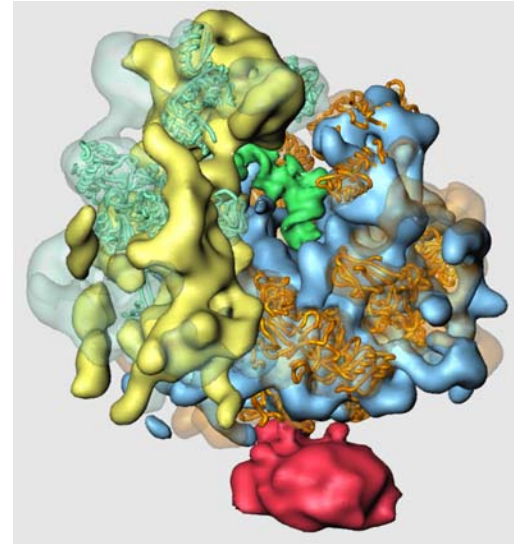


# Protein structure modeling and the Protein Structure Initiative

**András Fiser**



Department of Biochemistry and  
Seaver Foundation Center for Bioinformatics  
Albert Einstein College of Medicine  
New York, USA

# Sequence *versus* Structure

GDCAGDFKIWFYFGRTLLVAGAKDEFGAIDA

RTLAWYAGHLVAGAKDEFGGDFKIWFYFGAI

DFLLVAGAKDEFGKIWFYGGIDAWRTAGDC

HLVAGARTLAFGAIDWYAKDEFGGGDFKIW

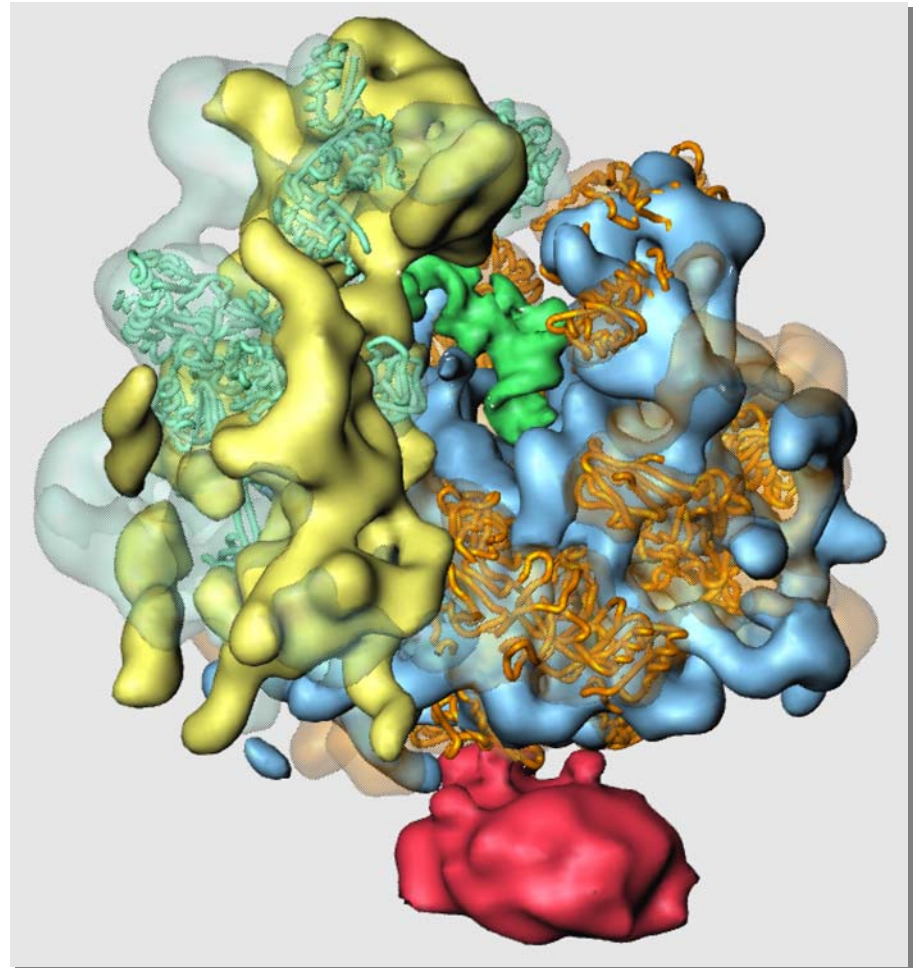
ARTHLVAGFGGGAIDWYFKIWFYAKLAFGDE

GCTAGCTTAAGGCCTTCATGATCTTCTGAG

AGGGCTCCTTCATGATAGCTTAAGGCTTAA

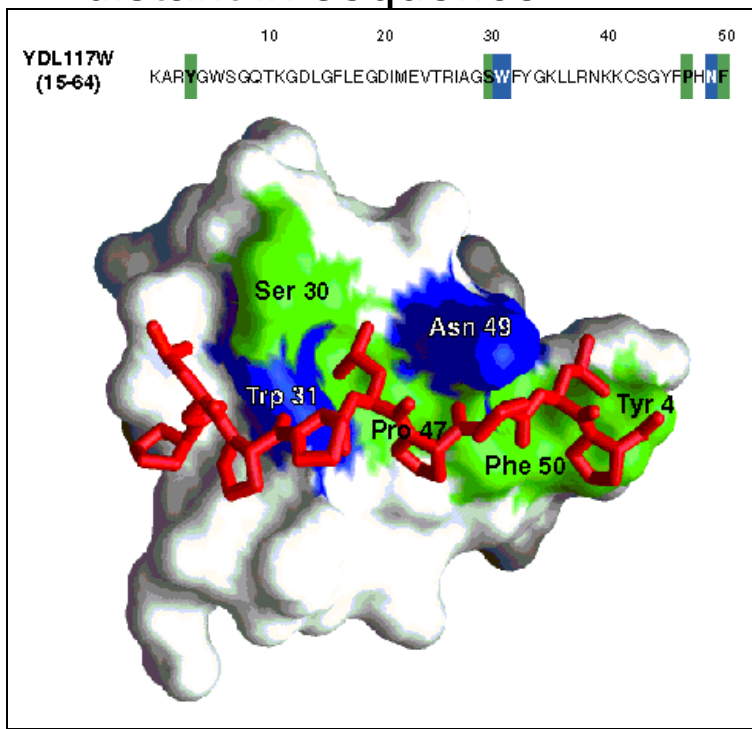
AGGCCTTCATGGGGTTAACATATCTTCTGA

CCTTCATGCTAGCTTAAGGGATCTTAACCG



# Why is it useful to know the structure of a protein not only its sequence?

- The biochemical function of a protein is defined by its interactions with other molecules.
- The 3D structure is more informative than sequence because interactions are determined by residues that are close in space but are frequently distant in sequence.



Evolution tends to conserve function and function depends more directly on structure than on sequence, **structure is more conserved in evolution** than sequence.

Patterns in space are frequently more recognizable than patterns in sequence.

# Function via Structure

Sequence

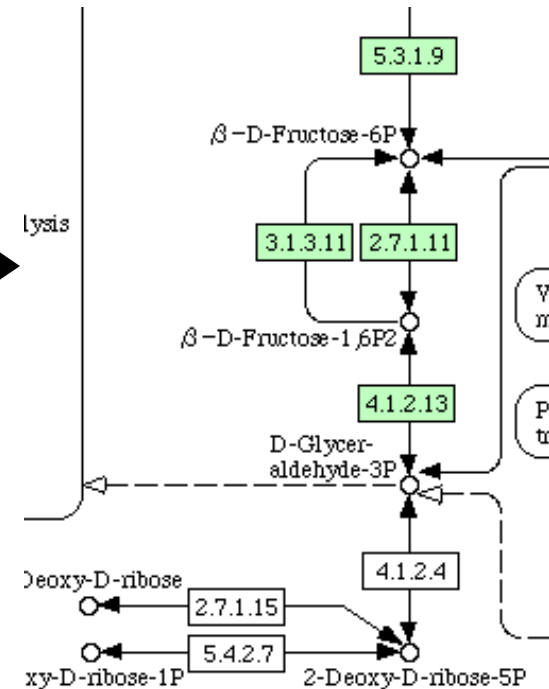
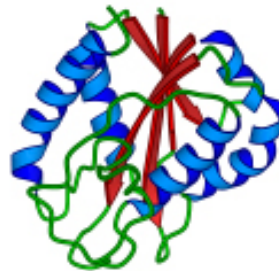


Structure



Function

GFCHIKAYTRLIM...

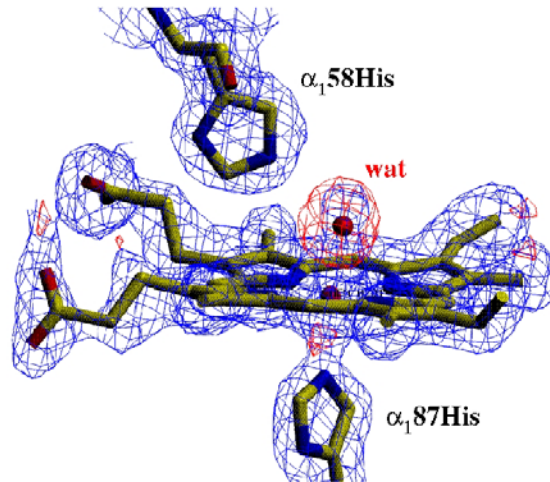
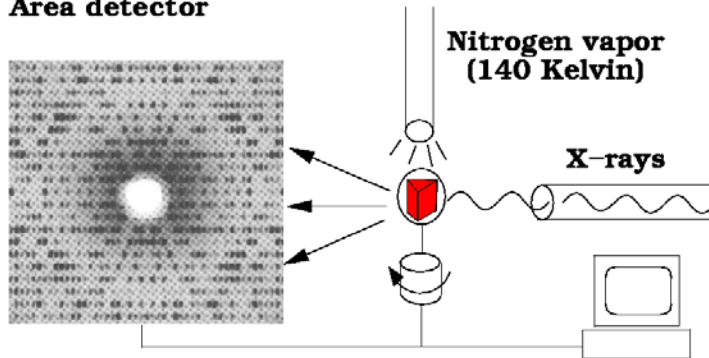




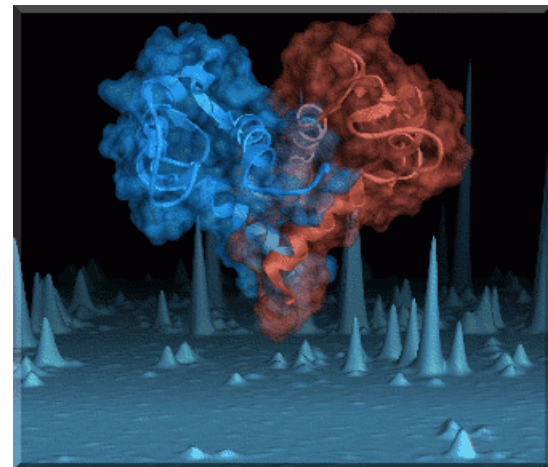
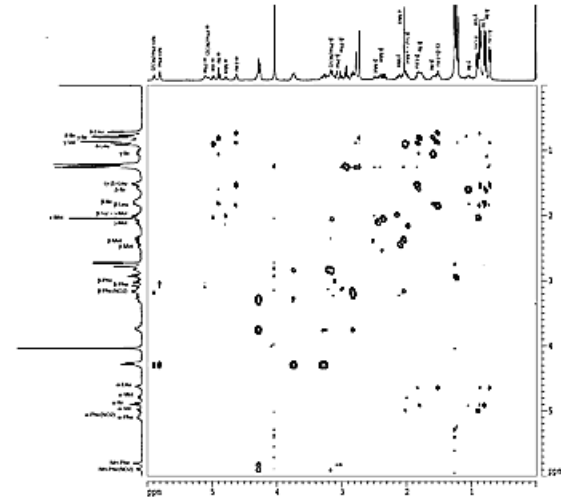
# Major experimental tools

## X-ray crystallography

Area detector



## Nuclear Magnetic Resonance



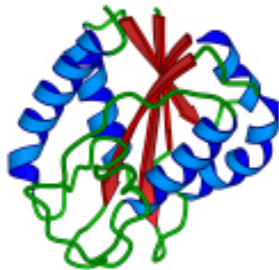
# Why Protein Structure Prediction?

	Y 2004	Y 2006
Sequences	1,900,000	millions
Structures	26,000	40,000

**We know the experimental 3D structure for about 1% of the protein sequences**

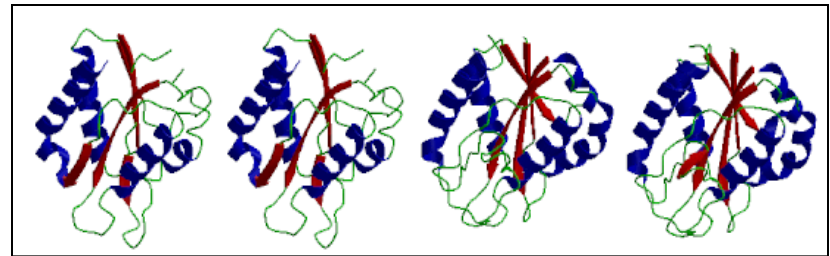
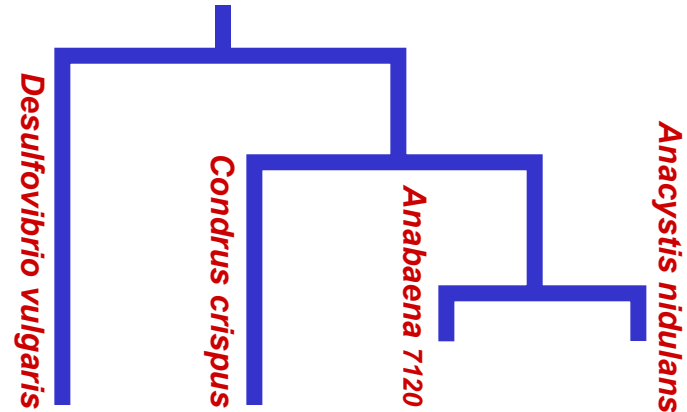
# Principles of Protein Structure

GFCHIKAYTRLIMVG...



folding

*Ab initio* prediction



evolution

Fold Recognition  
Comparative Modeling

# Protein structure modeling

## *Ab initio* prediction

**Applicable** to **any** sequence

**Not very accurate**, and attempted for proteins of <100 residues

*Accuracy and applicability are limited by our understanding of the protein folding problem*

## Comparative Modeling

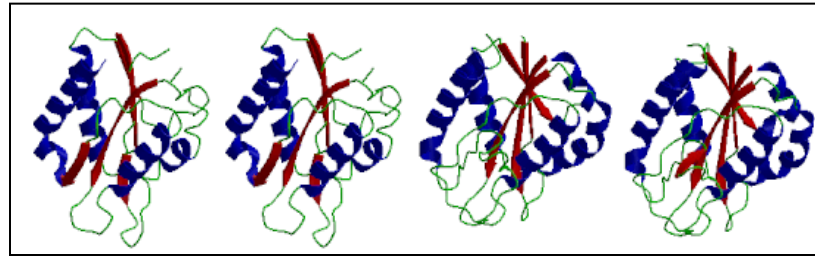
**Applicable** to those sequences **only** that share recognizable similarity to a template structure

**Fairly accurate**, typically comparable to a low resolution X-ray experiment.  
Not limited by size

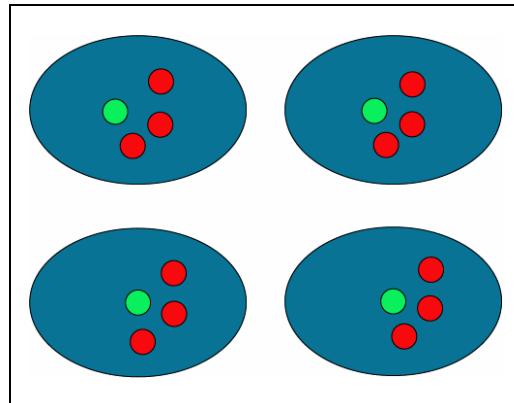
*Accuracy and applicability are rather limited by the number of known folds*

# What makes comparative modeling possible

I A small difference in the sequence makes a small difference in the structure

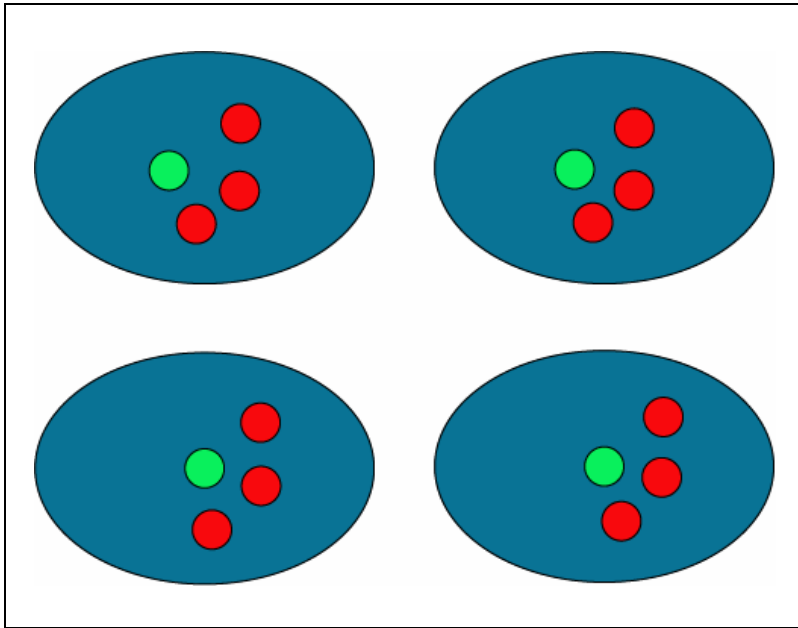


II Protein structures are clustered into fold families



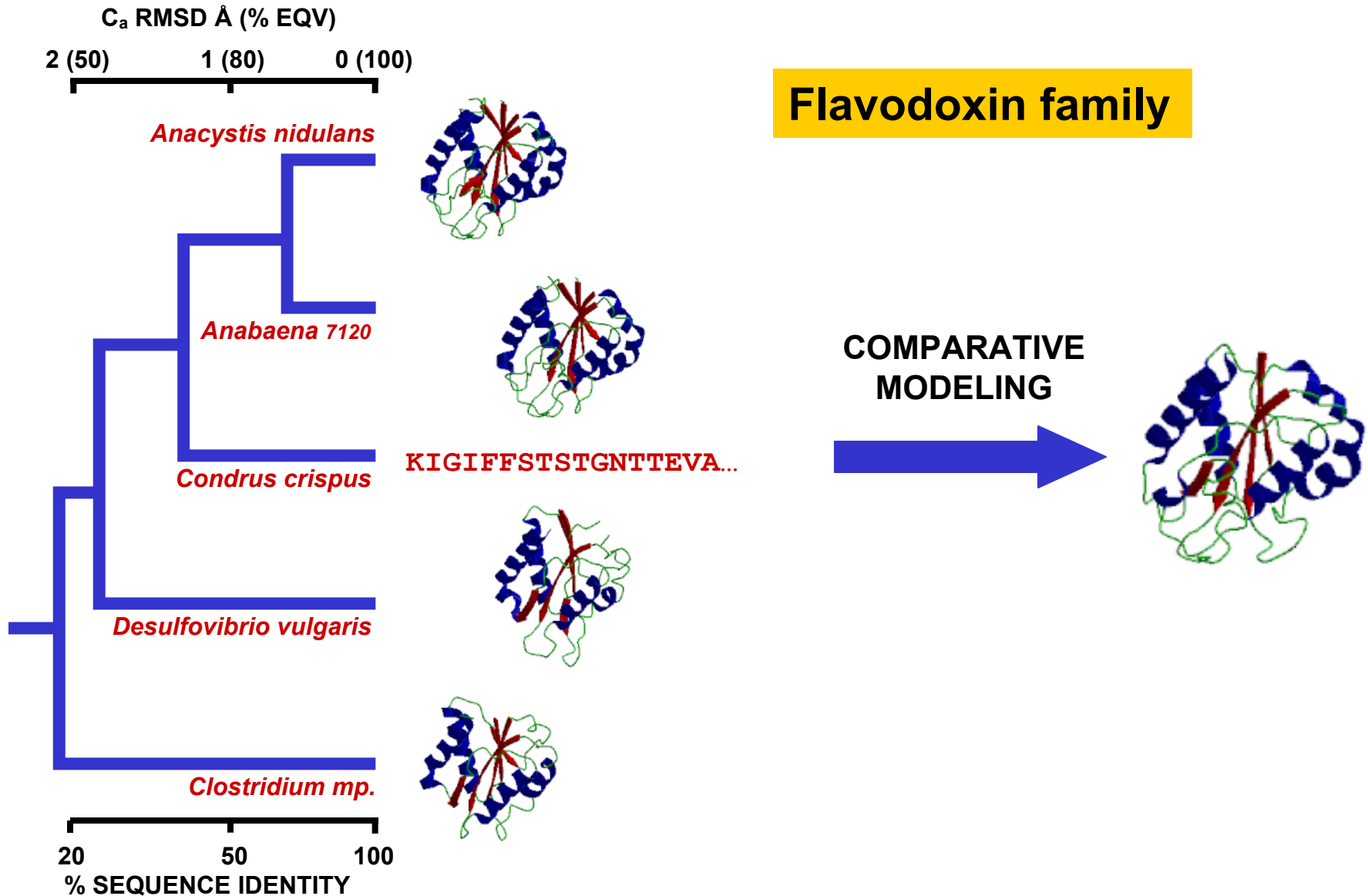
# Structural Genomics

Characterize most protein sequences (**red**) based on related known structures (**green**).



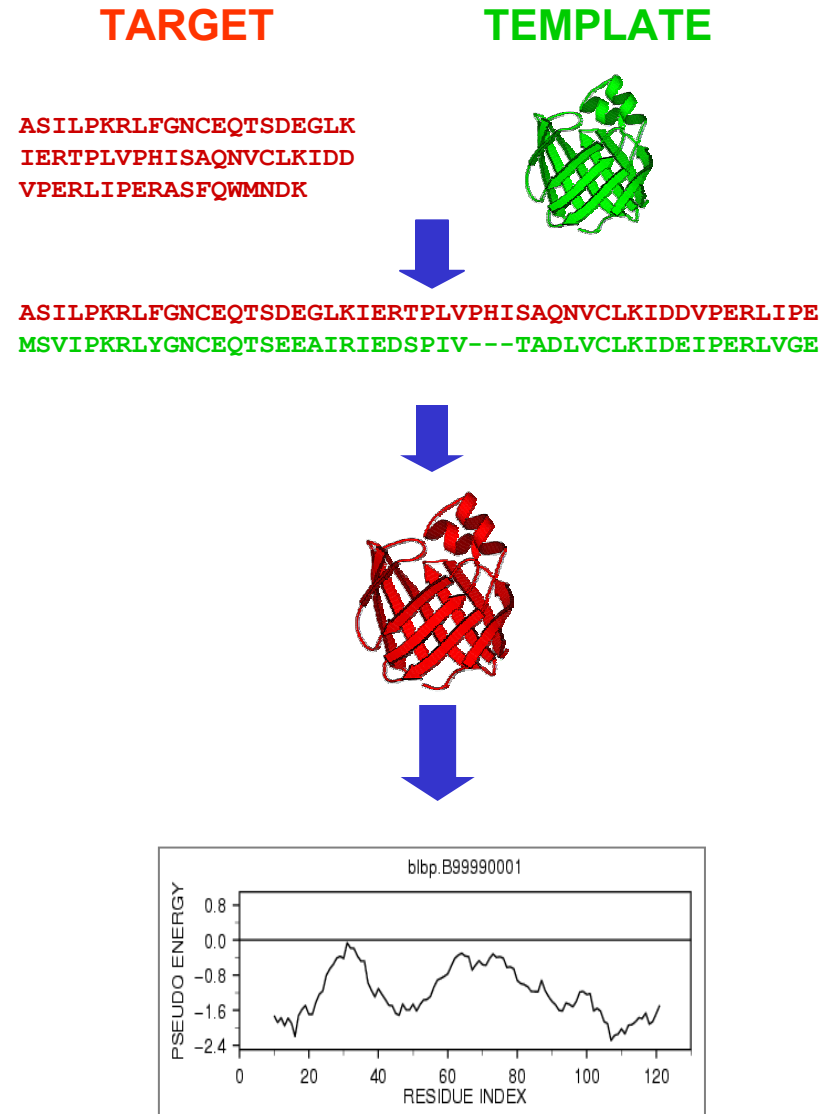
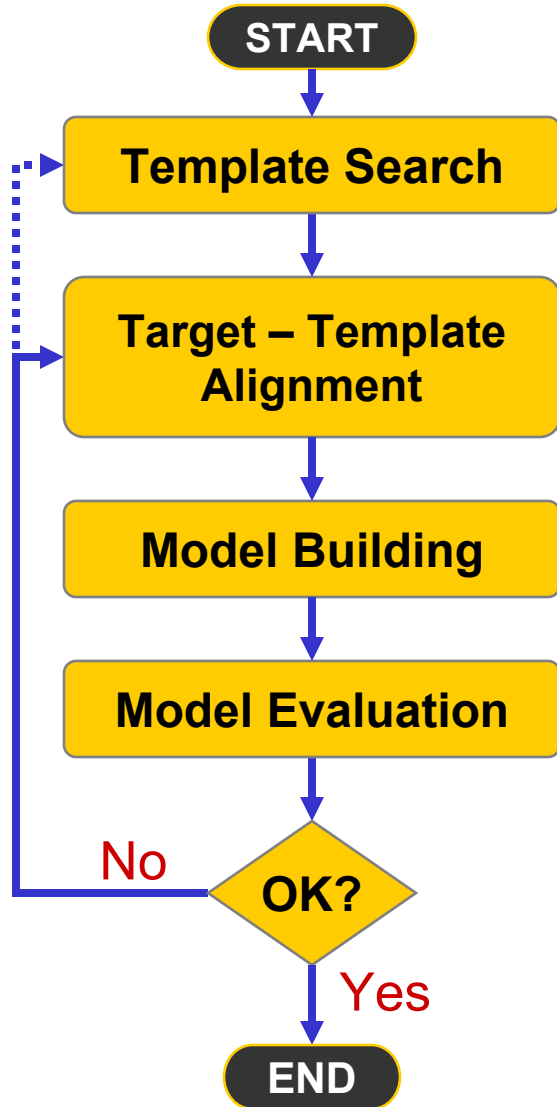
The number of “families” is much smaller than the number of proteins

# Comparative Protein Structure Modeling





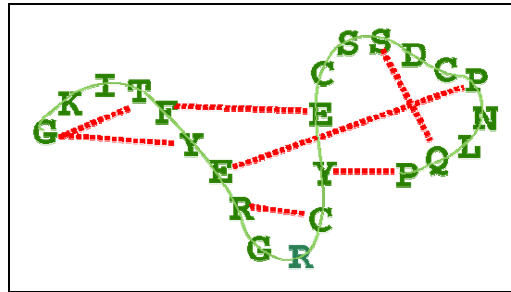
# Steps in Comparative Protein Structure Modeling



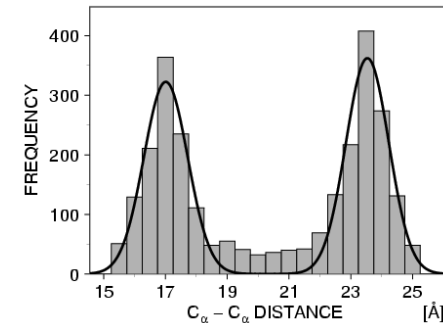
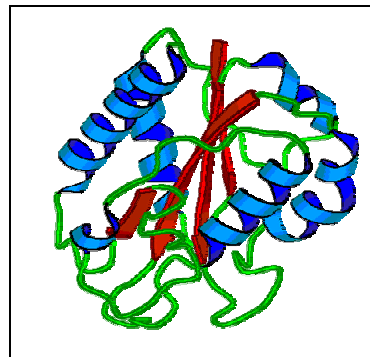
# Comparative Modeling by Satisfaction of Spatial Restraints (MODELLER)

3D **GKITYERGFQGH**CY**ESDC**-NL**Q**?...  
 SEQ GKITYERG---RCY**ESDC**PNL**Q**?...

## 1. Extract spatial restraints



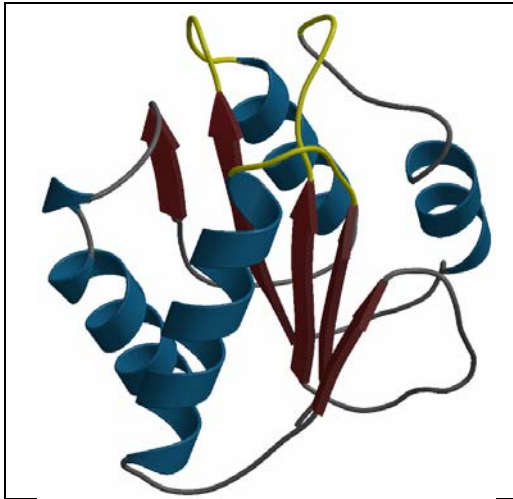
## 2. Satisfy spatial restraints



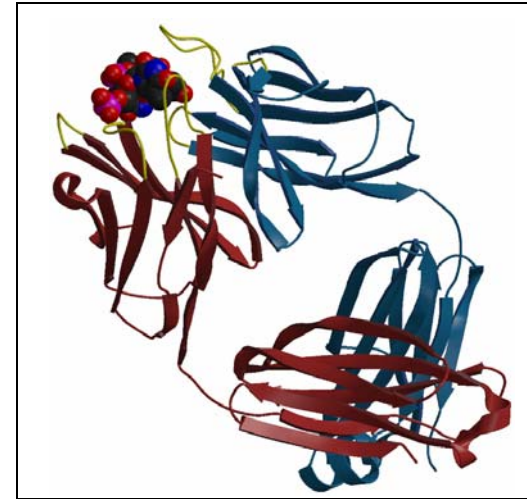
$$F(R) = P \prod_i p_i(f_i/l)$$

# Loop Modeling in Protein Structures

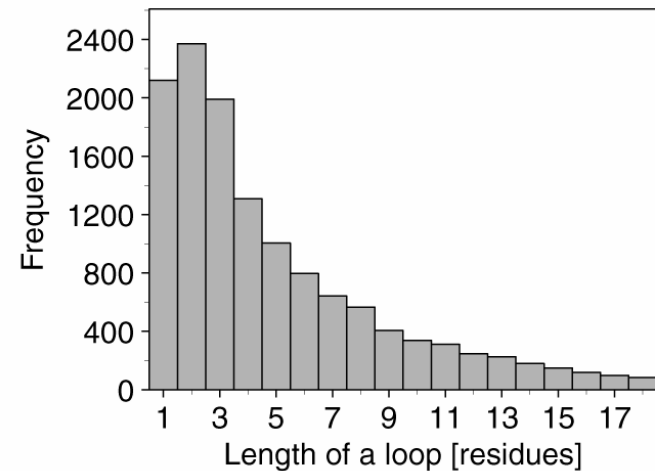
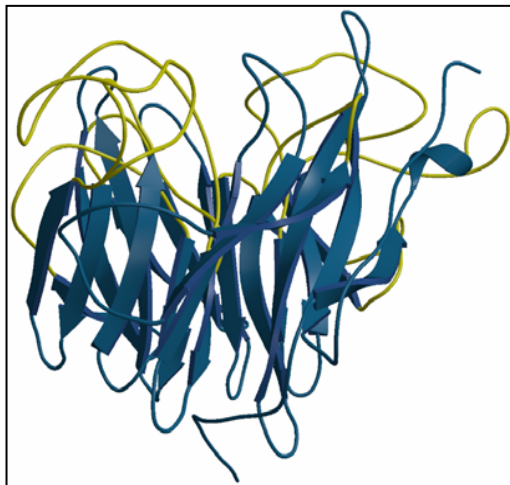
**$\alpha$ + $\beta$  barrel: flavodoxin**



**IG fold: immunoglobulin**



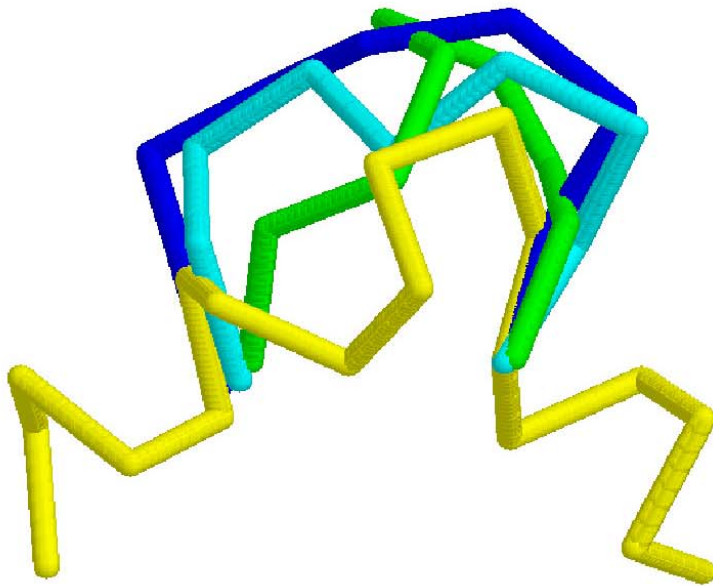
**antiparallel  $\beta$ -barrel**



# Loop modeling strategies

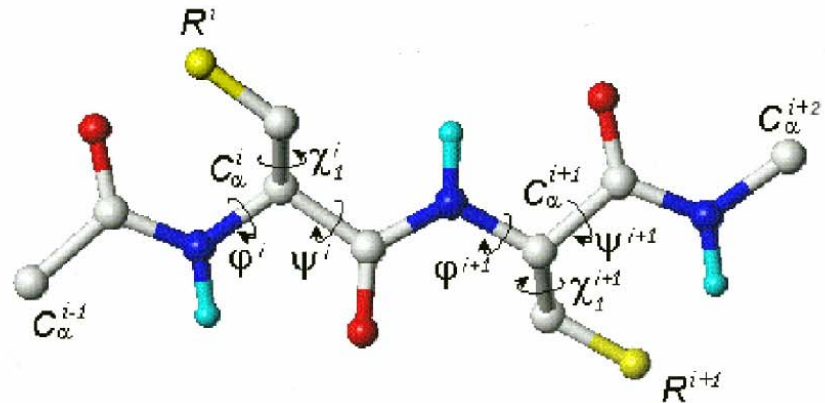
## Database search

“Comparative”



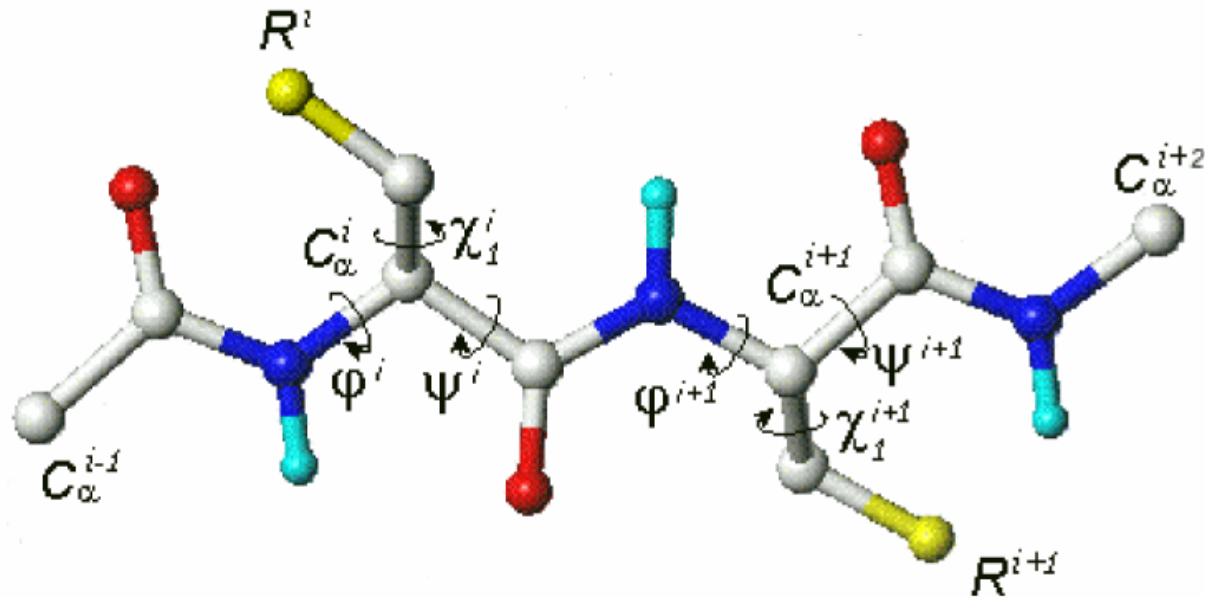
## Conformational search

“*ab initio*”



- even in DB search, the different conformations must be ranked
- database is complete only up to 4-6 residues
- loops longer than 4 residues need extensive optimization
- DB method is efficient for specific families (eg. Canonical loops in Ig's,  $\beta$ -hairpins etc)

# Loop Modeling by Conformational Search



1. Protein representation.
2. Energy (scoring) function.
3. Optimization algorithm.

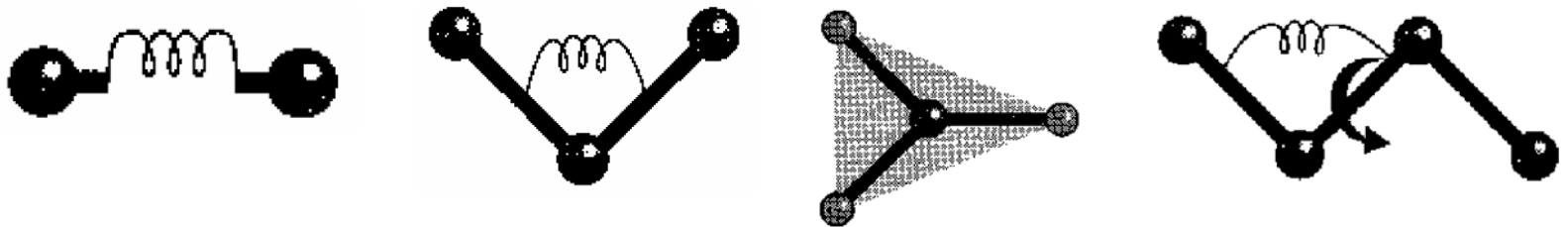
# Energy Function for Loop Modeling

The energy function is a sum of many terms:

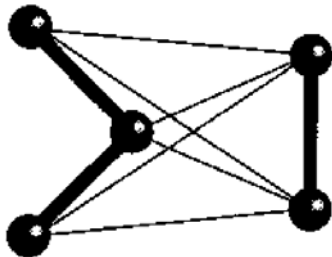
## 1) Statistical preferences for dihedral angles:



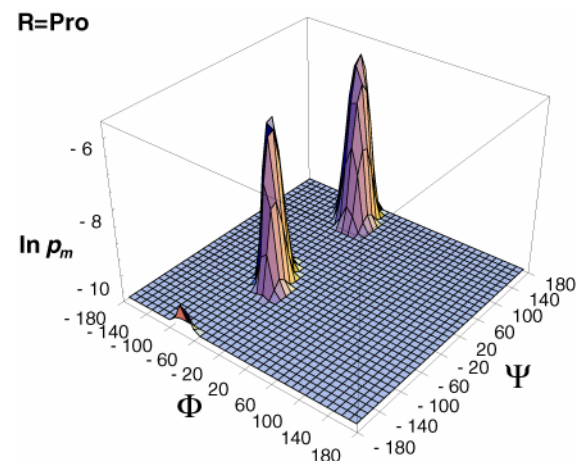
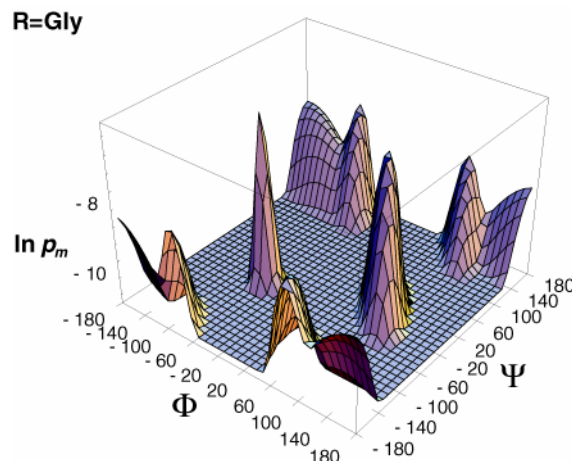
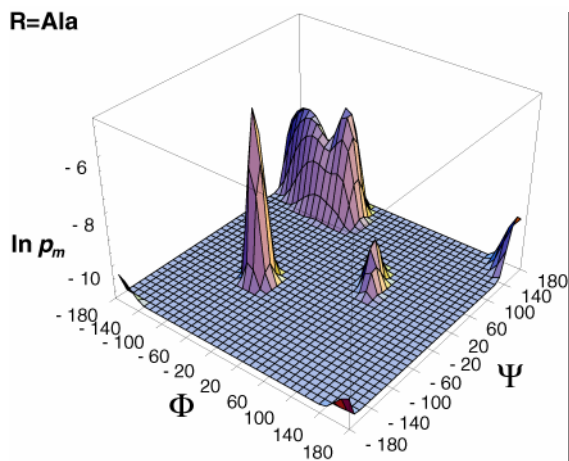
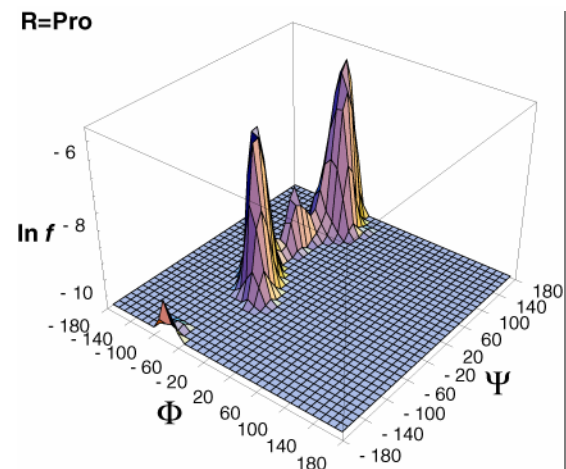
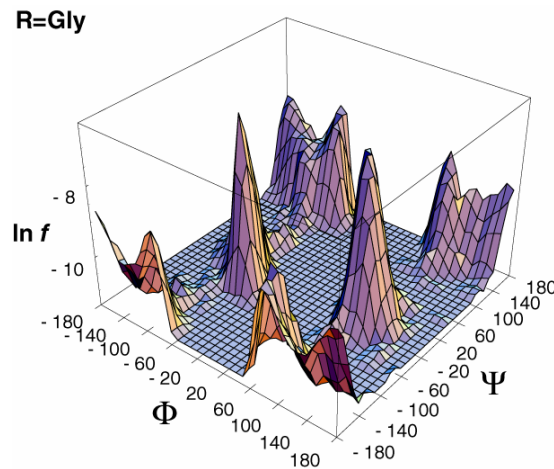
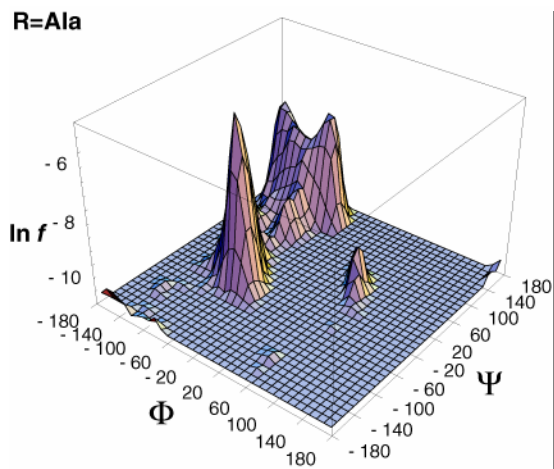
## 2) Restraints from the CHARMM-22 force field:



## 3) Statistical potential for non-bonded contacts:

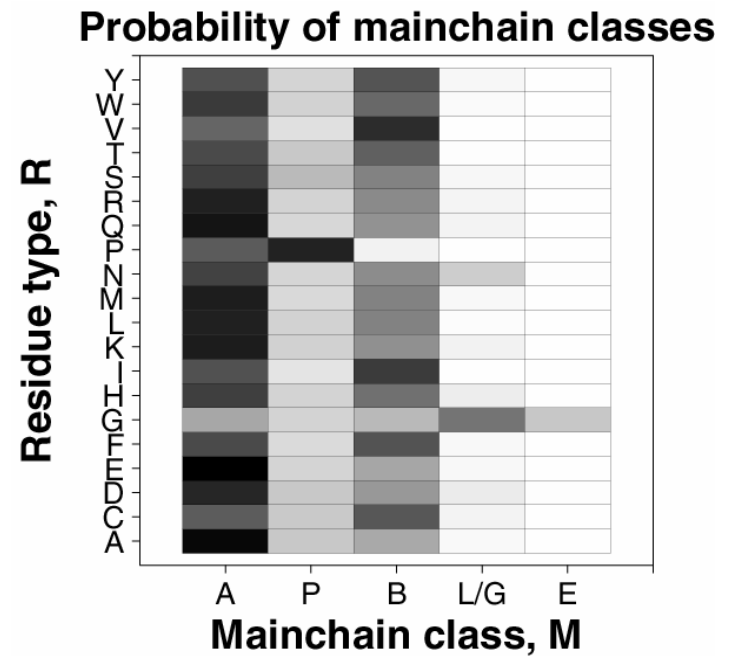
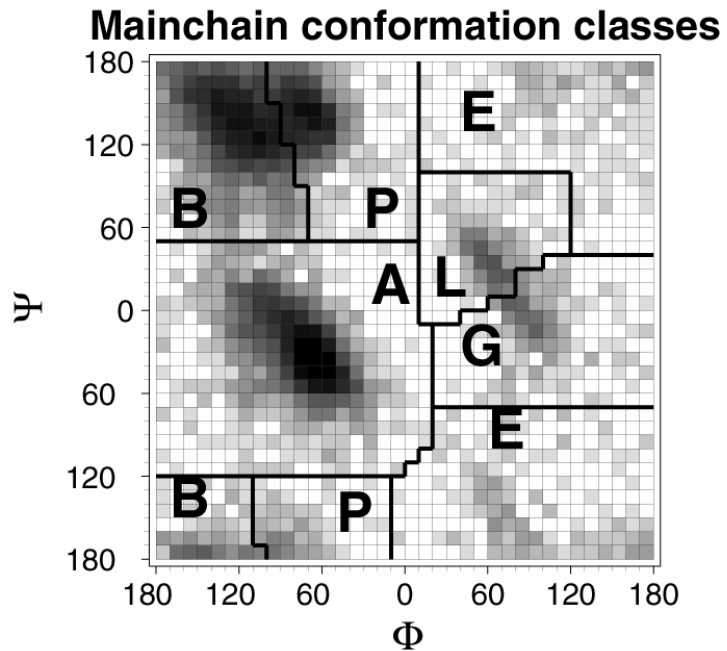


# Mainchain Terms for Loop Modeling



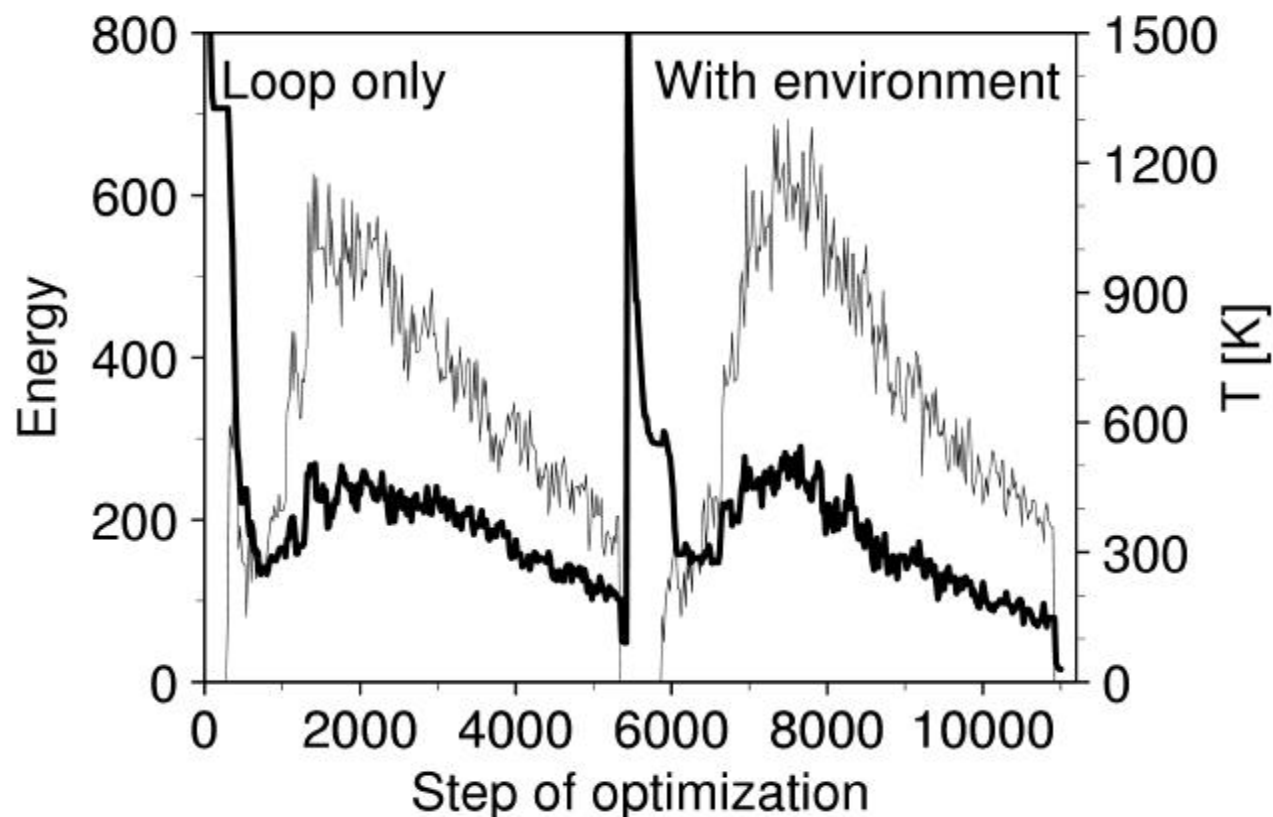


# Mainchain terms for loop modeling



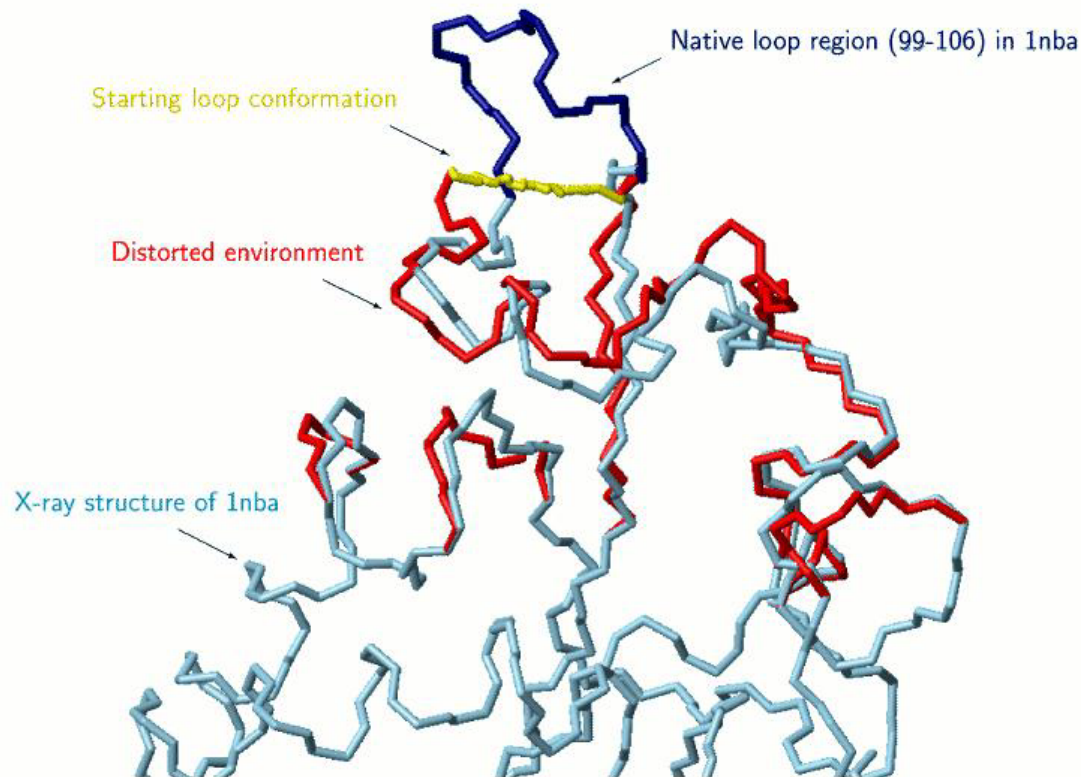
# Optimization of Objective Function

Many different combinations of objective function terms were explored. The objective function was optimized with a combination of conjugate gradients method and molecular dynamics simulation with simulated annealing in a two step process: without and with the environment.

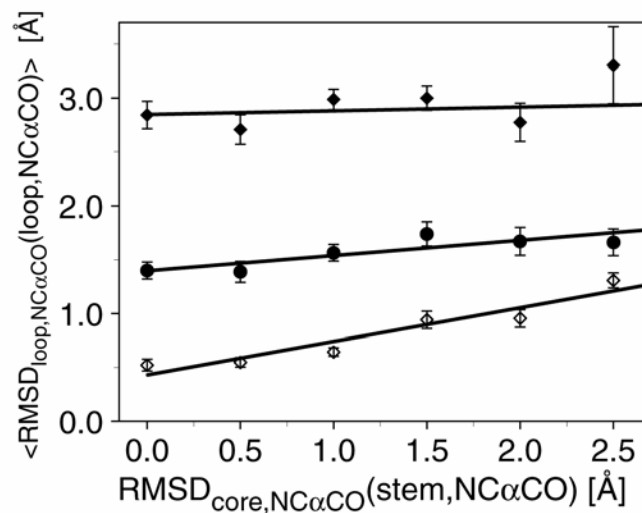
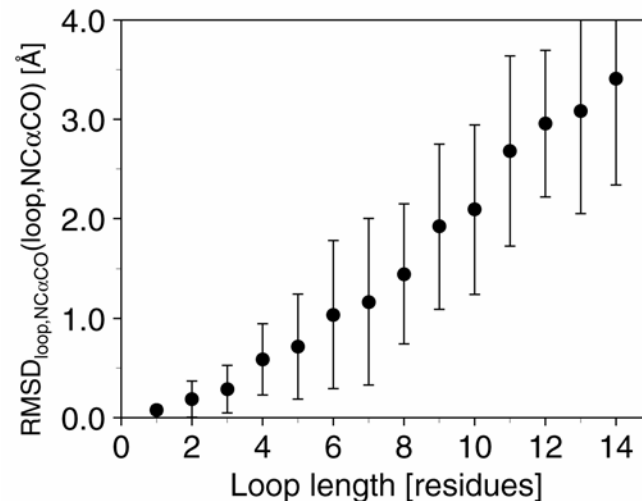
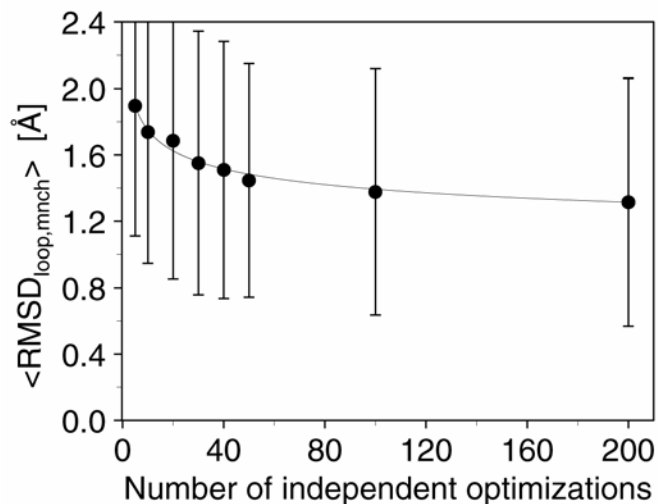


# Optimization of Objective Function

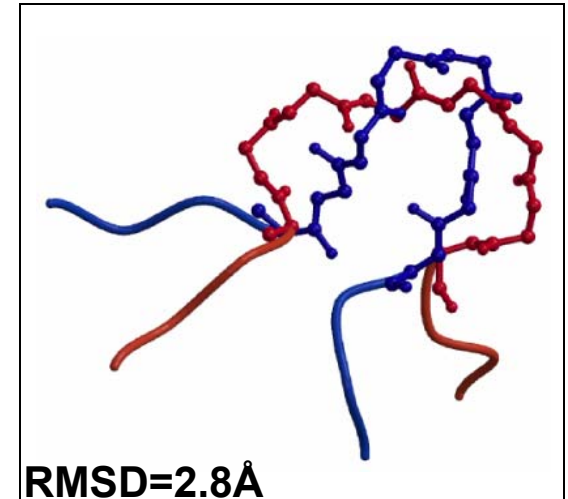
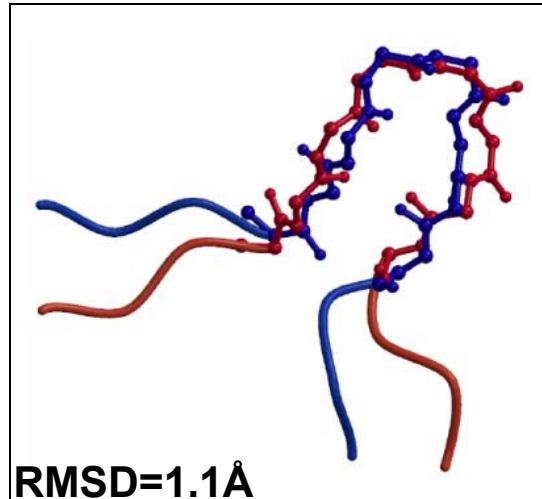
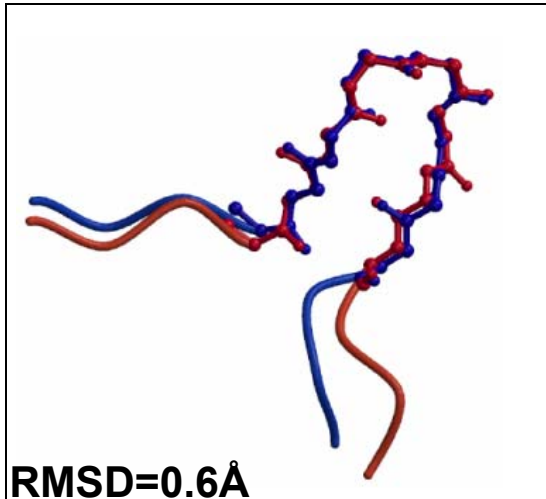
- Test set: 40 randomly selected loops of known structures, for each length from 1 to 14 residues.
- Starting conformation: Loop atoms were spaced evenly on a line spanning the two anchor regions, then randomized by  $\pm 5$  Å.
- To simulate real comparative modeling situations, performance of the loop modeling problem was determined by predicting loops in only approximately correct environment.



# Accuracy of loop models



# Accuracy of Loop Modeling



HIGH ACCURACY (<1Å)

50% (30%) of 8-residue loops

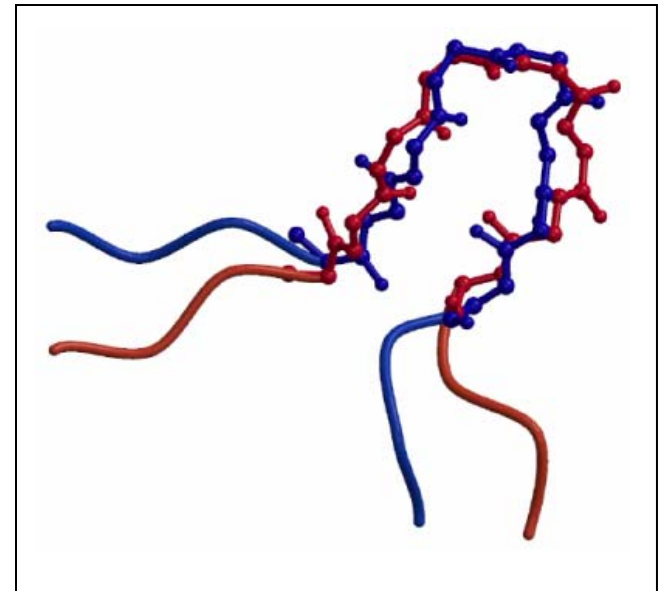
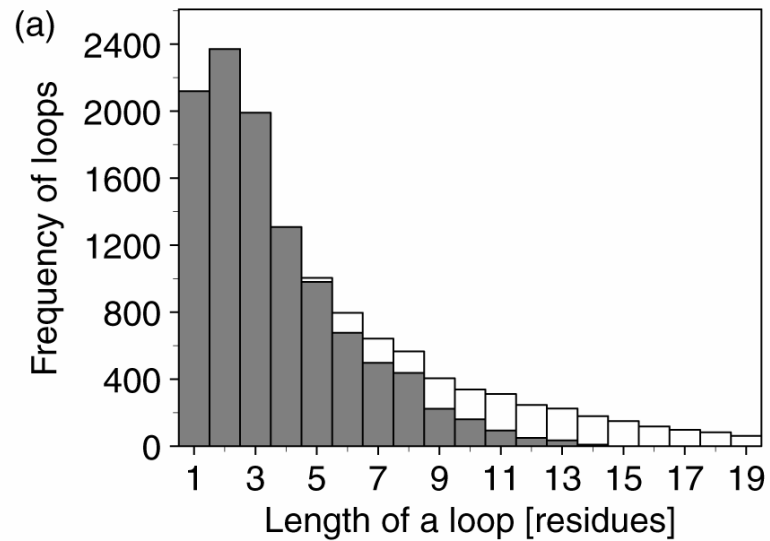
MEDIUM ACCURACY (<2Å)

40% (48%) of 8-residue loops

LOW ACCURACY (>2Å)

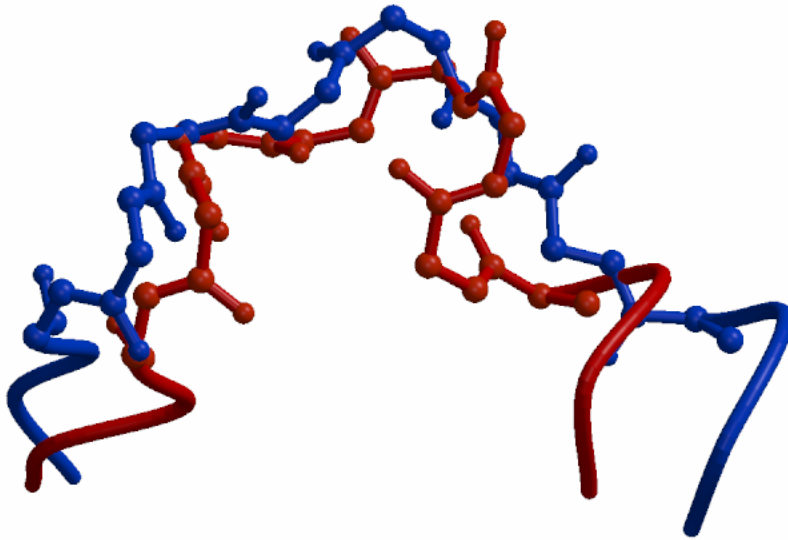
10% (22%) of 8-residue loops

# Fraction of Loops Modeled With at Least Medium Accuracy



# Problems in Practical Loop Modeling

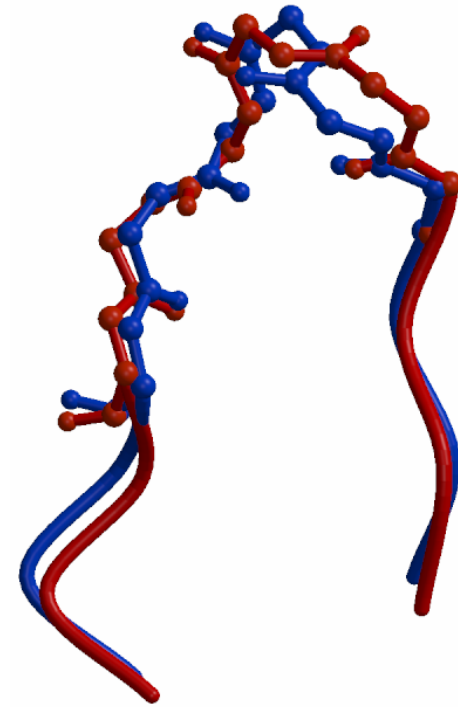
1. Decide which regions to model as loops.
2. Correct alignment of anchor regions & environment.
3. Modeling of a loop.



T0076: 46-53

$\text{RMSD}_{\text{mnch}}$  loop = 1.37 Å

$\text{RMSD}_{\text{mnch}}$  anchors = 1.52 Å



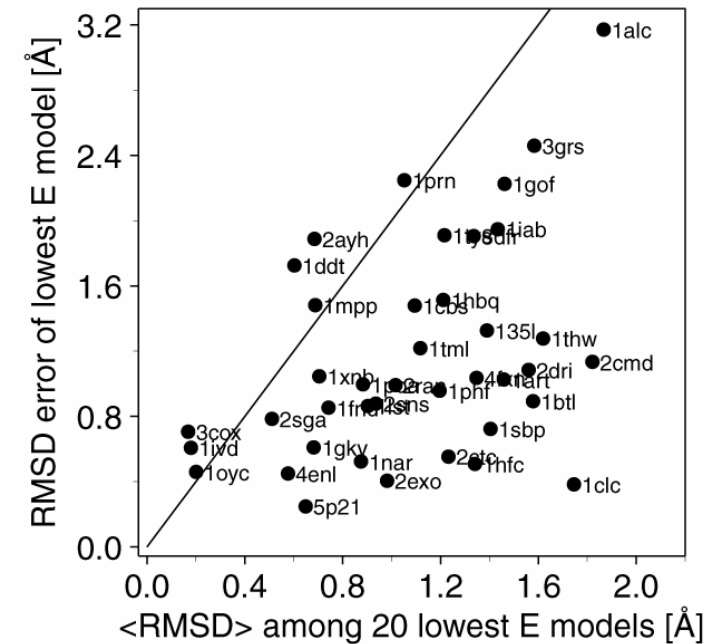
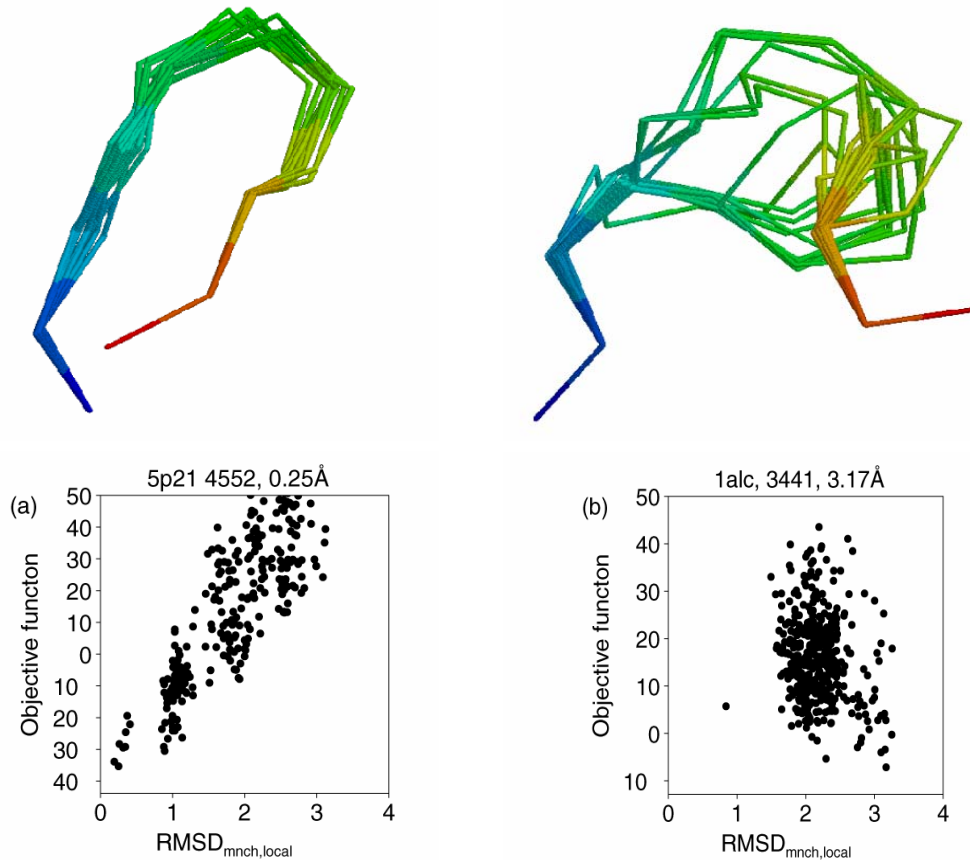
T0058: 80-85

$\text{RMSD}_{\text{mnch}}$  loop = 1.09 Å

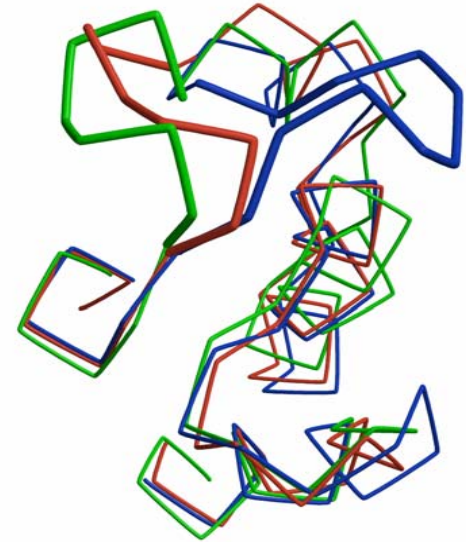
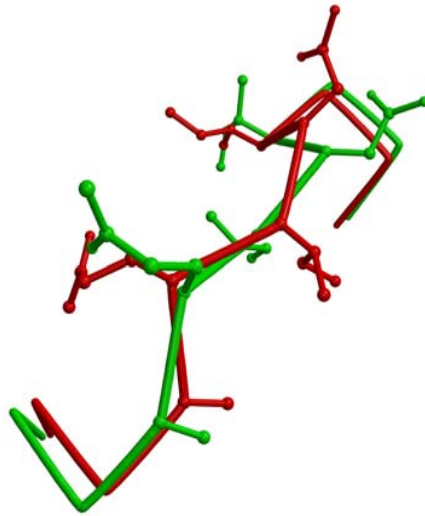
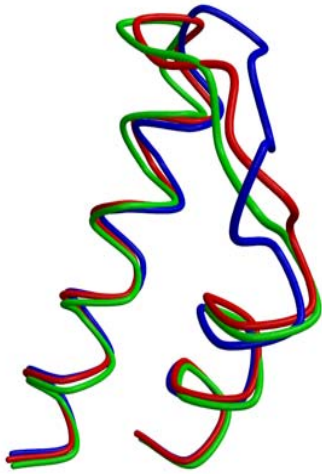
$\text{RMSD}_{\text{mnch}}$  anchors = 0.29 Å



# Assessing Accuracy of Loop Models



# Examples from CASP 3/4



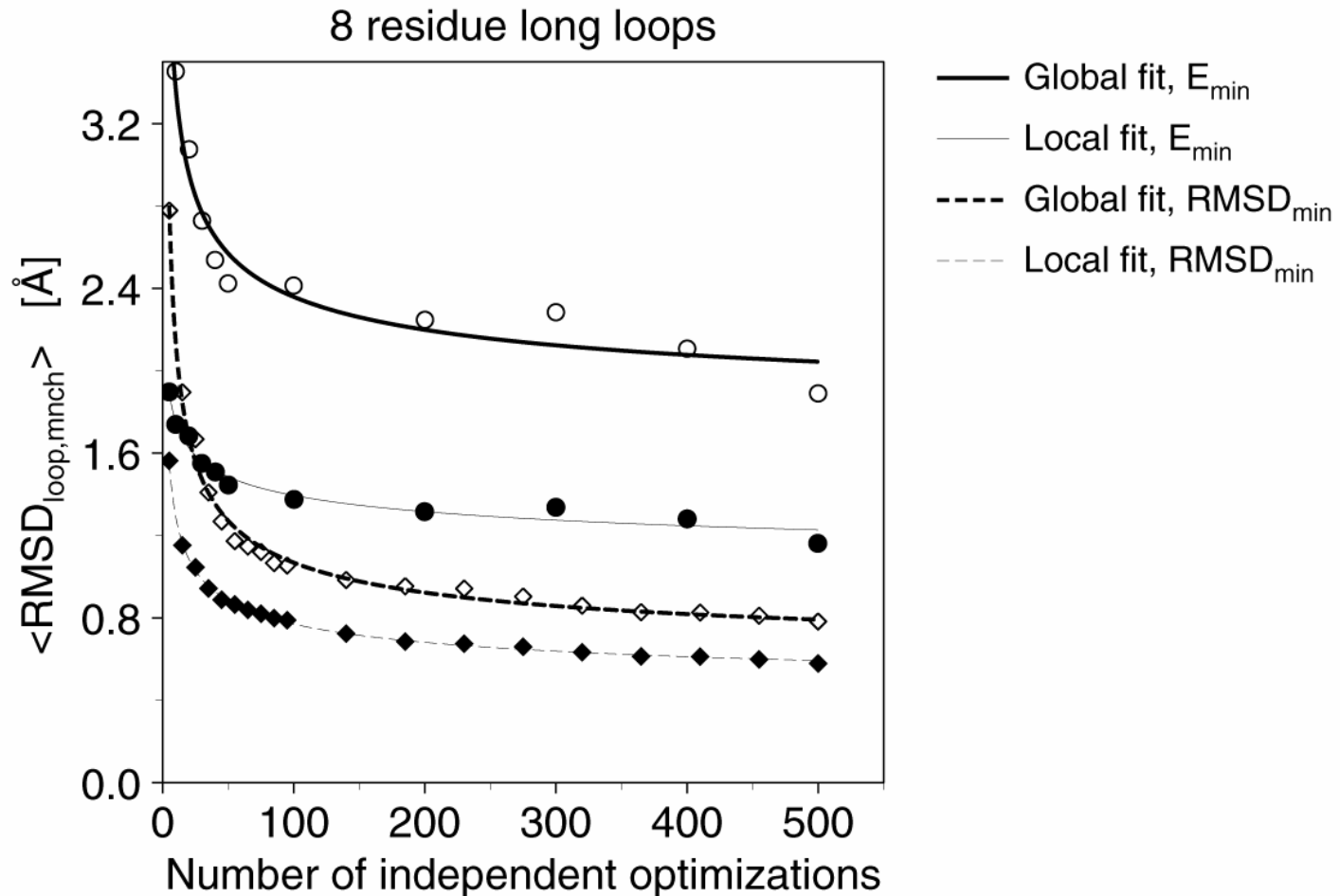
## RMSD

	3-3 anchor residues	global		local	
		MNCH	ALL	MNCH	ALL
<b>X-ray/model</b>	0.76	1.28	2.48	1.05	1.90
<b>X-ray/template</b>	2.38	2.69	4.22	1.94	3.53
<b>model/template</b>	2.10	2.75	4.29	1.64	2.92

*S. pombe* contractile ring protein Cdc4p.  
33 % seq id. 8 residue long loop,  
RMSD<sub>global</sub>: 3.64, RMSD<sub>local</sub>: 1.36Å

# **Adding solvent effect to loop modeling**

# Accuracy of loop models as a function of amount of optimization



# Refining with CHARMM/GB

sample loop models



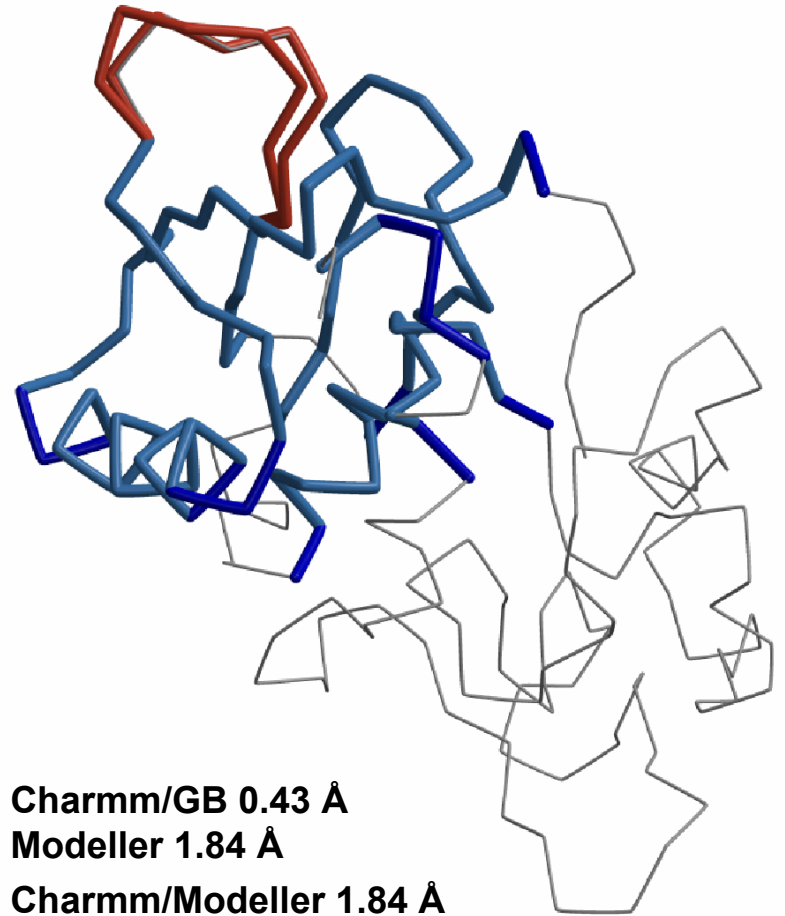
Select loop models



Select loop and its environment



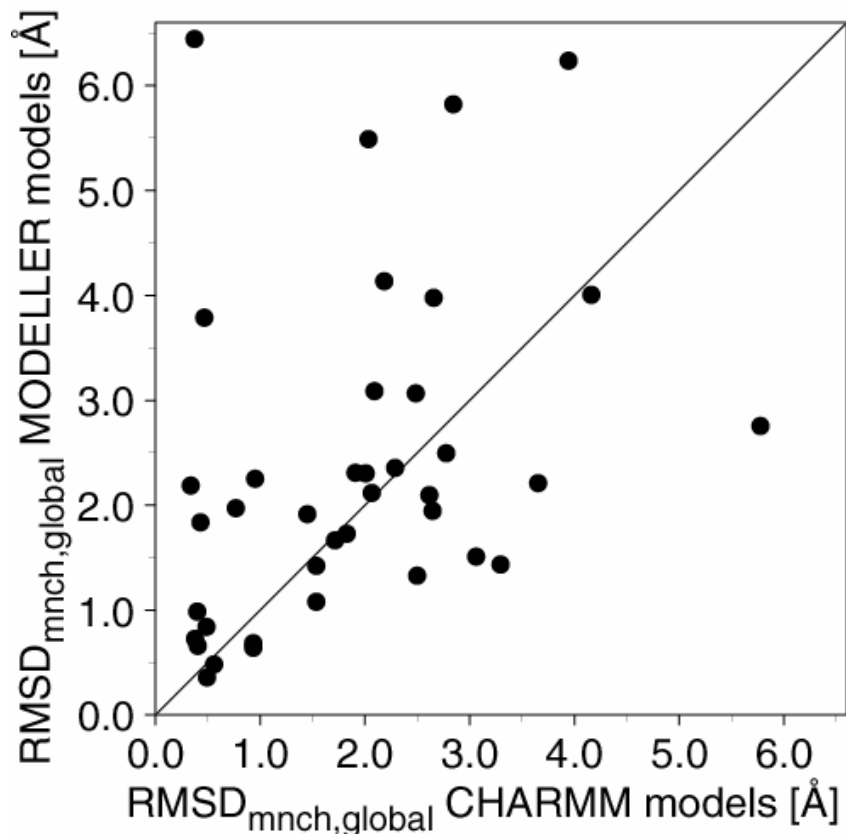
Minimize loop and its environment  
with CHARMM/GB



50 steps of the steepest descent relaxation,  
followed by 2000 steps of ABNR minimization  
or until convergence ( $D < 10^{-4}$  kcal/mol).

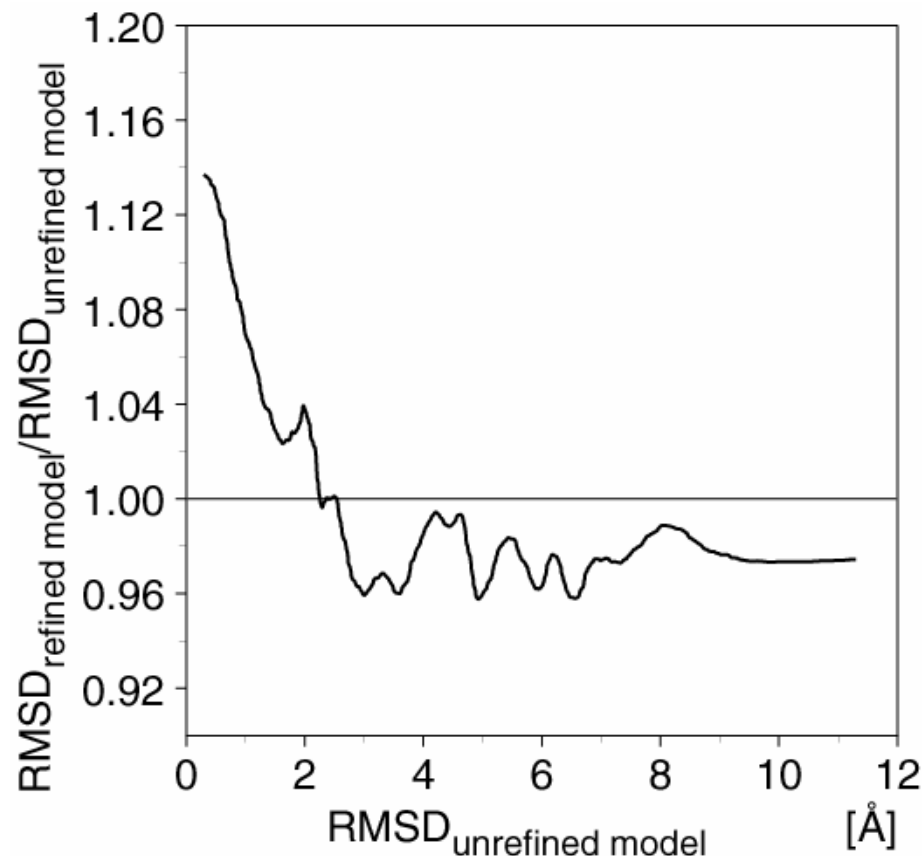
# Refining models with CHARMM/GB

## Improving ranking



$\langle \text{rmsd}_{\text{GLOBAL/LOCAL}} \rangle$   
 CHARMM/GB 1.87 / 1.07 Å  
 MODELLER 2.36 / 1.29 Å

## Improving model quality

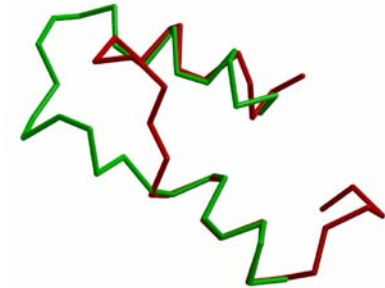
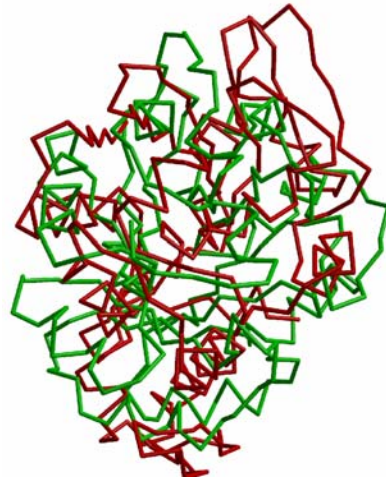
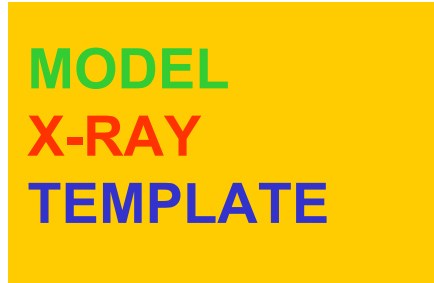


55%-63%  
 45%-37%

# Typical Errors in Comparative Models

Incorrect template

Misalignment



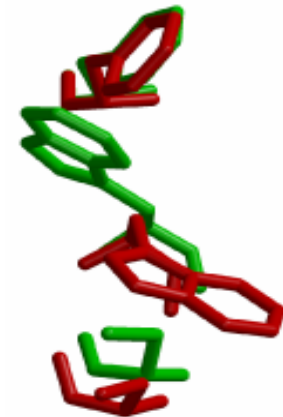
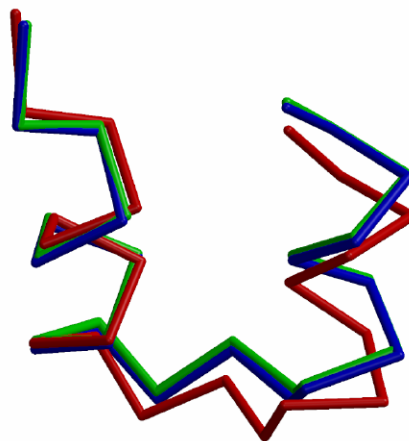
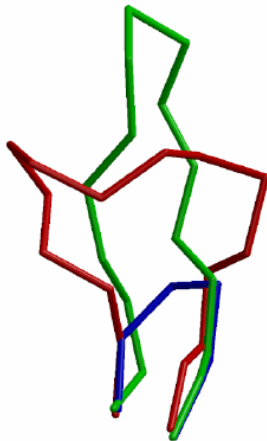
```
EDN  ---FPFQFTWAGMFETQHIIMTSQCOCTHAKGVINNYQRECKWNTPELLTTPAVVWVCCNENATCFEN
7RSA KETAAJFPERQHMDSSTAASSSNYCNQMKSRNLTKDECKPVNTPVHESLADVQAVCSQENIVAC-KN
      aaaaaaaaaa          bbbbbb aaaaaaaaaa

EDN  FTRENCHSGSQVPLIHCLNLTTPSPQISNCRYAQTFAMFVIYVACDNRDQREDPEQVFPVPHLDRII
7RSA -GQTNCYQSYSTMSTDCRETGSS--FYFNCAYETTDANKHIIIVACEGN-----FVYPVPHFDASV
      bbbb  bbbbbbbb  aaaaaaaaaaaaaaaaaa  bbbbbbbb
```

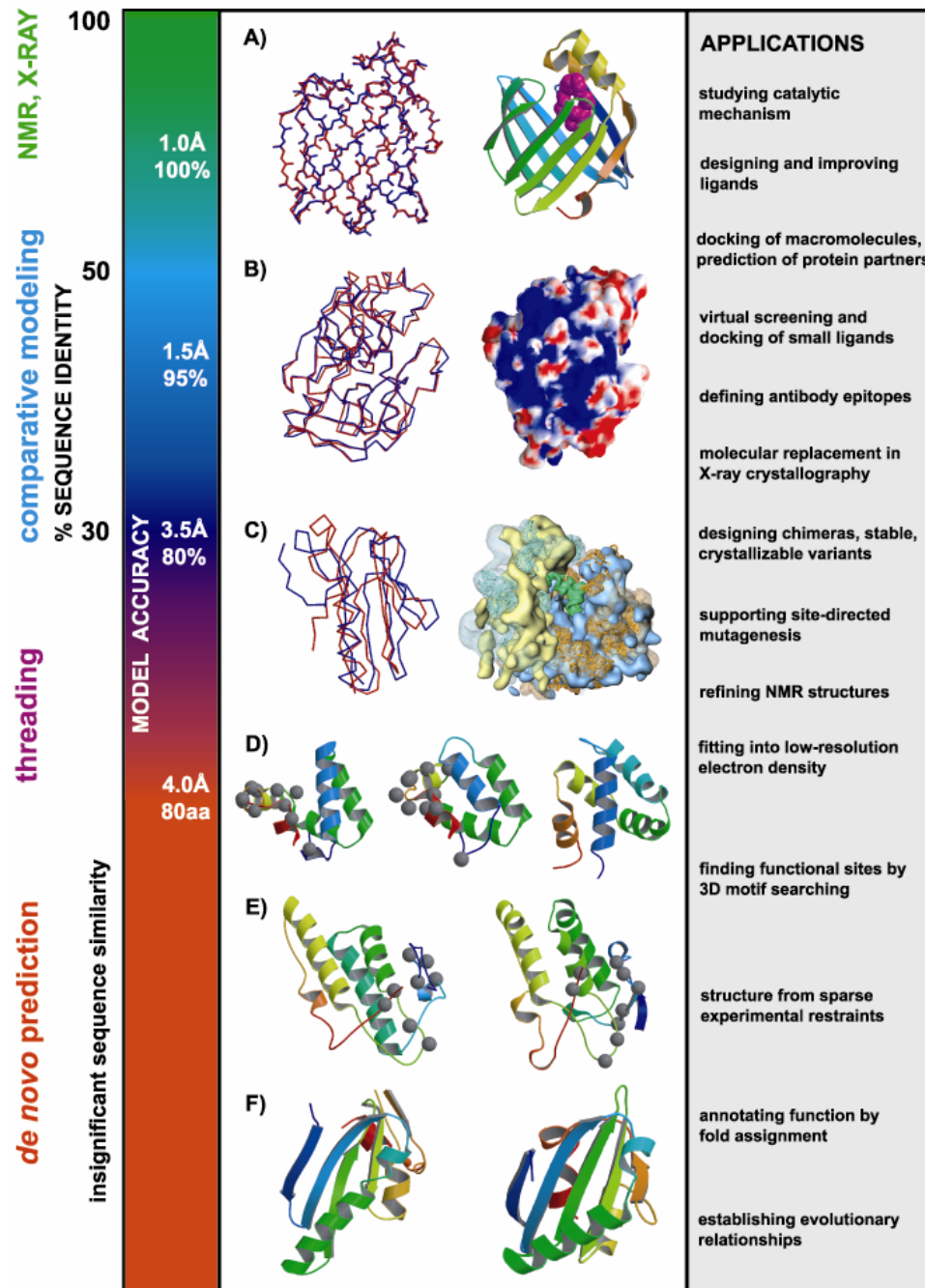
Region without a  
template

Distortion in correctly  
aligned regions

Side chain packing

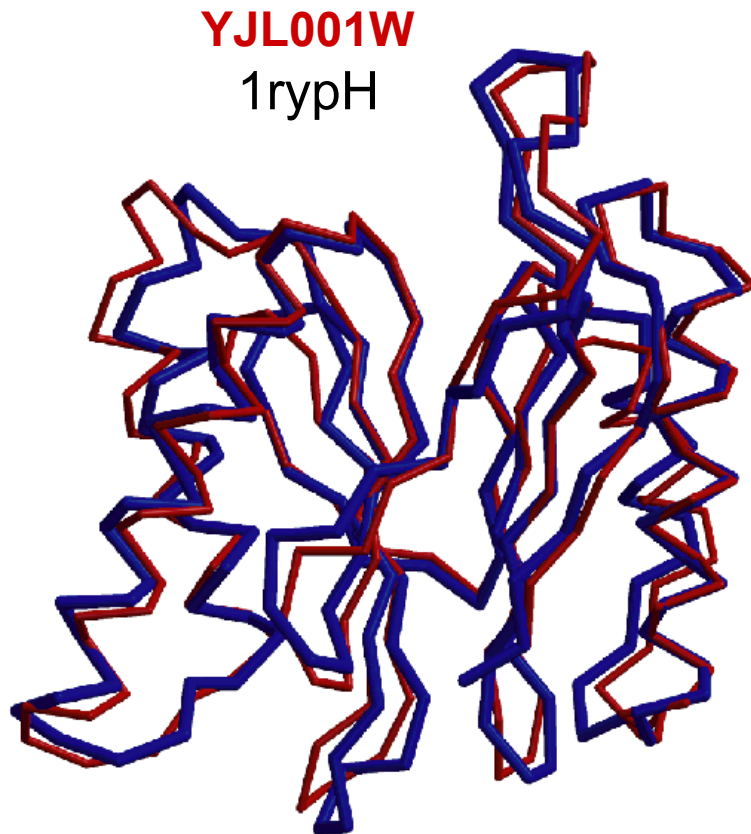




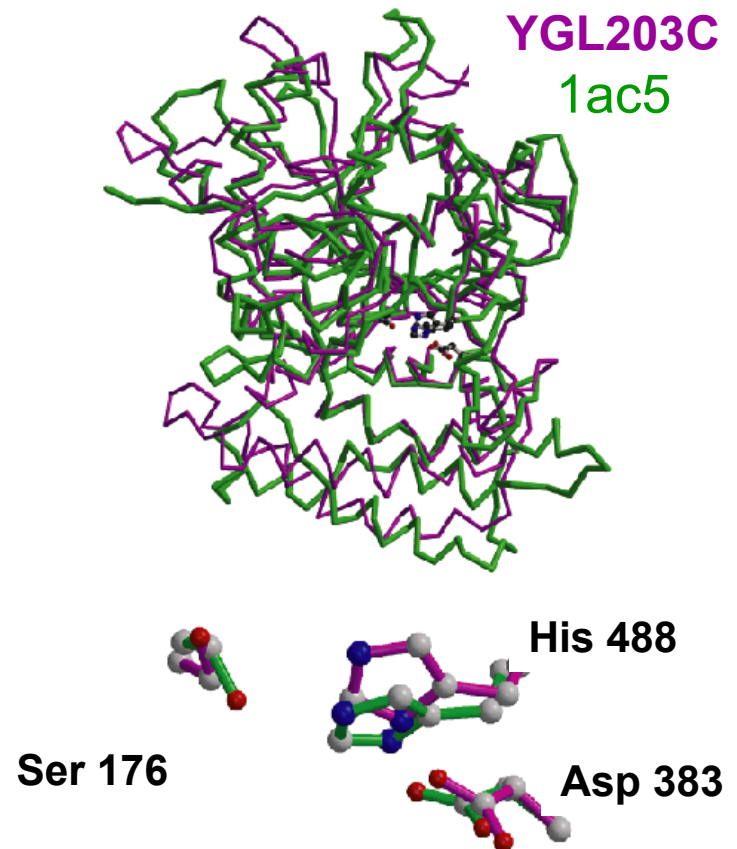


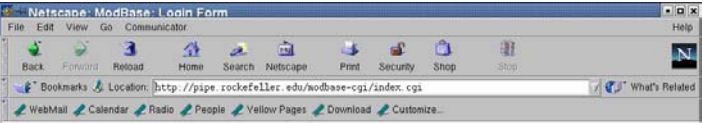
# Some Models Can Be Surprisingly Accurate (in Some Regions)

24% sequence identity



25% sequence identity





Welcome to MODBASE, a database of three-dimensional protein models calculated by [comparative modeling](#).

### General Information

## Glossary

#### Authors and acknowledgements

### Publications

[Related resources](#)

Users of ModBase are requested to cite this article in their publications:  
[ModBase, a database of annotated comparative protein structure models](#),  
 Ursula Pieper, Narayanan Eswar, Ashley C. Stuart, Valentin A. Blyn, Andrej Sali  
*Nucl Acids Res.* **30**, 255-259, 2002.

ModBase is maintained by [Ulrich Papier](mailto:Ulrich.Papier@rockefeller.edu) in the group of [Andres Sali](mailto:Andres.Sali@rockefeller.edu), Laboratories of Molecular Biophysics, Fels Family Center for Biochemistry and Structural Biology, The Rockefeller University, 1230 York Ave, New York, NY 10021. Please address all inquiries to [modbase@post.rockefeller.edu](mailto:modbase@post.rockefeller.edu).

© [The Rockefeller University](#)

837,698 [Reliable Models](#) or [PSI-BLAST Fold Assignments](#)  
for domains in 415,937 proteins. Last Update on 04/03/02.  
MODEBASE [statistics](#).

Enter SwissProt/TrEMBL/GenBank/PDB identifier or descriptor:

I Search

[Advanced Search](#)

**HELP**

**Academic login**      **User login**      **Logout**  
**Current logins:** *modbase*.

Some datasets are accessible freely without a login (ie, the "public" model set). Some datasets are available to academic users only (ie, our "SP/TR" model set). And some datasets require a specific username and password. For commercial access to the models, please contact [Structural Genomic Inc.](#)

MODBASE contains theoretically calculated models, not experimentally determined structures. The models may contain significant errors.

### SUMMARY Search Criteria

Summary		Search Criteria	
Keywords	dhfr		
Category	-		
Properties	( % Seq. Ident. and Model Size and Model Score )		
Ranges (min-max)	-30	-	-
Values			
Minimum	8.00	82	0.01
Average	25	181	0.79
Maximum	30.00	492	1.00

[View sequences that match the search criteria but could not be modeled.](#)

29 matches were found using the specified search criteria. Click on the links in the table header to resort your output

TARGET					MODEL DATA					TEMPLATE			
Model Fold Reliability	Sequence Based View	Select Sequence Database Links	Database Description	Organism	Protein Size	Model Segment	Size (aa)	E-value	Model Score	PDB (code)	Template Based View	Segment	Description
		<a href="#">NR_027412</a>	DHVRPOLATE REDUCTASE TYPE VIII (EC 1.5.1.3) (DHFR TYPE IBC) Unlabeled: 199A-2011 PFAM PRODOM	<i>Escherichia coli</i> <i>Shigella sonnei</i>	168	1-165	165	30.00	4e-32	1.00		1-158	DHVRPOLATE REDUCTASE
		<a href="#">NR_037362</a>	BIPHUNCTIONAL DHVRPOLATE REDUCTASE-INTERMEDIATE SYNTHASE (DHFR-18-BLUE)-DHVRPOLATE REDUCTASE (EC 1.5.1.3) Unlabeled: 199A-2011 PFAM PRODOM	<i>Leishmania major</i>	520	24-231	209	30.00	2e-49	1.00		1-186	DHVRPOLATE REDUCTASE (EC 1.5.1.3) (DHFR-18-BLUE POLATE 100F1)
		<a href="#">NR_046714</a>	BIPHUNCTIONAL DHVRPOLATE REDUCTASE-INTERMEDIATE SYNTHASE (DHFR-18-BLUE)-DHVRPOLATE REDUCTASE (EC 1.5.1.3) Unlabeled: 199A-2011 PFAM PRODOM	<i>Plasmodium chabaudi</i>	583	21-241	221	30.00	3e-38	1.00		2-193	SVPR3 DHVRPOLATE REDUCTASE
		<a href="#">NR_050861</a>	BIPHUNCTIONAL DHVRPOLATE REDUCTASE-INTERMEDIATE SYNTHASE (DHFR-18-BLUE)-DHVRPOLATE REDUCTASE (EC 1.5.1.3) Unlabeled: 199A-2011 PFAM PRODOM	<i>Plasmodium vivax</i>	623	25-237	203	30.00	1e-37	1.00		14-203	SVPR3 DHVRPOLATE REDUCTASE
Sequence based View	Select Sequence Database Links	Database Description			Organism	Protein Size	Model Segments - Schematics						
	<a href="#">NR_027412</a>	DHVRPOLATE REDUCTASE TYPE VIII (EC 1.5.1.3) (DHFR TYPE IBC) Unlabeled: PFAM PRODOM			<i>Escherichia coli</i> <i>Shigella sonnei</i>	169							
	<a href="#">Tr_O91801</a>	DHFRP2 PROTEIN (FRAGMENT) Unlabeled: PFAM PRODOM			<i>Homo sapiens</i>	121							

10 20 30 40 50 60 70 80 90 100  
 1vdrA: 2-157: LTVSVAALAEINWIGRDGELPUPSIADKKATSEIADDPVVLGWTITLREDDLPESAQIVNSRSESESVDTAHRAASVIAVAVIAANLDAETAYTGG  
 model: 2-156: KVSLLAKAKGVIGCPPIVSWAKQKLFKALTNGVLVGRKTIPTSGALPNRYTAVTVISGWTSNDONVVVFGSLEAMDLAEITQHVIVSGG  
 Similarity:::1  
 Signed by: Unsigned classes from local hard disk

PSI-BLAST Fold Assignments (left half) and Delinix Models (right half) are indicated in green.

\* Indicates an E-value from an unfiltered PSI-BLAST search when a filtered search does not result in a significant match.

**This Table displays all Models/Folds of this sequence**

Full View Download ModView

Database Synonyms for this Sequence (100% Sequence Identity)				
TrEMBL	<a href="#">Q93V12</a>	<i>Salmonella typhimurium</i>	cydohydroxide reductase	EFAN 0000
GI	<a href="#">133322</a>	Flacmid pl.M0229	dhfr product (AA.1 - 157)	----
GI	<a href="#">700750</a>	<i>Salmonella typhimurium</i>	dihydroxyolate reductase	----
GI	<a href="#">96785</a>	<i>Escherichia coli</i> plasmid p.M0229	811700 dihydroxyolate reductase (EC 1.5.1.3) type 1 - <i>Escherichia coli</i> plasmid p.M0229	----

MODEL DATA					LINKS				
Fold/Model Reliability	Size	Seq ID (%)	E-value	Model Score	Mod- TBL	3D Coor	3D Str	PDB	Mod View

(157 Residues)

3 -

**Trac** (2-156) DIHYDROFOLATE REDUCTASE - **CATH** 2.40.4.30.10.1.1 (99%) Subset: SPITR-2001

2 -

**TracA** (1-156) DIHYDROFOLATE REDUCTASE - **CATH** 2.40.4.30.10.5.2 (99%) Subset: SPITR-2001

File Edit View Go Communicator Help

Back Forward Reload Home Search Netscape Print Security Shop Stop

Bookmarks Go To: <http://pipe.rockefeller.edu/modloop/modloop/lind> What's Related

# MOD LOOP Modeling of Loops in Protein Structures

A. Fiser, R.K.G. Do and A. Sali, *Prot Sci.* 9, 1753-1773 (2000) [PDF](#)

**Upload your coordinate file :**

**Select loop segments :**

```
12:A:18:A:
54:B:62:B:
78:B:78:B:
```

**Number of iterations :**  Range(1 - 200)

**Name of your model :**

**Your e-mail (required) :**

**MODELLER key (required) :**

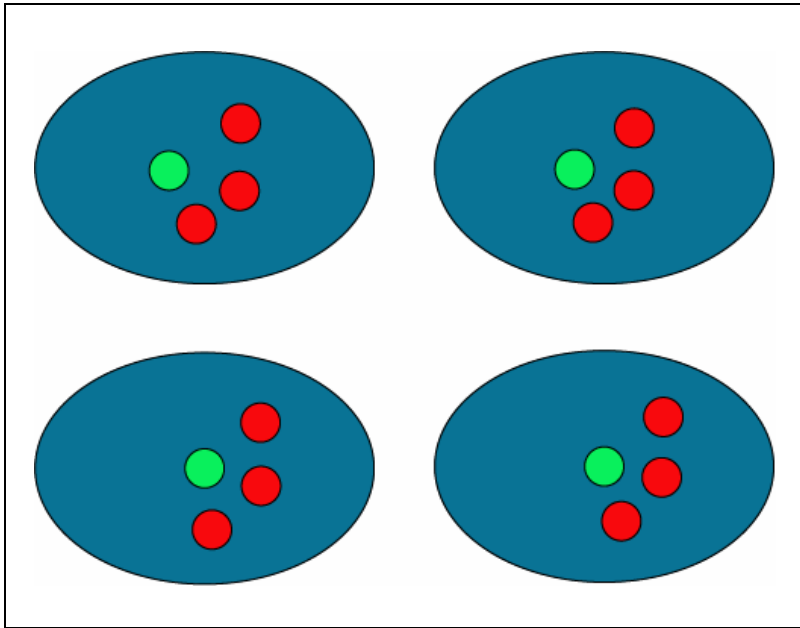
# Applications I.

## Structural genomics



# Structural Genomics

Characterize most protein sequences (**red**) based on related known structures (**green**).



The number of “families” is much smaller than the number of proteins

# Structural Genomics

- **Definition:** The aim of structural genomics is to put every protein sequence within a modeling distance of a known protein structure.
- **Size of the problem:**
  - There are a few thousand domain fold families.
  - There are ~20,000 sequence families (30% sequence id).
- **Solution:**
  - Determine protein structures for as many different families as possible.
  - Model the rest of the family members using comparative modeling

Burley et. al. Nat. Genet. 23, 151, 1999.

Sanchez et. al. Nat. Str. Biol. 7, 986, 2000



## New York Structural Genomics Research Consortium

### Mission Statement

To develop and use the technology for high-throughput structural and functional studies of proteins.

### Participating Research Groups

#### Albert Einstein College of Medicine

Mark R. Chance

Steve Almo

Anne Bresnick

Andras Fiser

#### Brookhaven National Laboratory

Robert Sweet

Jian-Sheng Jiang

F. William Studier

S. Swaminathan

#### Columbia University

Lawrence Shapiro

#### Structural Genomix, Inc

Stephen K. Burley

#### The Rockefeller University

Terry Gaasterland

#### University of California, San Francisco

Andrej Sali

#### Weill Medical College of Cornell

University

Christopher Lima

### Public Target Information

Public Target Progress Report

Download: Public Target Progress Report in XML Format

[Home](#)

[Proposal](#)

[Publications](#)

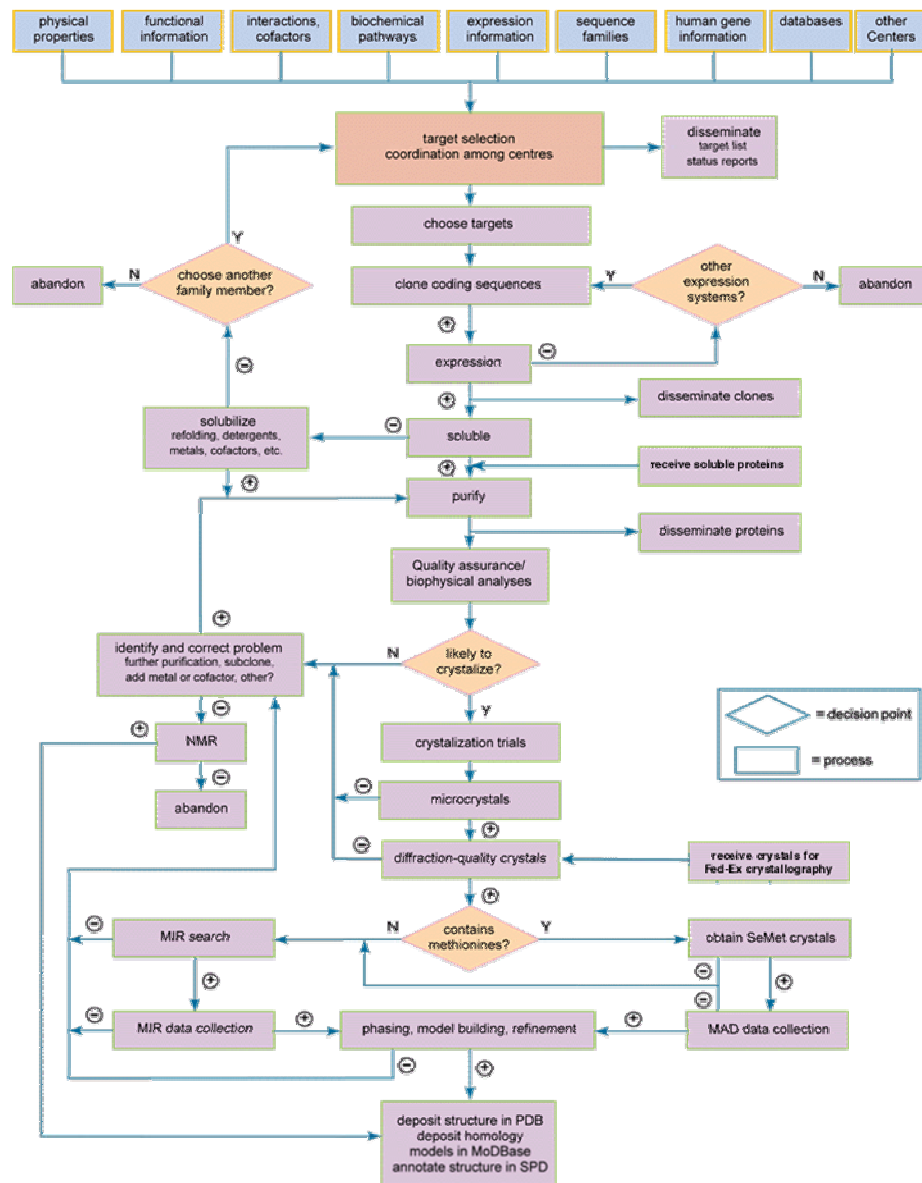
[Flowchart](#)

[IceDB](#)

[Tools](#)

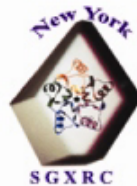
[Contact](#)





### Comparative Protein Structure Modeling with NYSGXRC Structures (January 23, 2003)

NYSGXRC SOLVED-STRUCTURE TEMPLATE					MODBASE	NYSGXRC ACCEPTABLE MODEL DATA					
Target_ID	Protein Name / Comment	GI or Swissprot Code	PDB Code	Protein Size	Total Models	Total Models	Min. Seq. ID	Max. Seq. ID	# Models >50% Seq. ID	# Models 30-50% Seq. ID	# Models <30% Seq. ID
<a href="#">T9</a>	Similar to putative GTP-binding protein	<a href="#">Q13998</a>	<a href="#">1NI3</a>	392	<a href="#">1253</a>	76	13	63	6	40	30
<a href="#">T132</a>	Putative cell cycle protein mesJ	<a href="#">P52097</a>	<a href="#">1NI5</a>	432	<a href="#">870</a>	150	11	56	1	15	134
<a href="#">T503</a>	Conserved hypothetical protein YDCE from Bacillus Subtilis	<a href="#">P96622</a>	<a href="#">1NE8</a>	116	<a href="#">52</a>	36	21	81	9	10	17
<a href="#">TBA</a>	unc78	<a href="#">TBA</a>	<a href="#">TBA</a>	1315	<a href="#">693</a>	0	0	0	0	0	0
<a href="#">TBA</a>	phlp1, a Major Timothy Grass Pollen Allergen	<a href="#">TBA</a>	<a href="#">1N10</a>	228	<a href="#">191</a>	162	20	92	18	53	91
<a href="#">T746</a>	phlp6	<a href="#">P43215</a>	<a href="#">TBA</a>	209	<a href="#">249</a>	11	25	47	0	10	1
<a href="#">P089</a>	Hypothetical 32.1 kDa protein in ADH3-RCA1 intergenic region	<a href="#">Q04299</a>	<a href="#">1NJR</a>	284	<a href="#">1</a>	0	0	0	0	0	0
<a href="#">P096</a>	Hypothetical 28.8 kDa protein in PSD1-SKO1 intergenic region	<a href="#">P53889</a>	<a href="#">1NKQ</a>	259	<a href="#">301</a>	161	14	43	0	82	79
<a href="#">TBA</a>	putative_thioesterase_(comA) polymer	<a href="#">TBA</a>	<a href="#">TBA</a>	138	<a href="#">479</a>	147	9	89	17	16	114
<a href="#">TBA</a>	hypothetical_protein_(yqeU) polymer	<a href="#">TBA</a>	<a href="#">TBA</a>	241	<a href="#">85</a>	76	18	81	6	16	54
<a href="#">T299</a>	URACIL-DNA GLYCOSYLASE FROM T. MARITIMA (Hypothetical protein TM0511)	<a href="#">Q9WYY1</a>	<a href="#">1L9G</a>	192	<a href="#">110</a>	84	13	49	0	49	35
<a href="#">P007</a>	Hypothetical 29.1 kDa protein in URA7-POL12 intergenic region	<a href="#">P38197</a>	<a href="#">1B54</a>	257	<a href="#">53</a>	44	27	43	0	34	10
<a href="#">P008</a>	Pyridoxamine 5'-phosphate oxidase	<a href="#">P38075</a>	<a href="#">1CI0</a>	228	<a href="#">1374</a>	1266	8	99	103	33	1130
<a href="#">P018</a>	Hypothetical 32.5 kDa protein YLR351C	<a href="#">P49954</a>	<a href="#">1F89</a>	291	<a href="#">302</a>	251	13	54	1	38	212
<a href="#">P044a</a>	L-allo-threonine aldolase	<a href="#">GI: 4982322</a>	<a href="#">1JG8</a>	343	<a href="#">1049</a>	923	10	46	0	13	910
<a href="#">P068</a>	Hypothetical 33.9 kDa esterase in SMC3-MRPL8 intergenic region	<a href="#">P40363</a>	<a href="#">TBA</a>	299	<a href="#">804</a>	119	10	52	1	14	104
<a href="#">P097</a>	Hypothetical 27.5 kDa protein in SPX19-GCR2 intergenic region	<a href="#">P40165</a>	<a href="#">1JZT</a>	246	<a href="#">1401</a>	6	14	41	0	4	2
<a href="#">P100</a>	Diphosphomevalonate decarboxylase	<a href="#">P32377</a>	<a href="#">1FI4</a>	396	<a href="#">422</a>	139	9	68	2	26	111
<a href="#">P102</a>	Glutathione synthetase Apo	<a href="#">Q08220</a>	<a href="#">1M0T</a>	491	<a href="#">140</a>	25	30	42	2	23	0
<a href="#">P102a</a>	Glutathione synthetase Lig	<a href="#">Q08220</a>	<a href="#">1M0W</a>	491	<a href="#">33</a>	30	32	39	0	30	0
<a href="#">P109a</a>	Isopentenyl-diphosphate delta-isomerase (IPP isomerase)	<a href="#">GI: 6225535</a>	<a href="#">1I9A</a>	182	<a href="#">947</a>	417	10	72	5	18	394
<a href="#">P111a</a>	Translation initiation factor 6	<a href="#">Q60357</a>	<a href="#">1G61</a>	228	<a href="#">37</a>	35	29	46	0	34	1



## IceDB Report Display

Target ID: P097  
Target Iteration: 1

[Go to Upload Report Page](#)  
[Download Report as PDF Format](#)  
**To make changes to other  
experimental types:**  
[Go to Add/Edit/Upload Data - Start  
Page](#)

**Target Identifier:** P097

**Protein Name:** HYPOTHETICAL 27.5 KDA PROTEIN IN SPX19-GCR2  
INTERGENIC REGION

**Organism:** *Saccharomyces cerevisiae*

**PDB Code:** 1JZT

**Rationale for Target Selection:** Unknown function and represents a domain from a number of protein families (ProDom accession PD005835).

**Method of Structure Determination:** Se-Met MAD method and NCS density averaging.

**Structure Description:** The structure of P097 is a three-layer a-b-a sandwich. The two molecules related by NCS in the asymmetric unit form a tightly packed dimer. Each monomer consists of eight b-strands and nine a-helices. The order of b-strands is 32145678 according to the SCOP classification.



**Comparisons of Structurally Similar Proteins in PDB:** P097 represents an unusual Rossmann fold. A typical NAD-binding Rossmann fold should have six b-strands forming an open twisted parallel b sheets in the middle and two a-helices on both sides (3LDH). The secondary structure of the Rossmann fold is b1-aA-b2-aB-b3-aC-b4-aD-b5-aE-b6 and the order of b-strands is 321456 (SCOP). The NAD-binding proteins usually consist of two domains: a dinucleotide (or

[Structural Genomics Initiatives](#)[Structural Genomics Link at the PDB](#)

## TargetDB

### Target Search for Structural Genomics

TargetDB is a target registration database that was originally developed to provide registration and tracking information for NIH P50 structural genomics centers. TargetDB has now been expanded to include target data from worldwide structural genomics and proteomics projects. The scope of TargetDB is to provide timely status and tracking information on the progress of the production and solution of structures.

Sequences from the NIH P50 and other structural projects have been loaded into the TargetDB database and can be searched using the form below. TargetDB is updated weekly. All targets are available for download in [XML format here](#).

A new [Target Status Query Feature](#) is now available, please click [here](#).

Target sequence lists are also maintained at the following sites:

| [BIGS](#) | [BSGC](#) | [BSGI](#) | [JCSG](#) | [MCSG](#) | [MSGP](#) | [NESG](#) | [NYSGRG](#) | [OPPF](#) | [PSF](#) | [RIKEN](#) | [S2F](#) | [SECSG](#) | [SGPP](#) | [TB](#) | [YSG](#) |

#### Using the Target Search Form:

- Enter text and/or select menu options in the form below to define the desired target search, select a result format, and press the SUBMIT button to execute the query.
- All form attributes are optional. If no options are entered a query will return all of the NIH target entries in the database.
- Click on any attribute name for an explanation and examples of the attribute.
- For a FASTA sequence comparison, enter the one-letter code sequence into the sequence text box.

[Project](#)[Target ID:](#) [Status:](#)[Site:](#)[Include](#)[Data](#)[From:](#)[Target](#)[Data](#)[Updated:](#)after   before   [Protein](#)[Name:](#)[Source](#)[Organism:](#)

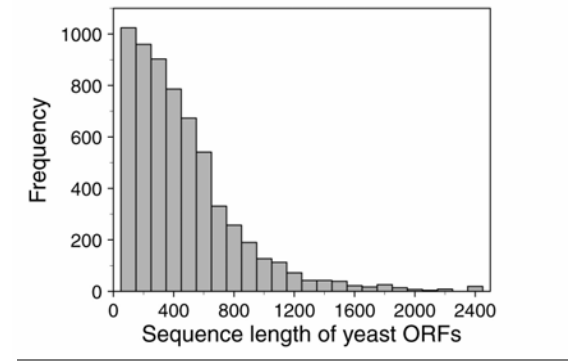
# Limitations of structural genomics: Quaternary structure

Proteins are modular, ~2.7 domains per protein.

Evolution shuffles domains.

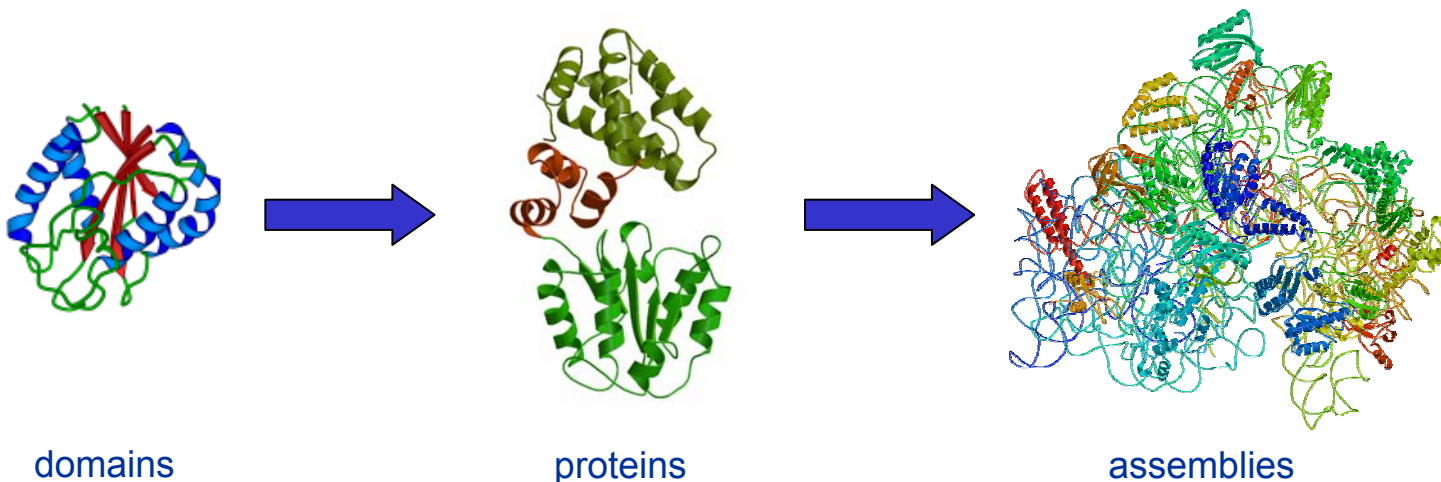
24% of domains/residues in 57% of proteins are modeled

Structural genomics is determining structures of domains, usually not proteins, definitely not assemblies.



Average Length		
Protein	Domain	Model
472	175	192

Thus, there is a great need for methods for docking of domains into proteins and of proteins into macromolecular assemblies.



# Acknowledgements

- Andrej Sali (Rockefeller/UCSF)
- Charlie Brooks (Scripps)
- Michael Feig (Scripps)
- Steve Burley (Rockefeller/SGX)
- Steve Almo (Albert Einstein Coll. Of Med.)
- Mark Chance ( Albert Einstein Coll. Of Med.)

## Lab

- Guiping Xu
- Eduardo Fajardo
- Dmitry Rykunov
- Brajesh Kumar Rai
- Narcis Fernando-Fuentes
- Rotem Rubinstein

# Reviews

Fiser, A.

Protein structure modeling in the proteomics era

*Exp. Rev. in Proteomics* . (2004) 1, 89-102

Fiser, A. and Sali, A.

MODELLER: Calculating and refining homology models

*Methods Enzymol.* (2003) 374,463-493

Fiser, A., Sanchez, R., Melo, F. and Sali, A.

Comparative protein structure modeling.

in *Computational Biochemistry and Biophysics*, (2001) pp. 275-312,

Marcel Dekker. Eds. M. Watanabe, B. Roux, A. MacKerell, and O. Becker.