "Max-Planck Institute for Developmental Biology". A "Quickfinder" search box is located in the top right.

The main navigation menu includes: Search, Alignment, Sequence Analysis, 2ary Structure, 3ary Structure, Classification, and Utils. Below this, a secondary menu lists various tools: CS-BLAST, FHMMER, HHpred, HHSenser, NucBLAST, PSI-BLAST, PatternSearch, ProtBLAST, and SimShiftDB.

The main content area is titled "HHpred - Homology detection & structure prediction by HMM-HMM comparison" with a "Help" link. A note states "HHpred now runs with HHsearch 1.6.0.0".

The "Input" section contains a large text area for "Paste protein sequence or multiple alignment", a "Browse..." button for "or upload a local file", and a "Select input format" dropdown menu currently set to "FASTA". "Reset form" and "Submit job" buttons are located at the bottom right of the input section.

The "Search Options" section includes:

- "Select HMM databases (hold Ctrl to select several)" with two scrollable lists:
  - Standard:** pdb70\_24Oct09, pdb\_on\_hold\_23Oct09, scop70\_1.71, scop70\_1.75, cdd\_23Oct09
  - Genomes:** Arabidopsis\_thaliana, Drosophila\_melanogaster, Homo\_sapiens, Mus\_musculus, Plasmodium\_falciparum
- "Max. PSI-BLAST iterations" set to 8
- "Score secondary structure" with radio buttons for "yes", "no", and "predicted vs predicted only" (selected)
- "Alignment mode" with radio buttons for "local" (selected) and "global"
- "Realign with MAC" checkbox (unchecked)

The left sidebar contains a "MAX-PLANCK-GESELLSCHAFT" logo, a "Show results of job:" section with a "Show results" button, and a "Recent jobs:" section with "Select all", "Deselect all", "Clear sel. Jobs", and "Delete sel. Jobs" buttons. Below this is a status table with columns for "queued", "running", "done", and "error".

<http://toolkit.tuebingen.mpg.de/hhpred>

# Servers - GenThreader

psipred : index

http://bioinf.cs.ucl.ac.uk/psipred/

UCL Department Of Computer Science  
Bioinformatics Group

Search Group

UCL Home >> Departments of Computer Science >> Bioinformatics Group >> psipred

## The PSIPRED Protein Structure Prediction Server

The PSIPRED Protein Structure Prediction Server aggregates several of our structure prediction methods into one location. Users can submit a protein sequence, perform the prediction of their choice and receive the results of the prediction via e-mail. You may select one of three prediction methods to apply to your sequence:

- PSIPRED - a highly accurate method for protein secondary structure prediction
- MEMSAT and MEMSAT-SVM - our widely used transmembrane topology prediction method
- and one of GenTHREADER, pGenTHREADER and pDomTHREADER - sequence profile based fold recognition methods. [More...](#)

**For queries regarding PSIPRED:** [psipred@cs.ucl.ac.uk](mailto:psipred@cs.ucl.ac.uk)

### Choose Prediction Method

- Predict Secondary Structure (PSIPRED v3.0)
- Predict Transmembrane Topology (MEMSAT2 & MEMSAT-SVM)
- SVM Prediction of TM Topology and Helix Packing (MEMPACK) - **NEW!**
- Fold Recognition (GenTHREADER - quick)
- Fold Recognition (pGenTHREADER - quick) - requires pre-computed secondary structure)
- Fold Recognition (pDomTHREADER - annotates multiple domain on chains)

[Help...](#)

### Input Sequence (single letter amino acid code)

[Help...](#)

If you wish to test these services follow this link to retrieve a [test fasta sequence](#).

### Filtering Options

<http://bioinf.cs.ucl.ac.uk/psipred/>

# 2D Threading Disadvantages\*

- **Reliability is not 100% making most threading predictions suspect unless experimental evidence can be used to support the conclusion**
- **Does not produce a 3D model at the end of the process**
- **Doesn't include all aspects of 2° and 3° structure features in prediction process**
- **PSI-BLAST may be just as good (faster too!)**

# Making it Better

- **Include 3D threading analysis as part of the 2D threading process -- offers another layer of information**
- **Include more information about the “coil” state (3-state prediction isn't good enough)**
- **Include other biochemical (ligands, function, binding partners, motifs) or phylogenetic (origin, species) information**

# 3D Threading Servers

- **Generate 3D models or coordinates of possible models based on input sequence**
- **Loopp (version 4)**
  - <http://cbsuapps.tc.cornell.edu/loopp.aspx>
- **Phyre**
  - <http://www.sbg.bio.ic.ac.uk/~phyre/index.cgi>
- **All require email addresses since the process may take hours to complete**

# phyre

Version 0.2

Protein Homology/analogY Recognition Engine

## New Phyre server scores highly in CASP8 competition. [Results](#)

Phyre has been highlighted in [June 2009 Nature PSI Knowledgebase](#)

The Phyre webserver is for **Academic use only**  
For in-house and/or commercial use please click [here](#)

[Other tools available from our lab \(function prediction, docking, etc.\)](#)

E-mail Address

Optional Job  
description

Amino Acid  
Sequence

Quick Phyre Search

**APPLICATIONS**

(click on a category below to access programs)

Show all Hide all

- Sequence analysis
- Sequence alignment
- Population genetics
- Protein structure
- LOOPP**
- MODELLER**
- MSR Biomedical
- Other
- Links

# LOOPP @ BioHPC

version 4.0

The **LOOPP** (Learning, Observing and Outputting Protein Patterns) server. LOOPP is a fold recognition program based on the collection of numerous signals, merging them into a single score, and generating atomic coordinates based on an alignment into a homologue template structure. The signals we are using include straightforward sequence alignment, sequence profile, threading, secondary structure and exposed surface area prediction. For more information please refer to the [LOOPP home page](#).

**LOOPP HAS BEEN UPGRADED.** We have upgraded LOOPP to the newest version just in time for CASP8. The new version is much better than the previous one, it uses new scoring technique as well as upgraded database. The improvements are described in the upcoming *Proteins* paper.

**NOTE ABOUT CASP8.** This server is participating in CASP8 experiment. All LOOPP predictions for CASP8 targets with the default server parameters are available online [here](#). Please don't submit CASP8 targets again with default parameters - just use the link!

If you have any comments or questions about LOOPP please contact us at [loopp@tc.cornell.edu](mailto:loopp@tc.cornell.edu).

It may take several hours to run this program.

Calculations will be carried out on the BioHPC compute cluster at [CBSU](#). You will receive e-mail notifications when the job is submitted, when it starts, and when it is finished. Output will be available via links embedded in the notification e-mails. For more information about this program and BioHPC interface in general, please visit our [Frequently Asked Questions](#) page.

Please acknowledge us in all publications and presentation of work that used our resources using the following [text](#).

E-mail:  (only guests need to use this field, registered users should log in)

Job name:  (please, no spaces, special characters etc., underscore is OK)

Please notify me about LOOPP updates

Input sequence (amino acids one-letter code only, no names, numbers etc):

**MISCELLANEOUS**

- [Subscribe](#)
- [Apps Home](#)
- [Clusters Status](#)
- [Applications Statistics](#)
- [BioHPC Home](#)
- [CBSU Home](#)
- [CBSU ftp server](#)
- [CBSU SeqDB](#)
- [CTC Windows Bioinformatics Applications](#)
- [DISTRUCT](#)
- [T-REX \(T-RFLP manager\)](#)



# Outline

- **Secondary Structure Prediction**
- **Threading (1D and 3D threading)**
- **Ab initio Structure Prediction**

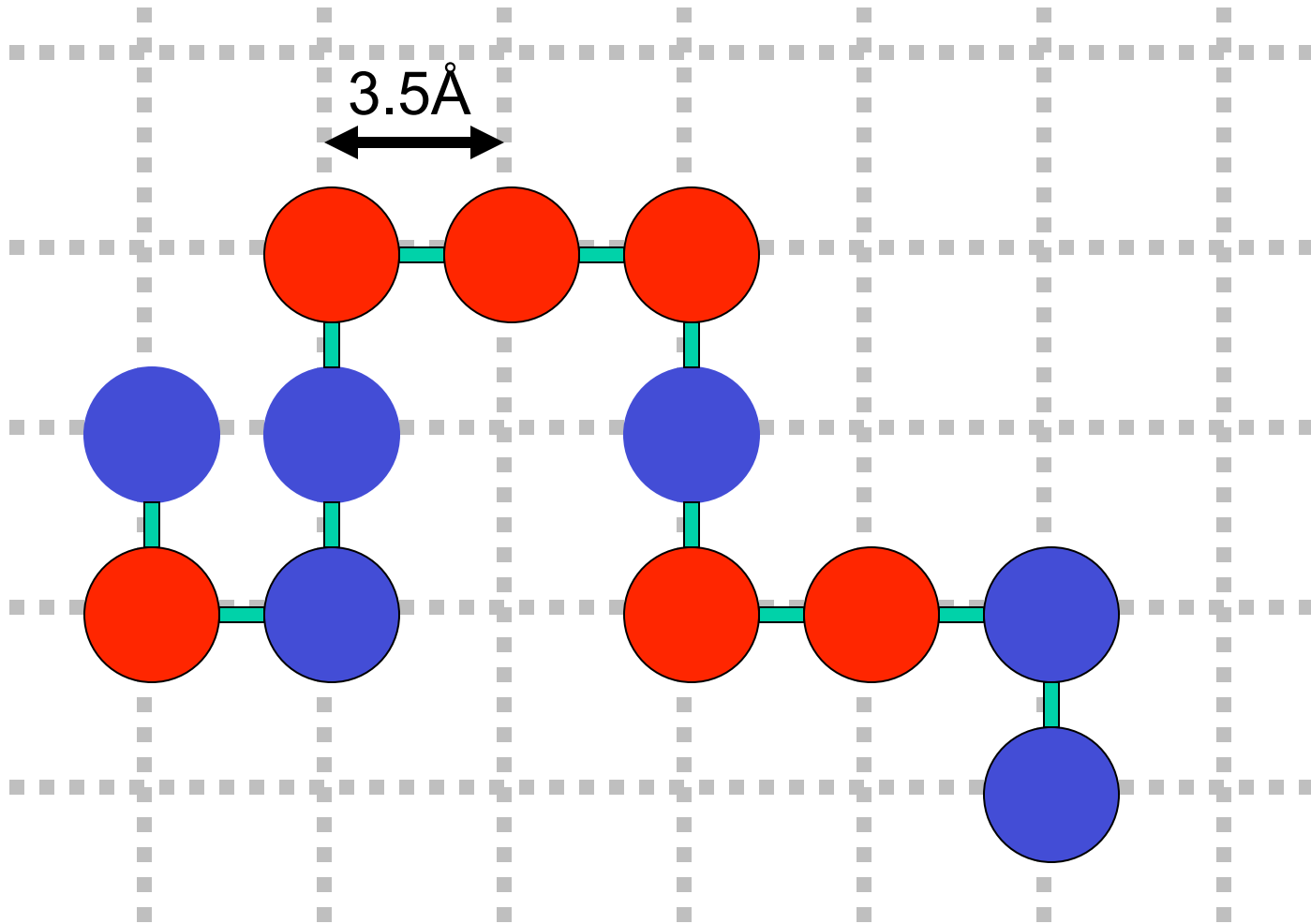
# Ab Initio Prediction\*

- **Predicting the 3D structure without any “prior knowledge”**
- **Used when homology modelling or threading have failed (no homologues are evident)**
- **Equivalent to solving the “Protein Folding Problem”**
- **Still a research problem**

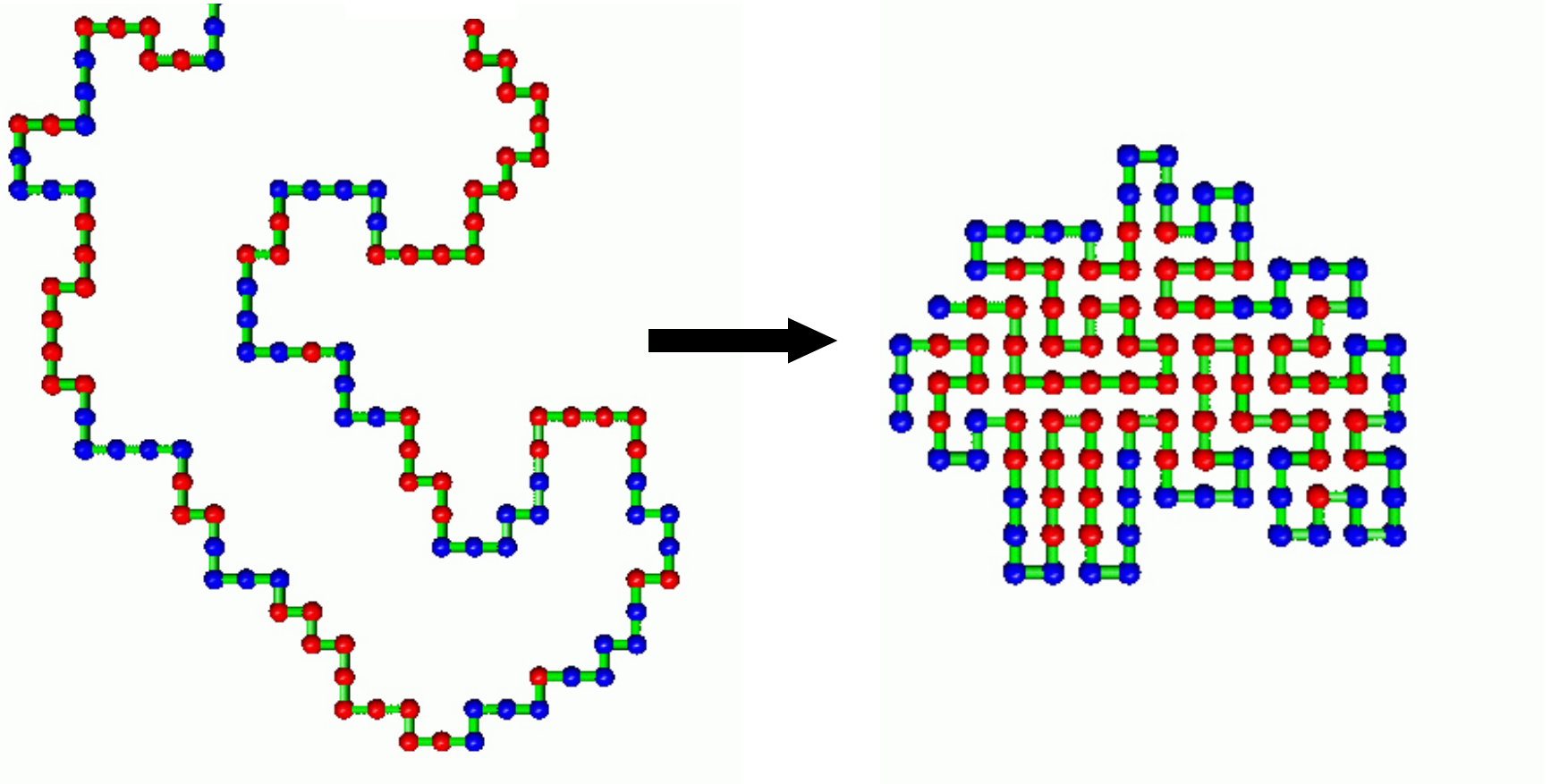
# Ab Initio Folding\*

- **Two Central Problems**
  - Sampling conformational space ( $10^{100}$ )
  - The energy minimum problem
- **The Sampling Problem (Solutions)**
  - Lattice models, off-lattice models, simplified chain methods, parallelism
- **The Energy Problem (Solutions)**
  - Threading energies, packing assessment, topology assessment

# A Simple 2D Lattice



# Lattice Folding



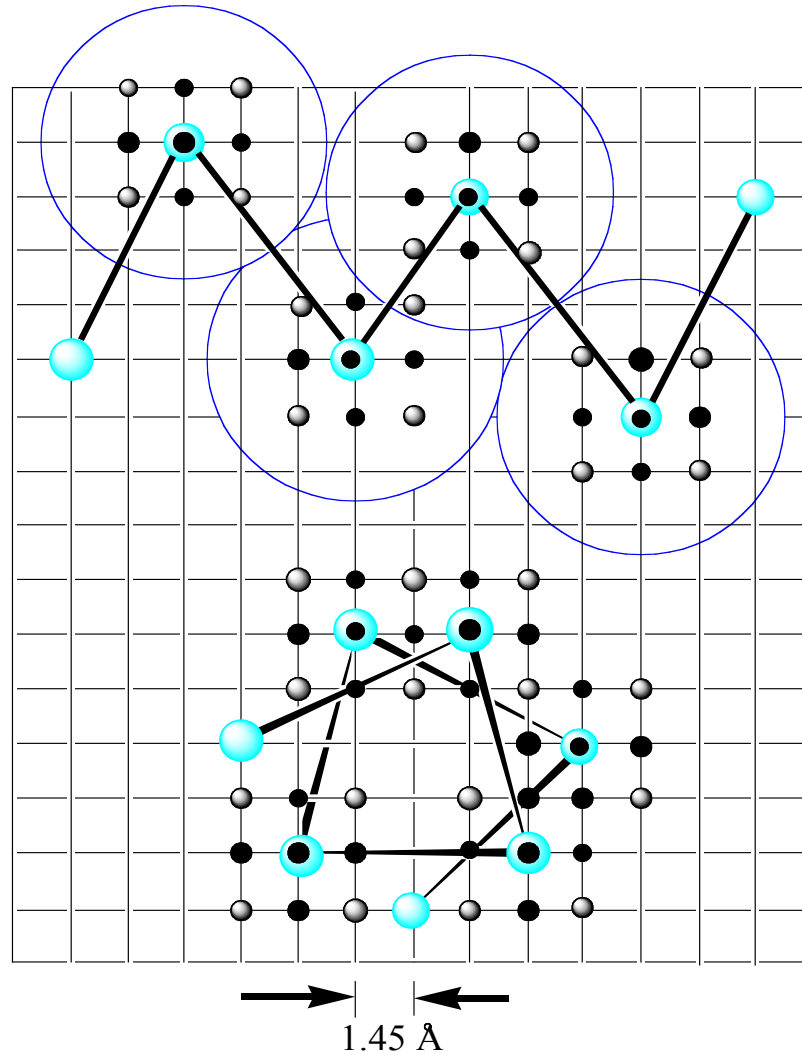
# Lattice Algorithm

- ***Build a “n x m” matrix (a 2D array)***
- ***Choose an arbitrary point as your N terminal residue (start residue)***
- ***Add or subtract “1” from the x or y position of the start residue***
- ***Check to see if the new point (residue) is off the lattice or is already occupied***
- ***Evaluate the energy***
- ***Go to step 3) and repeat until done***

# Lattice Energy Algorithm

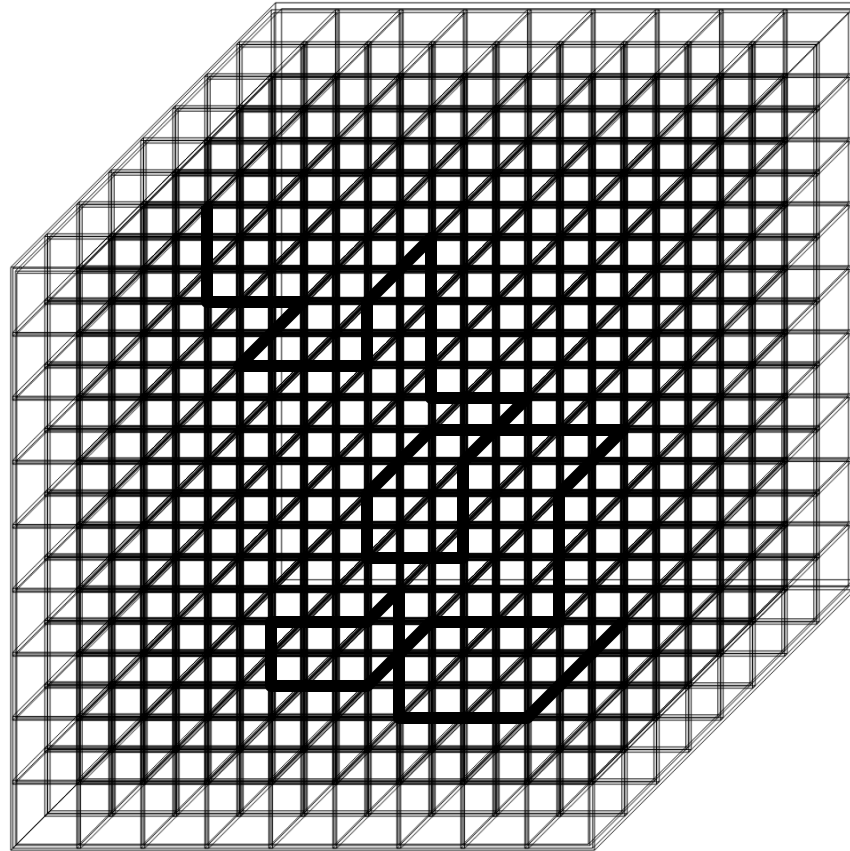
- *Red = hydrophobic, Blue = hydrophilic*
- *If Red is near empty space  $E = E+1$*
- *If Blue is near empty space  $E = E-1$*
- *If Red is near another Red  $E = E-1$*
- *If Blue is near another Blue  $E = E+0$*
- *If Blue is near Red  $E = E+0$*

# More Complex Lattices

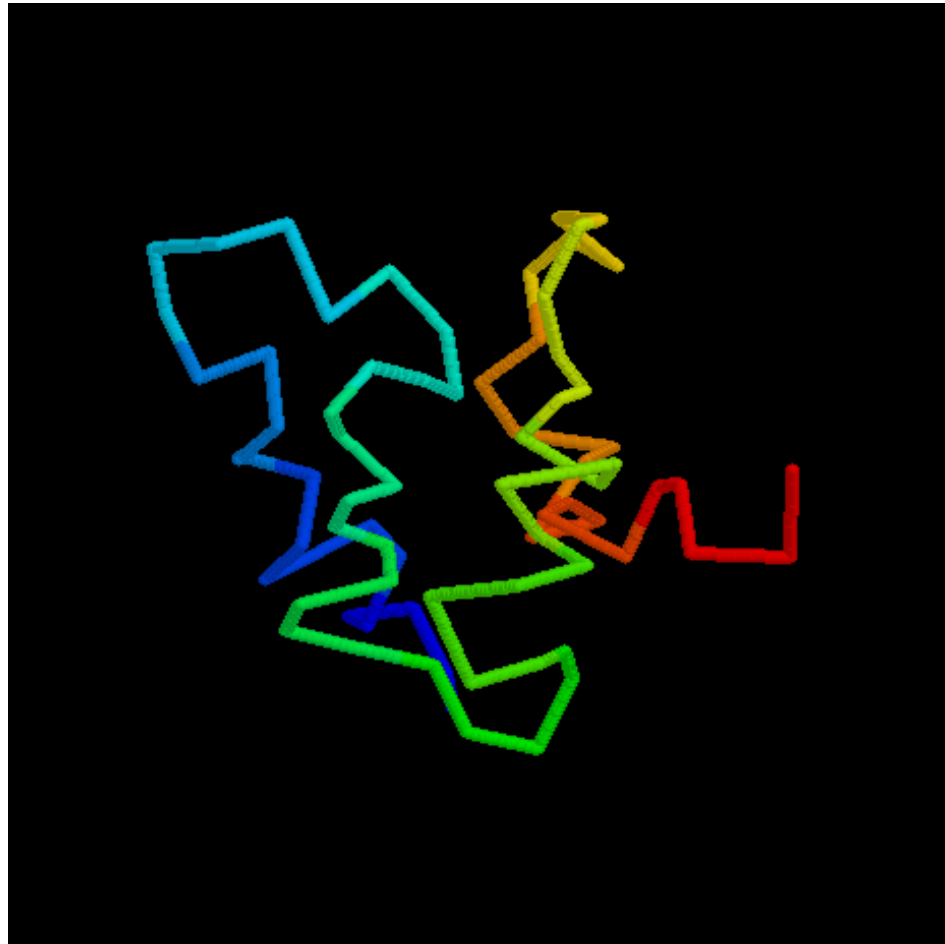




# 3D Lattices



# Really Complex 3D Lattices



J. Skolnick

# Lattice Methods\*

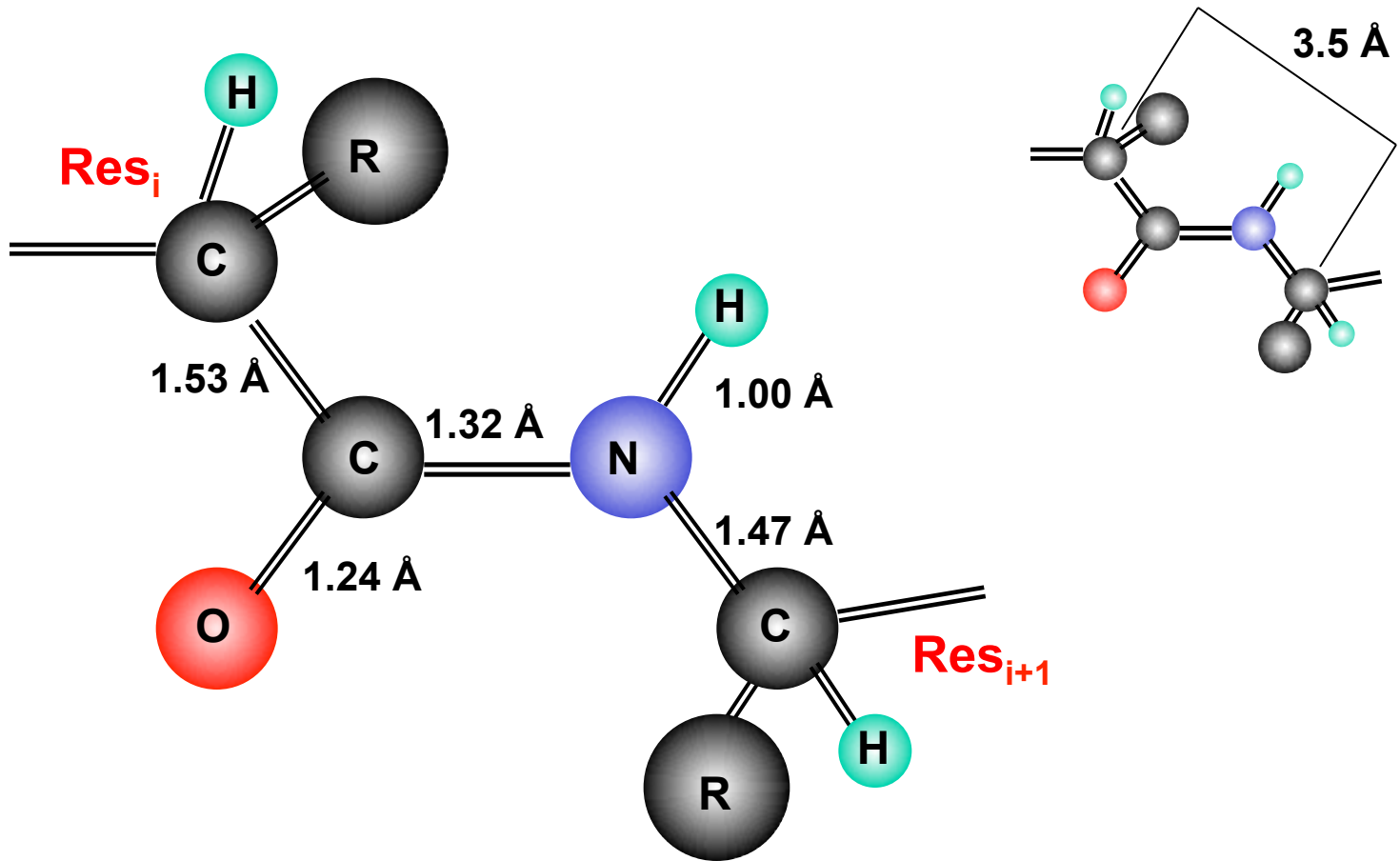
## Advantages

- Easiest and quickest way to build a polypeptide
- Implicitly includes excluded volume
- More complex lattices allow reasonably accurate representation

## Disadvantages

- At best, only an approximation to the real thing
- Does not allow accurate constructs
- Complex lattices are as “costly” as the real thing

# Non-Lattice Models



# Best Method So Far...\*

Target 74



native



model 4

Target 77



native



model 4

Target 79



native



model 4

Target 56



native



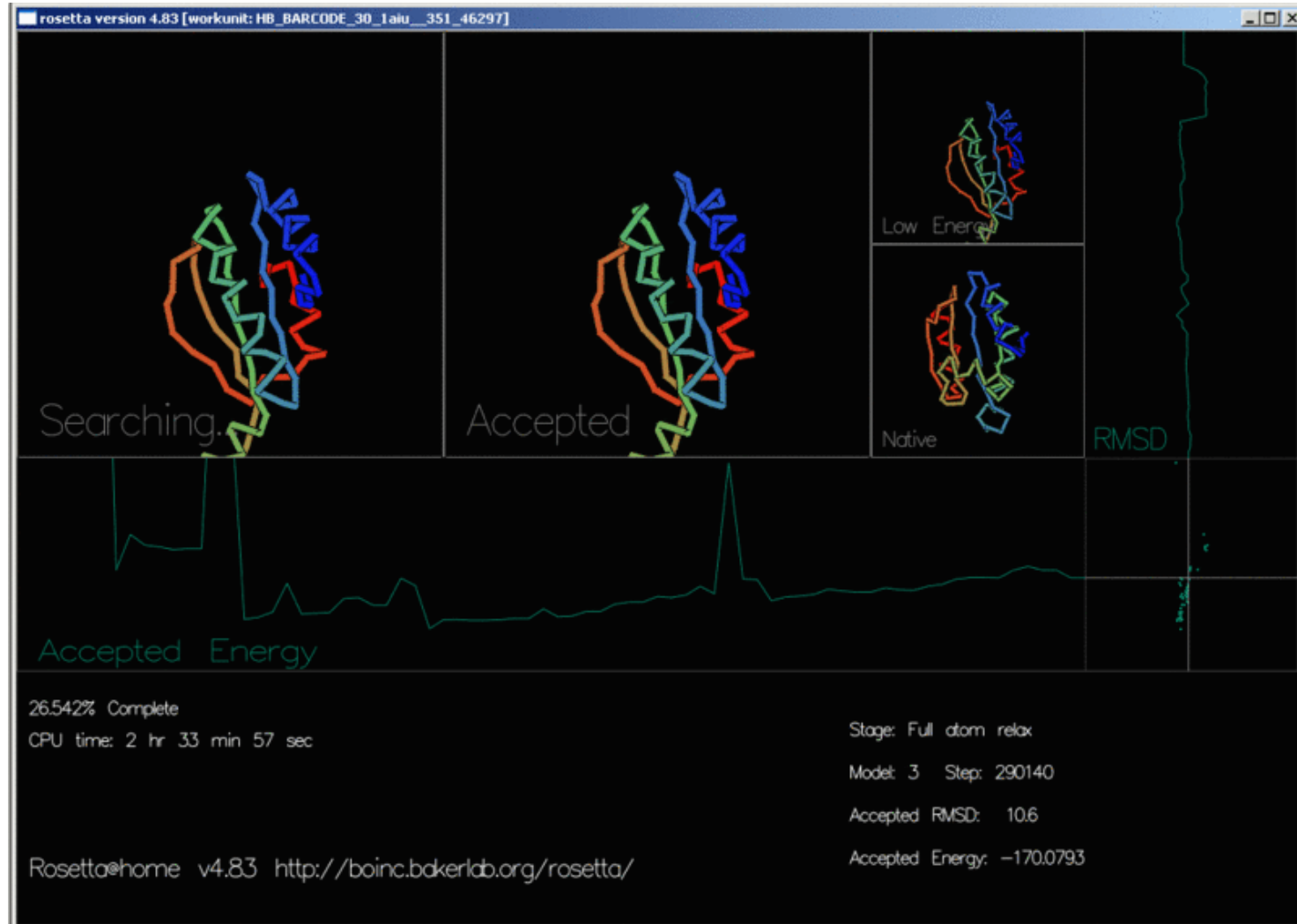
model 4

**Rosetta - David Baker**

# Rosetta Outline\*

- **Assembles proteins using “fragment assembly” of known protein fragments**
- **Fragments are 3-9 residues long**
- **Fragments identified via PSI-BLAST**
- **Starts with extended chain and then randomly changes conformation of selected regions based on fragment matches**
- **Evaluates energy using Monte Carlo**

# Rosetta in Action



# Robetta & Rosetta

The screenshot shows a Microsoft Internet Explorer browser window displaying the Robetta website. The browser's address bar shows the URL <http://robetta.bakerlab.org/>. The website header includes the Robetta logo and the text "Full-chain Protein Structure Prediction Server". The main content area is divided into two columns. The left column features two protein structure models: "T134: Homology Modeling" and "T148: De Novo". Each model is shown in two views: "Model 1" and "Native". The "T134: Homology Modeling" models are labeled with "dom 1" and "dom 2" and show N and C termini. The "T148: De Novo" models also show "dom 1" and "dom 2" and N and C termini. Below these models, the text reads "examples of predictions by Robetta in CASP-5". The right column contains navigation links under three main sections: "REGISTRATION" with links for "[ Register / Update ]" and "[ Login ]"; "DOCUMENTATION" with links for "[ Docs / FAQs ]" and "[ News ]"; and "SERVICES" with sub-sections: "Domain Parsing & 3-D Modeling" (homology modeling, *ab initio* structure prediction, and structure prediction using NMR constraints) with links for "[ Submit ]" and "[ Queue ]"; "Interface Alanine Scanning" with links for "[ Submit ]" and "[ Queue ]"; and "Fragment Libraries" with links for "[ Submit ]" and "[ Queue ]". At the bottom of the right column, there is a "RELATED SITES" section with a link for "Rosetta Design Server". The browser's status bar at the bottom shows "Done" and "Internet". The Windows taskbar at the very bottom displays the Start button and several open applications: "CBW Proteomics\_5...", "Proteomics2006", "1.2Proteinfeatures...", "Proteomics3.2.ppt", and "Robetta: full-chain ...". The system clock shows "11:31 PM".

<http://robetta.bakerlab.org/>



# Robetta

- **Allows users predict 3D structures using Rosetta ab-initio method and to do homology modelling too**
- **Requires considerable computational resources (now hosted at Los Alamos supercomputer facility)**
- **Requires that users register and login (to track mis-use and abuse)**

# Another Approach...

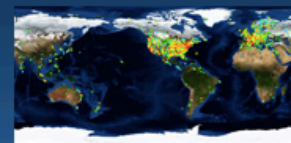
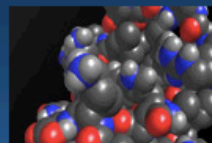
## Distributed Folding

- **Attempt to harness the same computational power as BlueGene but by doing on thousands of PC's via a screen saver**
- **Three efforts underway:**
  - <http://folding.stanford.edu/>
  - <http://boinc.bakerlab.org/rosetta/>
- **You can be part of this exp' t too!**



# Folding@home

## distributed computing



Home Download Guides FAQ Stats Science Results Awards About Us

Main  
News  
Forum  
(Help)  
Donate

العربية  
中文  
Dansk  
Deutsch  
English  
Español  
Français  
Italiano  
Ελληνικά  
日本語  
한국말  
Lietuvių  
Magyar  
Nederlands  
Norsk  
Occitan  
فارسی  
Polski  
Português  
Русский  
Suomeksi  
Svenska  
Turkish  
Tiếng Việt

## Our goal: to understand protein folding, misfolding, and related diseases

### What is protein folding?

Proteins are biology's workhorses -- its "nanomachines." Before proteins can carry out these important functions, they assemble themselves, or "fold." The process of protein folding, while critical and fundamental to virtually all of biology, in many ways remains a mystery.

[Download Folding@home](#)

### Protein folding is linked to disease, such as Alzheimer's, ALS, Huntington's, Parkinson's disease, and many Cancers

Moreover, when proteins do not fold correctly (i.e. "misfold"), there can be serious consequences, including many well known diseases, such as Alzheimer's, Mad Cow (BSE), CJD, ALS, Huntington's, Parkinson's disease, and many Cancers and cancer-related syndromes.

### You can help scientists studying these diseases by simply running a piece of software.

Folding@home is a distributed computing project -- people from throughout the world [download](#) and run software to band together to make one of the largest supercomputers in the world. Every computer takes the project closer to our goals. Folding@home uses novel computational methods coupled to distributed computing, to simulate problems millions of times more challenging than previously achieved.

### What have we done so far?

We have had several successes. You can read about them on our [Science](#) page, on our [Awards](#) page, or go directly to our [Results](#) page.



### Want to learn more?

Click on the links on the left for downloads or more information. You can also download our [Executive Summary](#), which is a PDF suitable for distribution. One can also help by [donating funds to the project](#), via Stanford University.

Make a gift now!   
[Click here to make a gift >](#)

# Rosetta@home

Protein Folding, Design, and Docking



## What is Rosetta@home?



**Rosetta@home** needs your help to determine the 3-dimensional shapes of proteins in research that may ultimately lead to finding cures for some major human diseases. By running the Rosetta program on your computer while you don't need it you will help us speed up and extend our research in ways we couldn't possibly attempt without your help. You will also be helping our efforts at designing new proteins to fight diseases such as HIV, Malaria, Cancer, and Alzheimer's (See our [Disease Related Research](#) for more information). Please [join us](#) in our efforts! **Rosetta@home is not for profit.**

[ login/out ]

Site search

### Join Rosetta@home

1. [Rules and policies](#)
2. [System requirements](#)
3. [Download, install, and run BOINC](#)  
(enter the project URL: <http://boinc.bakerlab.org/rosetta/>)
4. [A welcome from David Baker](#)
5. [Donate](#)

### About

- [10 reasons why users crunch Rosetta@home](#)
- [Quick Guide to Rosetta@home and Its Graphics](#)
- [Play the interactive rosetta game, FoldIt!](#)
- [Rosetta@home FAQ](#)
- [Rosetta@home Science FAQ](#)
- [Disease Related Research](#)
- [Research Overview](#)
- [Publications](#)
- [News & Articles about Rosetta](#)
- [David Baker's Rosetta@home Journal](#)
- [Rosetta@home promo video](#)
- [Technical news](#)

### Returning participants

- [Your account](#) - view stats, modify preferences
- [Results](#) - view your results
- [Teams](#) - create or join a team
- [Applications](#)

### User of the day



[Herbert surft](#)

...nur damit die Rechner auch was zu tun haben...

### Server Status as of 26 Oct 2009 22:52:09 UTC

[ Scheduler running ]

Total queued jobs: **360,803**

In progress: 396,146

Successes last 24h: 232,018

Users [↓](#) (last day [↓](#)) : 269,289 (+79)

Hosts [↓](#) (last day [↓](#)) : 804,167 (+380)

Credits last 24h [↓](#) : 9,712,698

Total credits [↓](#) : 8,333,557,026

TeraFLOPS estimate: 97.127

Oct 26, 2009

**Predictor of the day:** Congratulations to [dookie\\_2k3](#) (Team [Asturias-Team](#)) for predicting the lowest energy structure for workunit [lr8\\_combine\\_smooth\\_torsion\\_it00\\_rama03\\_A\\_rln\\_1tig\\_14898\\_0](#) !

[...more](#)

[XML](#) Available as an [RSS feed](#)..

### News

Oct,14, 2009

The minirosetta application has been updated to version 1.98. For details and to report bugs, go to [this thread](#).

Sep 16, 2009

Our filesystem became bogged down late last night. Thanks to Keith, our systems administrator, the project is back online.

Sep 11, 2009

# D.W. Shaw Research Institute (MD for 3D Structure Prediction)



D. E. Shaw Research

www.deshawresearch.com/chiefscientist.html

About D. E. Shaw Research Chief Scientist Members of the Lab Joining the Lab Publications Resources Contacting Us

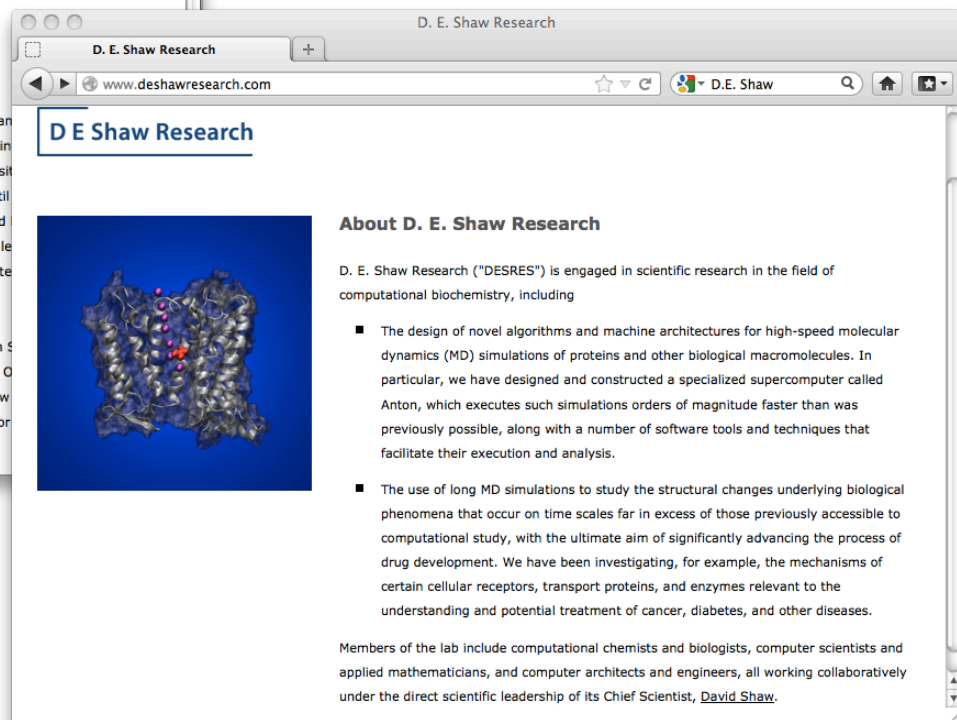
## D E Shaw Research



**Chief Scientist**

David E. Shaw serves as Chief Scientist of D. E. Shaw Research and Research Fellow at the Center for Computational Biology and Bioinformatics at Columbia University. He received his Ph.D. from Stanford University and was a member of the faculty of the Computer Science Department at Columbia until he joined the D. E. Shaw group in 1988. Since 2001, Dr. Shaw has devoted his research efforts in the field of computational biochemistry. Although he has been actively involved in research efforts in his role as Chief Scientist, his focus is largely on the involvement in operational and administrative management.

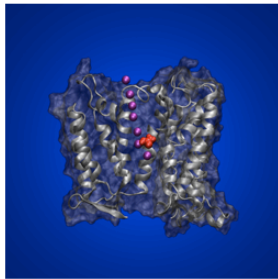
Dr. Shaw was appointed to the President's Council of Advisors on Science and Technology by President Clinton in 1994, and again by President Obama in 2009. He is a member of the National Academy of Engineering, and is a fellow of the National Academy of Arts and Sciences and of the American Association for Artificial Intelligence.



D. E. Shaw Research

www.deshawresearch.com

## D E Shaw Research



**About D. E. Shaw Research**

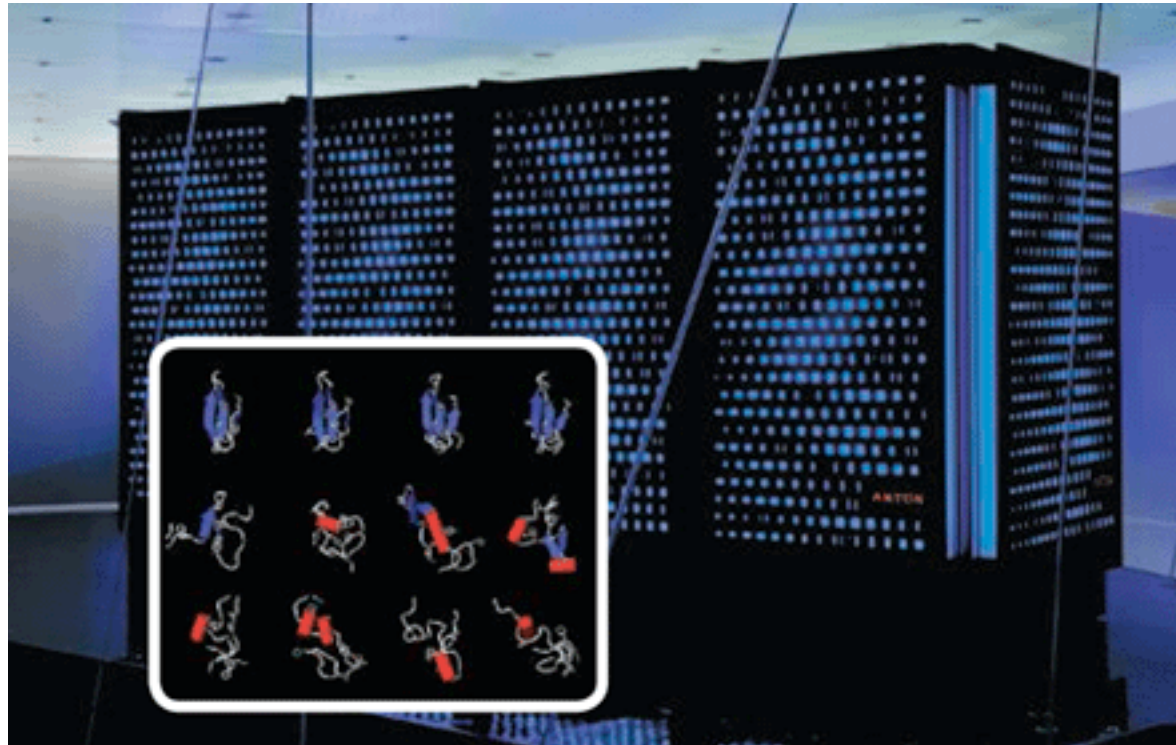
D. E. Shaw Research ("DESRES") is engaged in scientific research in the field of computational biochemistry, including

- The design of novel algorithms and machine architectures for high-speed molecular dynamics (MD) simulations of proteins and other biological macromolecules. In particular, we have designed and constructed a specialized supercomputer called Anton, which executes such simulations orders of magnitude faster than was previously possible, along with a number of software tools and techniques that facilitate their execution and analysis.
- The use of long MD simulations to study the structural changes underlying biological phenomena that occur on time scales far in excess of those previously accessible to computational study, with the ultimate aim of significantly advancing the process of drug development. We have been investigating, for example, the mechanisms of certain cellular receptors, transport proteins, and enzymes relevant to the understanding and potential treatment of cancer, diabetes, and other diseases.

Members of the lab include computational chemists and biologists, computer scientists and applied mathematicians, and computer architects and engineers, all working collaboratively under the direct scientific leadership of its Chief Scientist, [David Shaw](#).

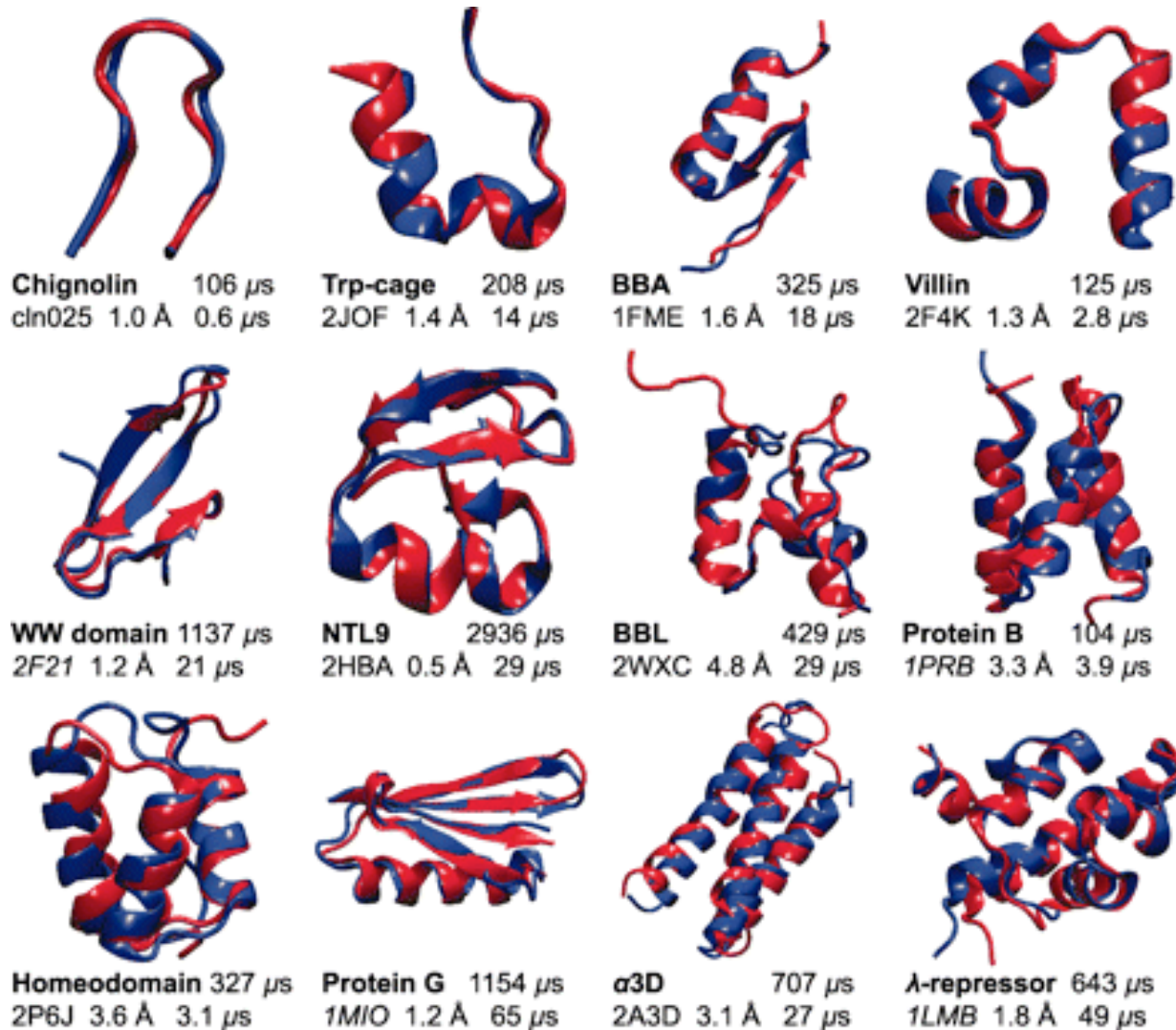
<http://www.deshawresearch.com/>

# David E. Shaw Institute



The Anton Supercomputer – 100 X faster than any other supercomputer for protein folding simulations

# How Well Does Anton Do?



# Summary

- **Structure prediction is still one of the key areas of active research in bioinformatics and computational biology**
- **Significant strides have been made over the past decade through the use of larger databases, machine learning methods and faster computers**
- **Ab initio structure prediction remains an unsolved problem (but getting closer)**