



LABORATORY
of THEORY of
BIOPOLYMERS

Proteins - structural bioinformatics (2)

<http://biocomp.chem.uw.edu.pl>

How many proteins?

- Just 150 AA protein - 10^{195} sequences
- Eukaryotic protein universe $\sim 10^{12}$
- Prokaryotic – much more, difficult to estimate
- 12-14 thousand of known protein families cover about 60% of known proteins
- 5000 – 20,000 of possible folds (about 1500 currently known)

Tertiary Structure and the “Hydrophobic Effect”

What would this protein look like when properly folded?

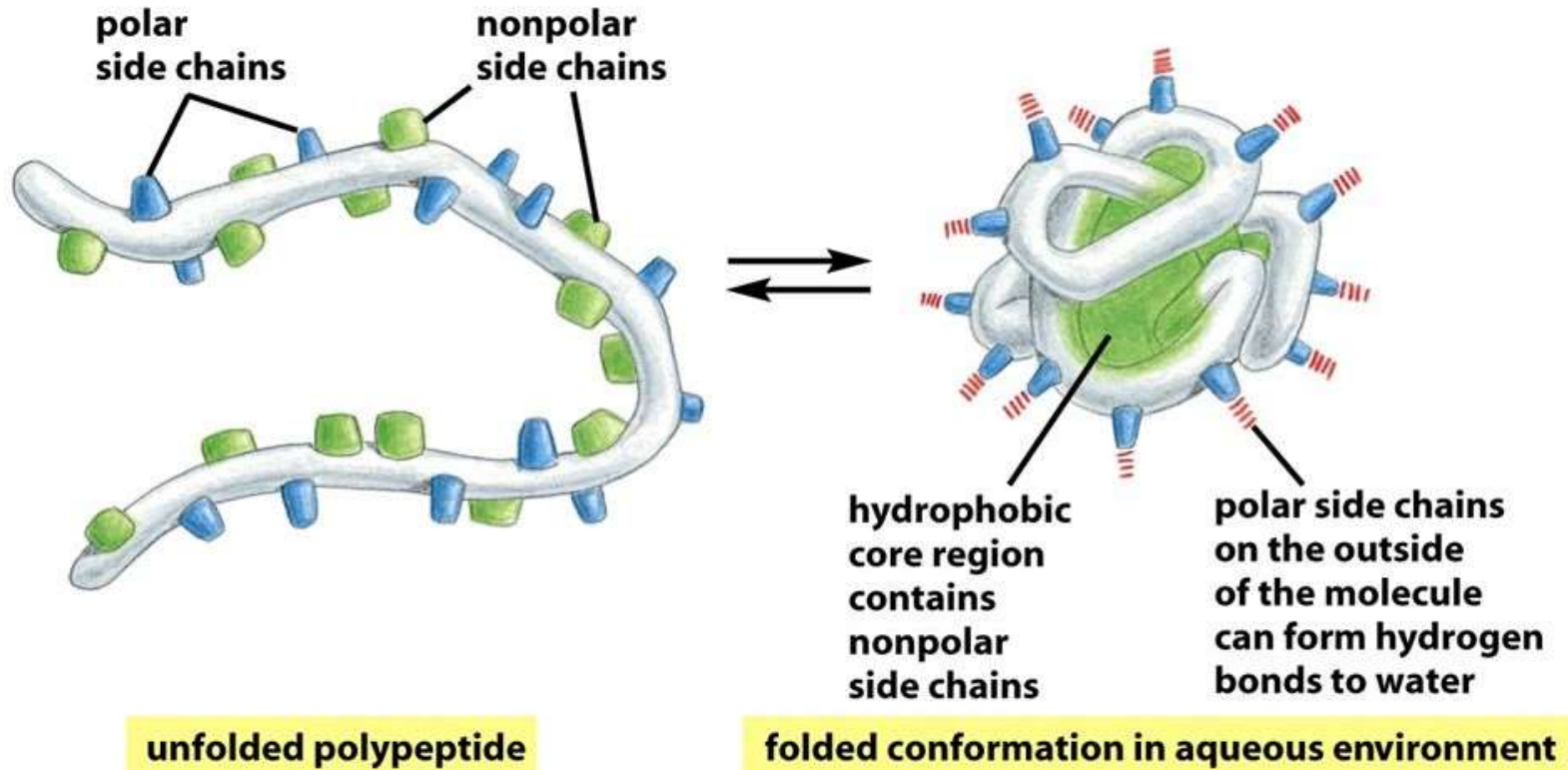
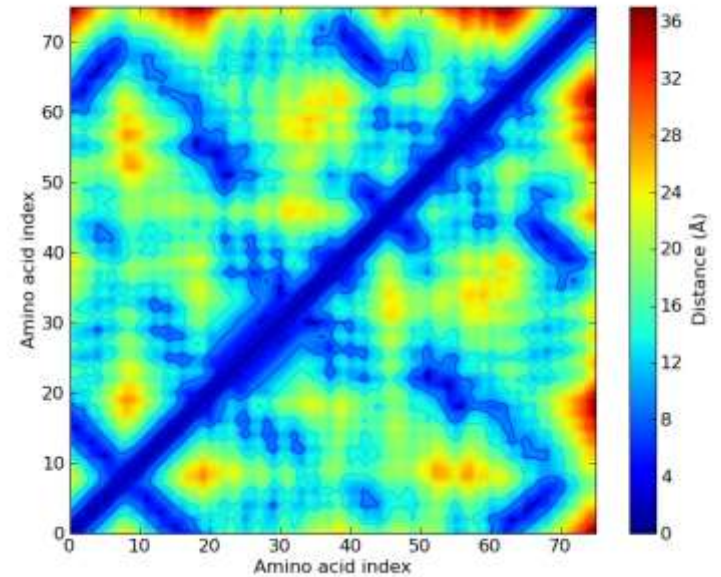
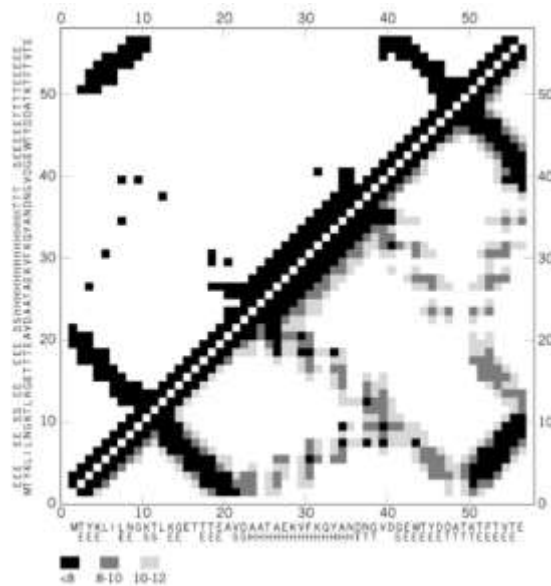
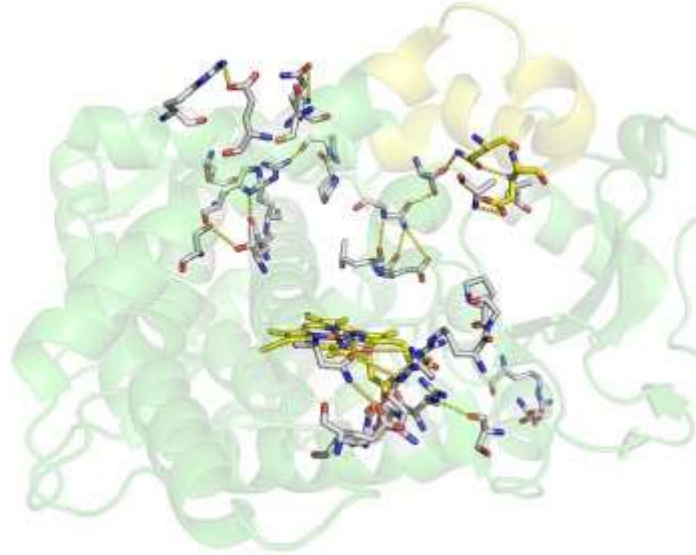
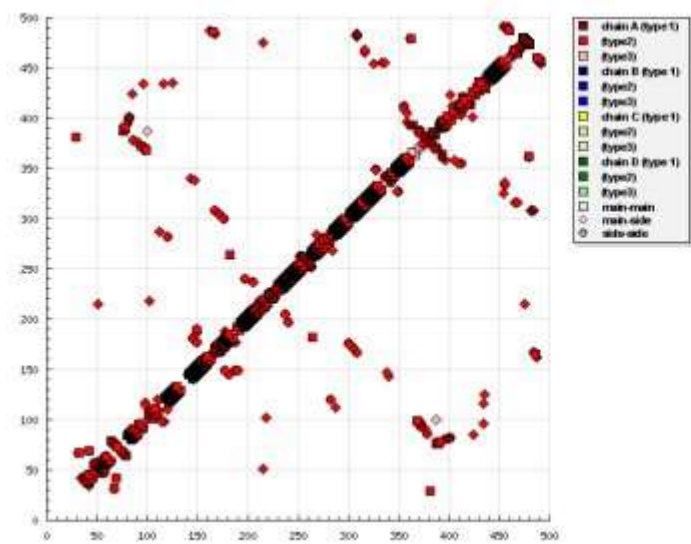
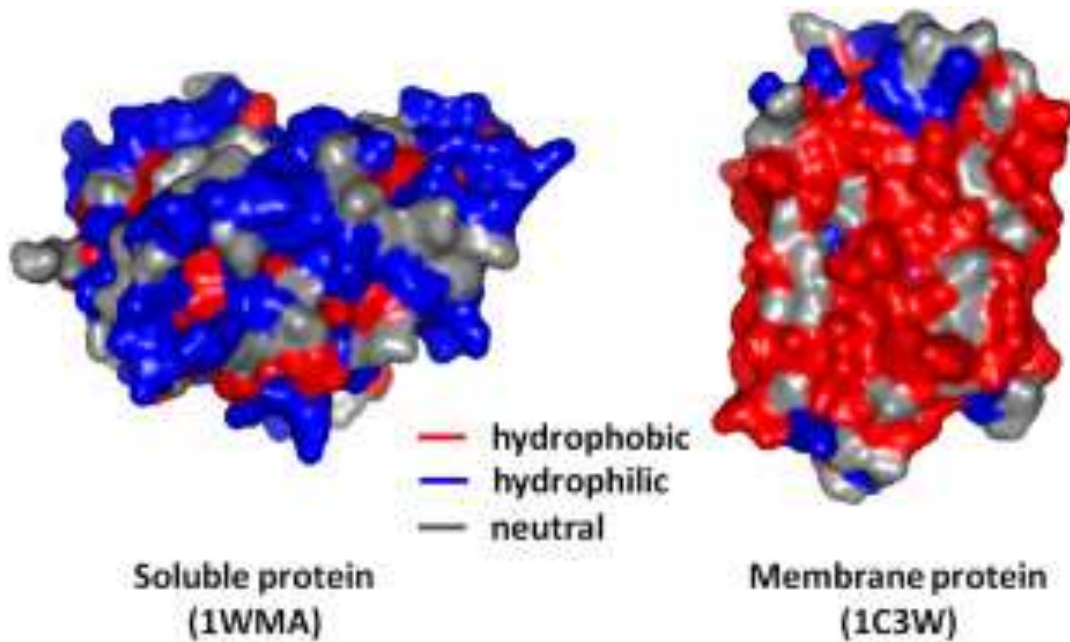


Figure 3-5 *Molecular Biology of the Cell* (© Garland Science 2008)

Side chain packing



Hydrophobic effects



Tertiary Structure and the “Hydrophobic Effect”

What would this protein look like when properly folded?

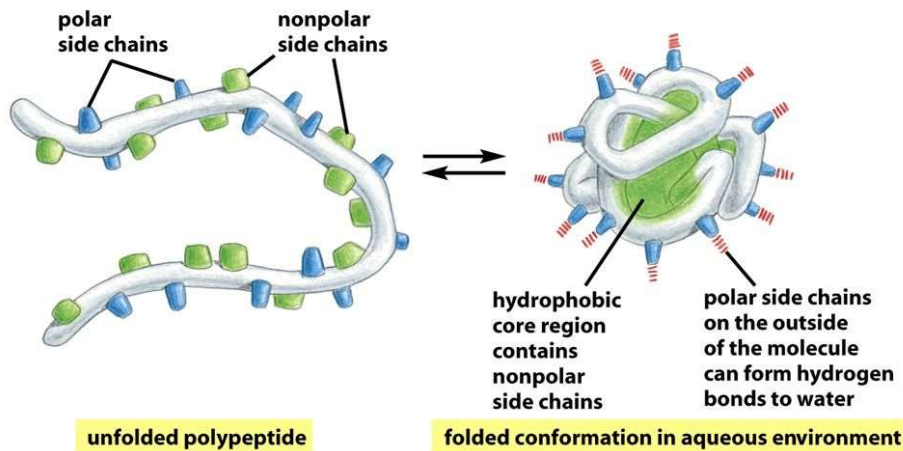
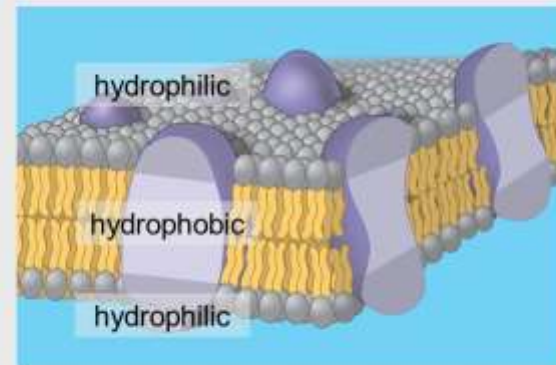


Figure 3-5 Molecular Biology of the Cell (© Garland Science 2008)

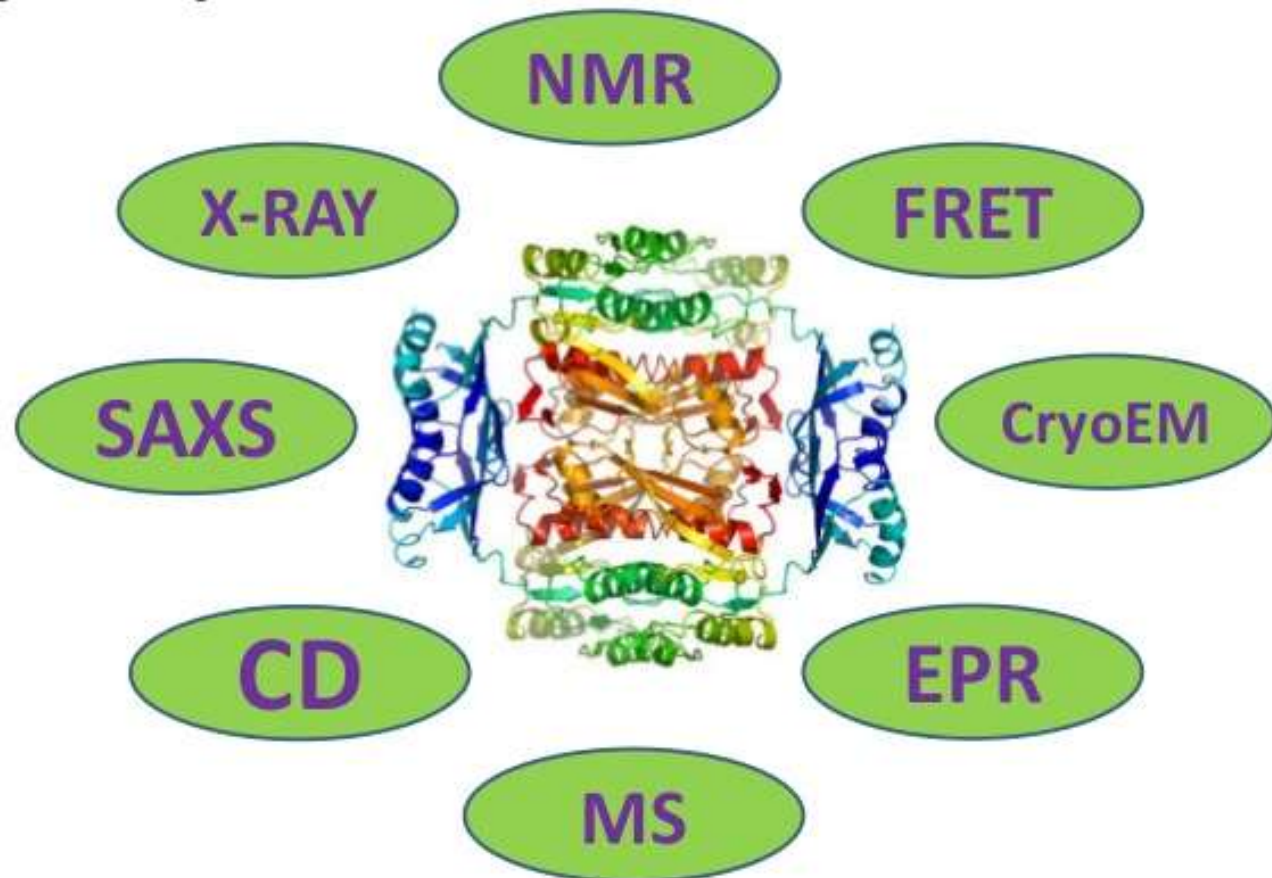
Membrane structure

- cell membrane – amphipathic - hydrophilic & hydrophobic



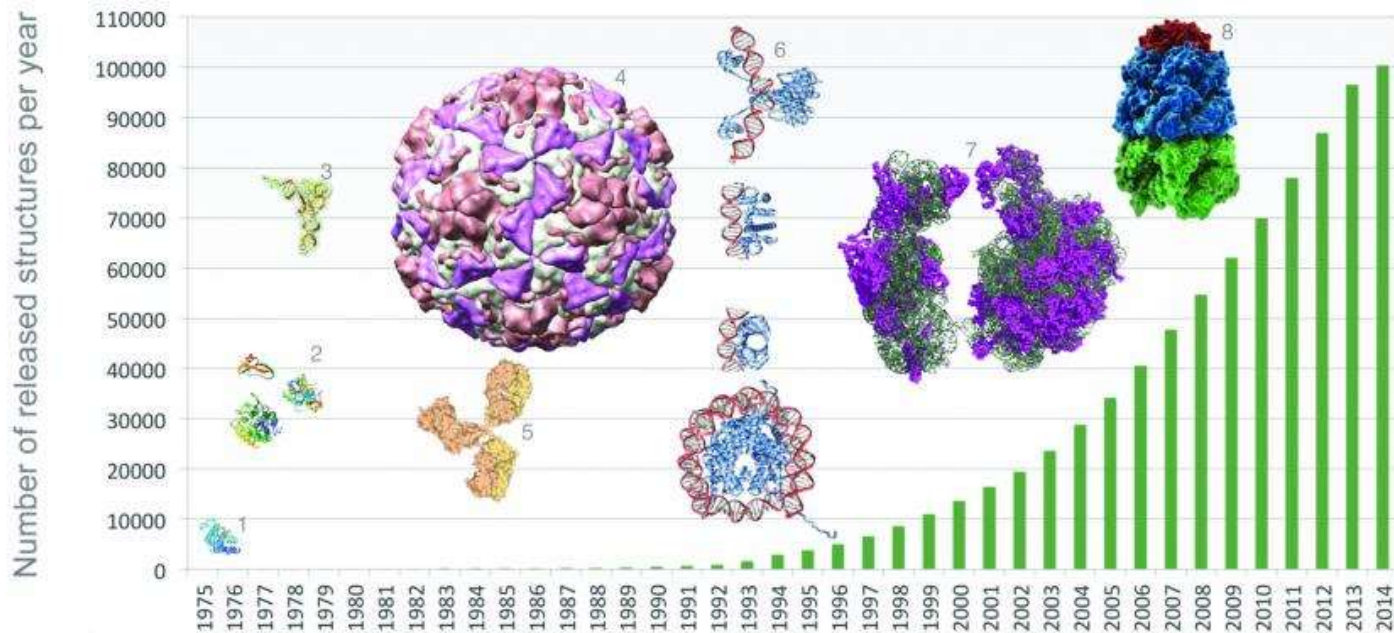
- membrane proteins that are inserted, also amphipathic

Hybrid protein structure determination

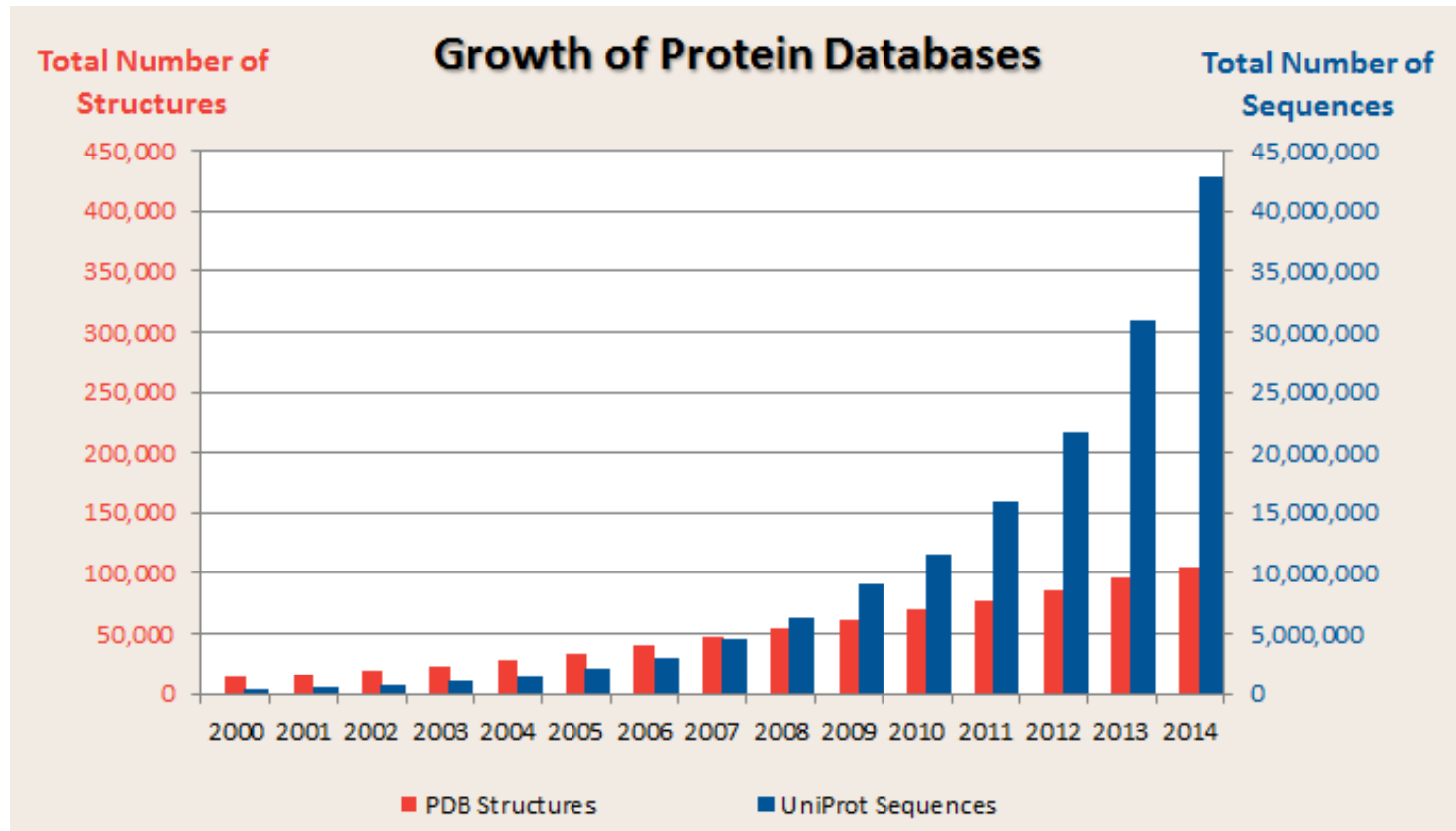




Research Collaboratory for Structural
Bioinformatics:
Rutgers and UCSD/SDSC



Sequence - structure



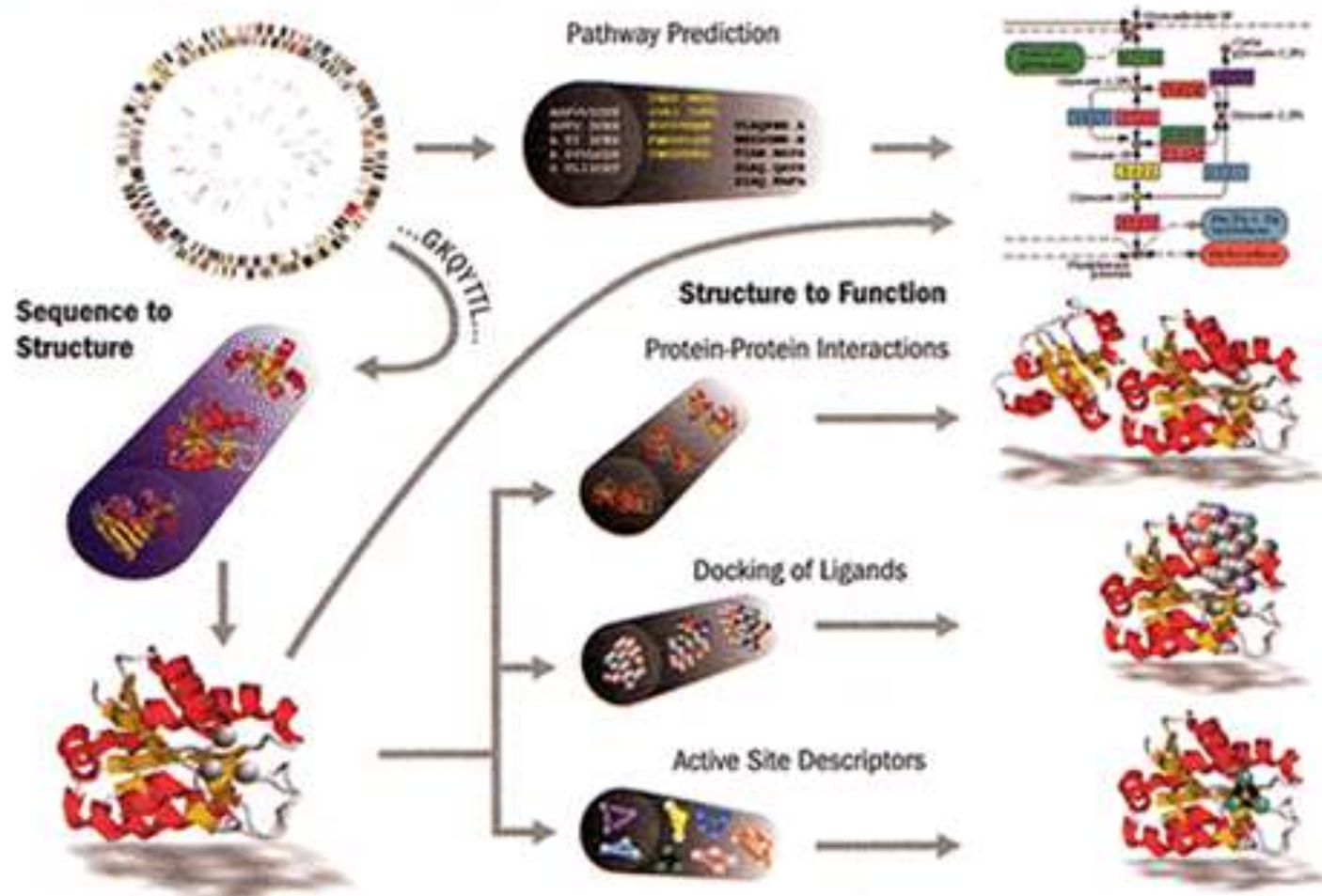
Protein Data Bank (PDB) - 140 000 protein structures

UniProtKB/TrEMBL sequence database - 133 507 323 nonredundant entries . Nov. 2018

Integrated Microbial Genomes & Microbiomes(IMG/M)database of 51 775 423 466 genes

(Coding genes *E. coli* - 4000, yeast – 6000, human, about -20000)

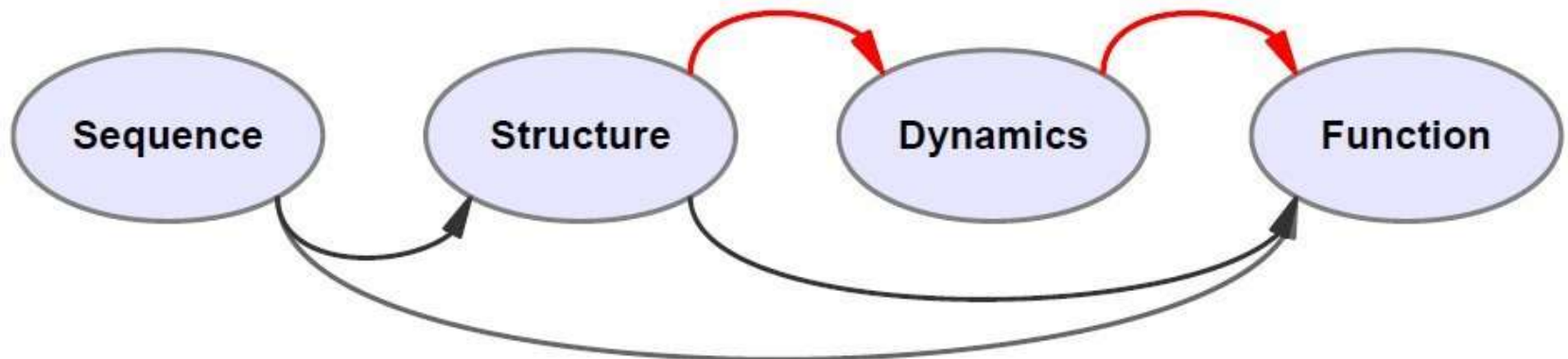
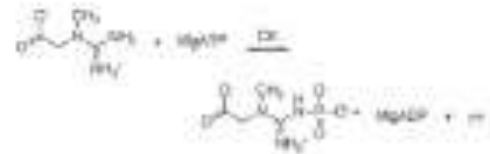
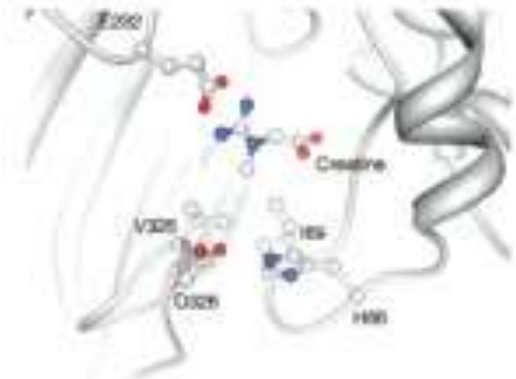
The Sequence-to-Structure-to-Function Paradigm



All the potential open reading frames (ORFs) in a protein sequence are threaded through a library of previously solved template protein structures. If a template is found, the structure is scanned for a match to a known active site. Alternatively, ligands can be virtually docked to identify the active site. Threading can also be used to identify potential interacting partners in the genome, or assist ORF pathway assignment.

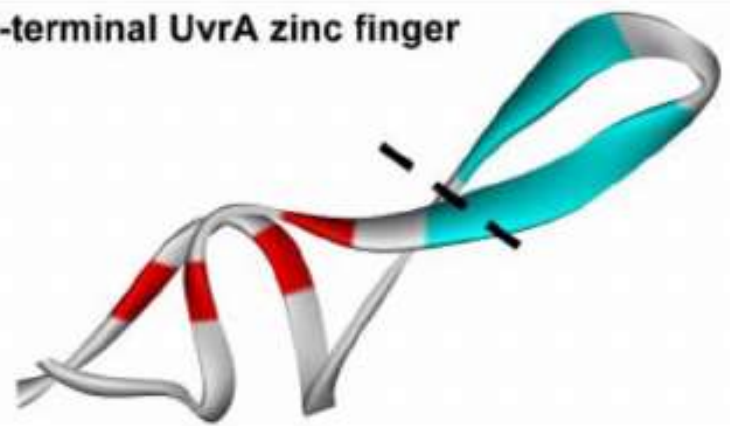
Sequence → Structure → Function

MPFGNTHNKFKL
NYKPEEEYPDLSK
HNNHMAKVLTL
LYKCLRDKETPSGF
TVDDVIQTGVDNP
GHPFIMTVGCVAG
DEESYEYVFKELFDPI
ISDRHGGYKPTD...



Structure – Comparative modeling

| | |
|----------|---|
| A: ECOLI | GRFSFNVRGGR <u>CEACQGDGVIK</u> VEMHFLPDIYVP--- <u>CDQCKGKRYNRETLE</u> |
| RHIME | GRFSFNVKGGRCEACQGDGVIKIEMHFLPDVYVT---CDVCHGKRYNRETLD |
| TREPA | GRFSFNVPGGRCEHCKGDGVITIEMNFLPDVYIT---CDVCHGTRFNRETLA |
| HELPJ | SRFSFNVKGGRCEK <u>CQGDGIKI</u> EMHFLPDVLVQ---CDSCKGAKYNPQTLE |
| BCACA | GRFSFNVKGGRCEACHGDGIKIEMHFLPDVYVP---CEVCHGKRYNRETLE |
| ZnG A | GRFSFNVKGGRCEACHGDGII-----G-----VP---CEVCHGKRYNRETLE |
| Ydj1 | GRGGKKGAVKK <u>CTSCNGQGI</u> KFVTRQMGPMIQRFO <u>TECDVCHGTGDIIDPKD</u> |
| B: Ydj1 | C: C-terminal UvrA zinc finger |



Comparative Modeling--Basic Protocol

42

1. **Identification** of homologue for target sequence
2. **Alignment** of target sequence to template sequence and structure
3. **Side-chain modeling**, copy the backbone of the template and model the new side chains onto this backbone
4. **Loop modeling**, for insertions and deletions in the alignment
5. **Refinement of model** -- moving template closer to target
6. **Assessment** of (predicted) model quality
7. **Using the model** to explain experiments and guide new ones

Sources of errors

- experimental errors and uncertainties in X-ray, NMR

1Å
100%



- side-chain packing
- mis-placed side-chains

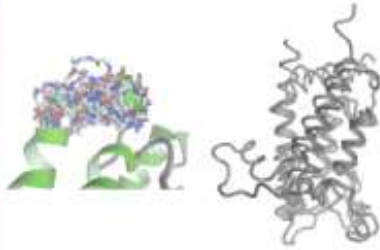
1.5Å
95%



- modeling of loop regions (insertions and deletions)

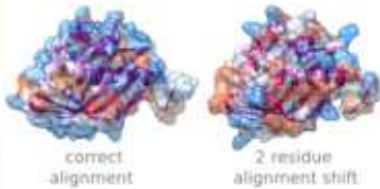
60%

- distortions of aligned regions



- alignment errors

3Å
40%



- sub-optimal template selection

>3Å
<30%

- model may even have the wrong fold



Applications

- studying catalytic mechanism / function

- structure-based drug design, ligand docking

- structural support for mutagenesis studies

- molecular replacement

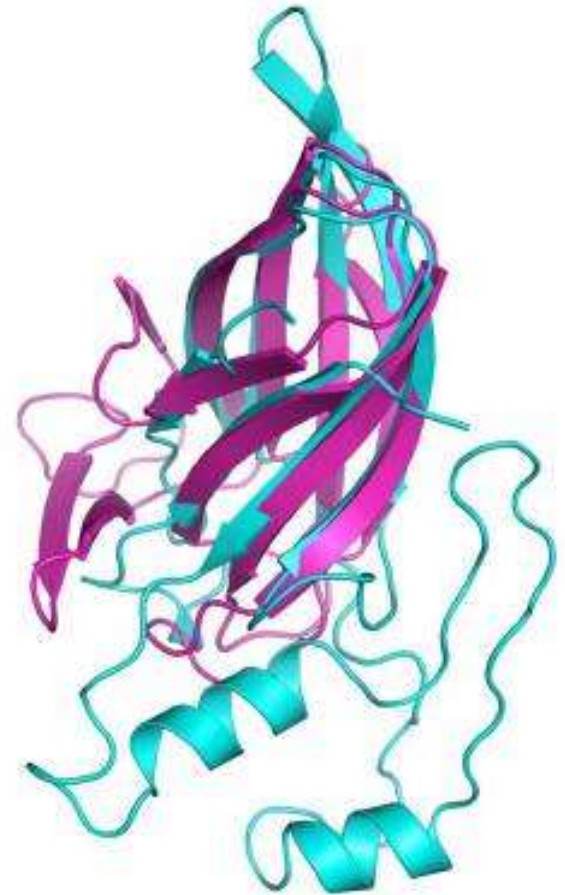
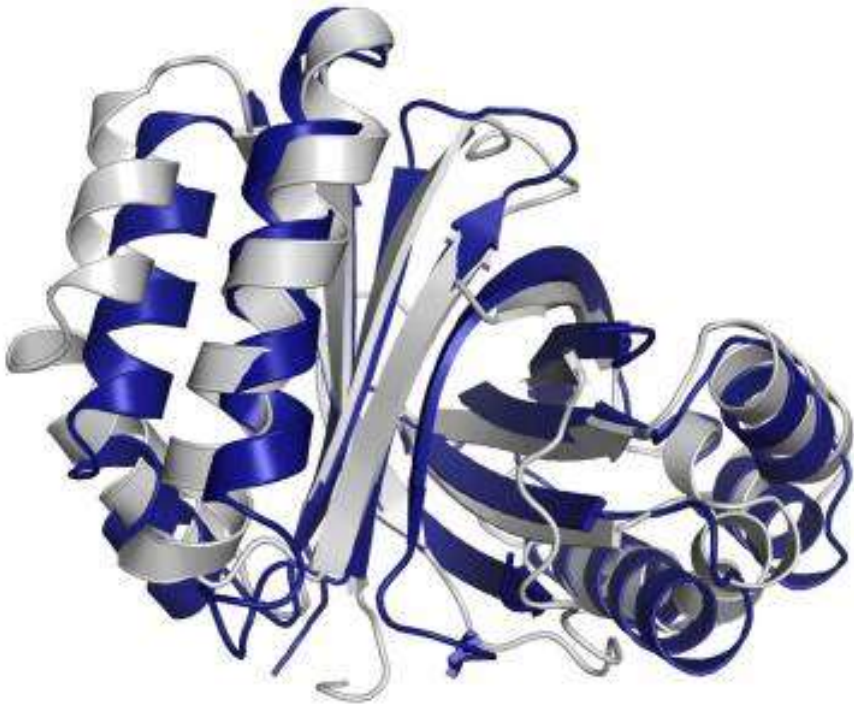
- integrative modeling

- modeling into low-resolution density maps

- domain boundaries

- identification of structural motives

Comparative (homology) modeling



Both cases (A,B) represent extremely distant homologies with sequence identity on the level of 10–12%

A

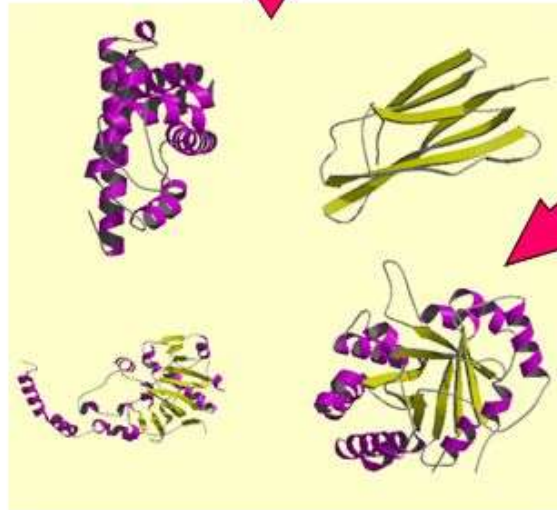
B

Protein folding problem

PRIMARY STRUCTURE (amino acid sequence)

VHLTPEEKSAVTALWGKLVNDE
VGGEALGRLLVVYPWTQRFFE
SFGDLSTPDVAVMGNPKVKAHG
KKVLGAFSDGLAHLNLDLGTFA
TLSELHCDKLHVDPENFRLLGN
VLVLCVLAHHFGKEFTPPVQAA
YQKVVAGVANALAHKYH

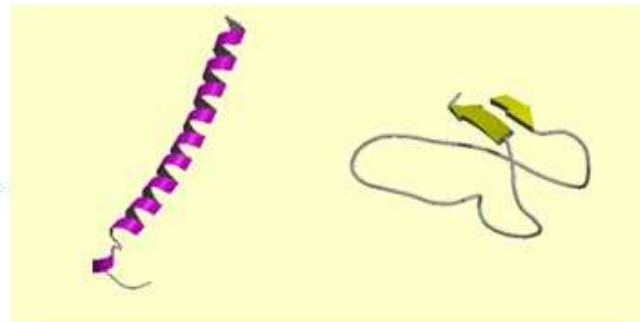
*1-step
process*



TERTIARY STRUCTURE (fold)

*Each protein sequence
“knows” how to fold into its
tertiary structure. We still do
not understand how and why*

SECONDARY STRUCTURE (helices, strands)



*2-step
process*

*The 1-step process is based on a
hydrophobic collapse; the 2-step
process, more common in forming
larger proteins, is called the
framework model of folding*



Post-translational or co-translational folding

5 ribosomes reading same RNA sequentially

Growing polypeptide chains

Complete polypeptide

(Initiator codon)

AUG

5'

UAG

3' mRNA

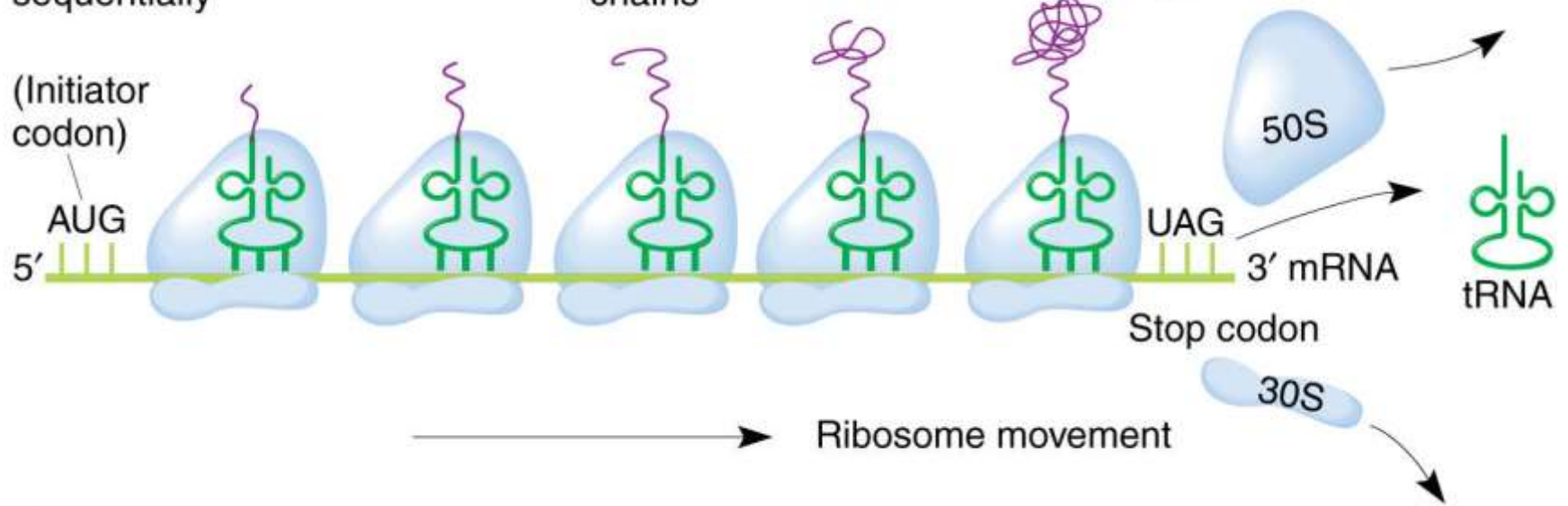
Stop codon

50S

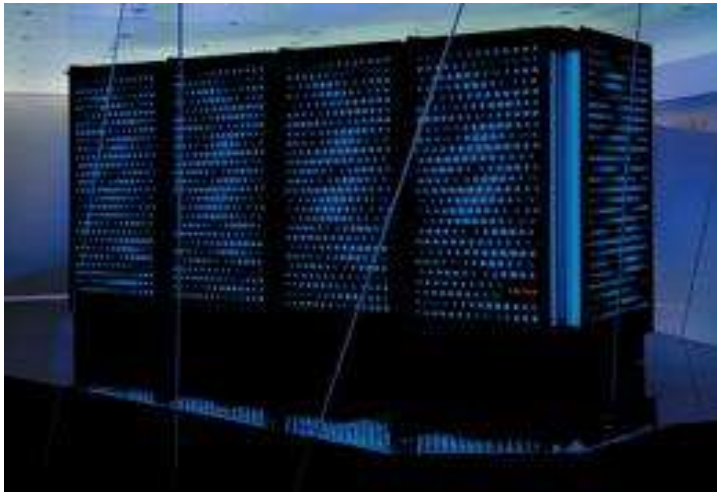
30S

tRNA

Ribosome movement



Protein folding problem - the Holy Grail of the structural biology

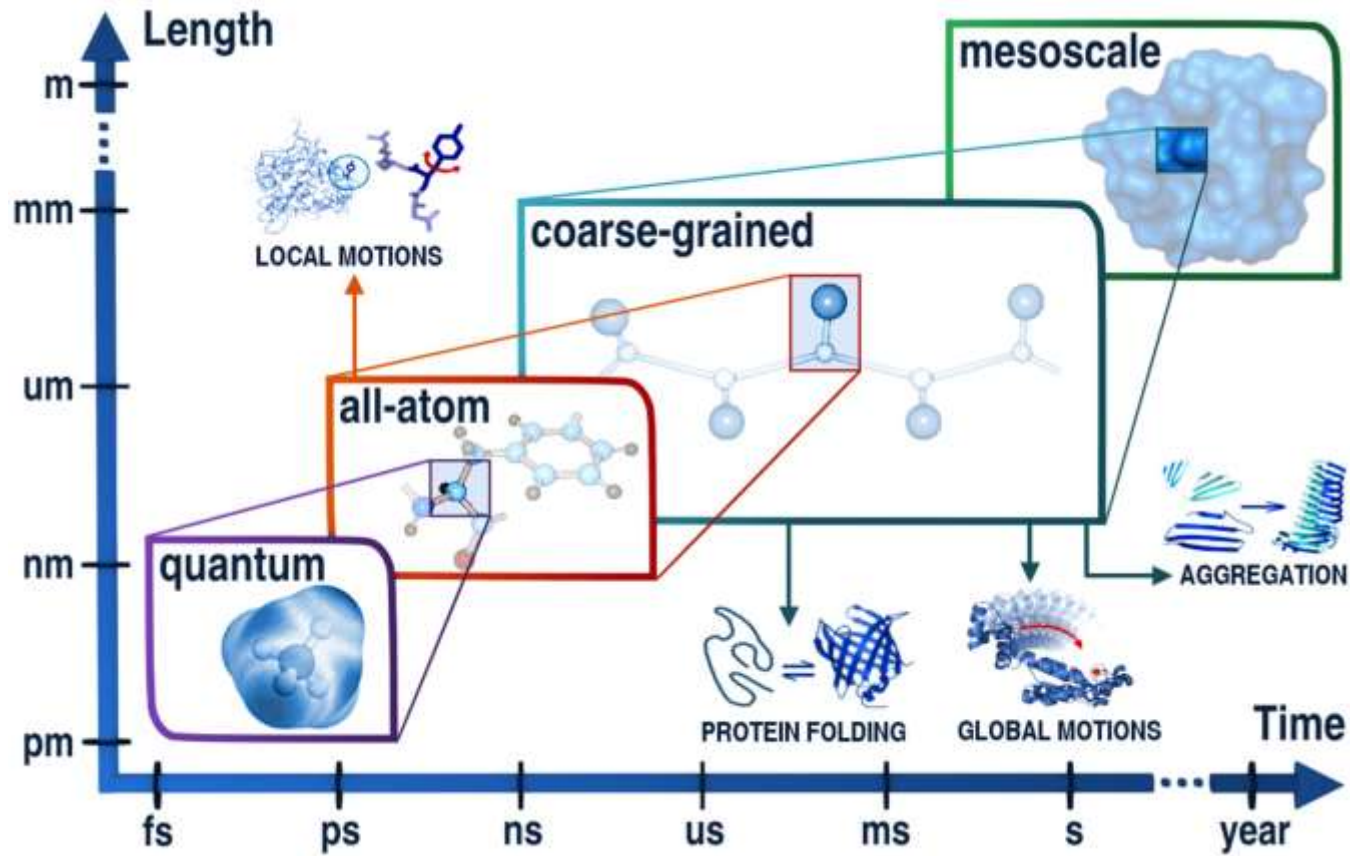


Anton
David E. Shaw Research

All-atom MD with explicit water
- milliseconds of folding process
of a small protein.

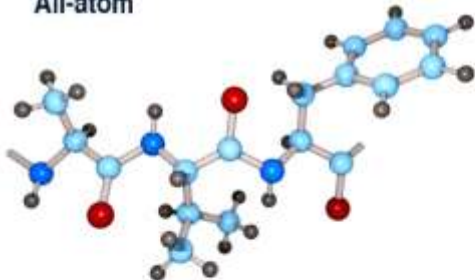
For realistic modeling of larger
biomolecular systems, including
flexible protein-protein docking, **we
need much faster simulations.**

How to solve the Holy Grail problem

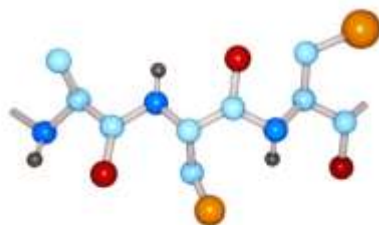


How to solve the Holy Grail problem – Multiscale Modeling

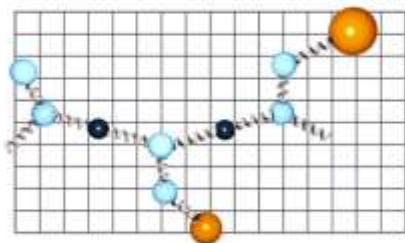
All-atom



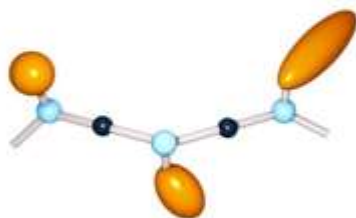
Rosetta CEN



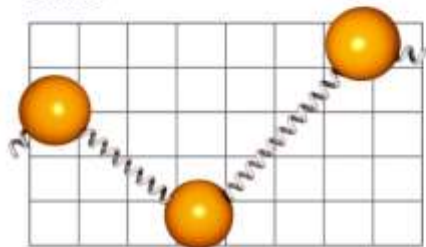
CABS



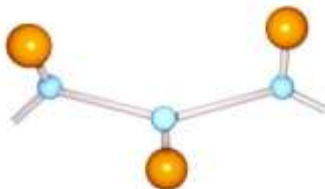
UNRES



SICHO




Levitt-Warshel




Nobelpriset 2013


The Nobel Prize in Chemistry 2013



Martin Karplus
Université de Strasbourg,
France and Harvard
University, Cambridge,
MA, USA



Michael Levitt
Stanford University School of
Medicine, CA, USA



Arieh Warshel
University of Southern
California, Los Angeles, CA,
USA

for "the development of multiscale models for complex chemical systems"

CABS model

C_{α} - C_{β} -Side chain

High-coordination lattice

Statistical force-field

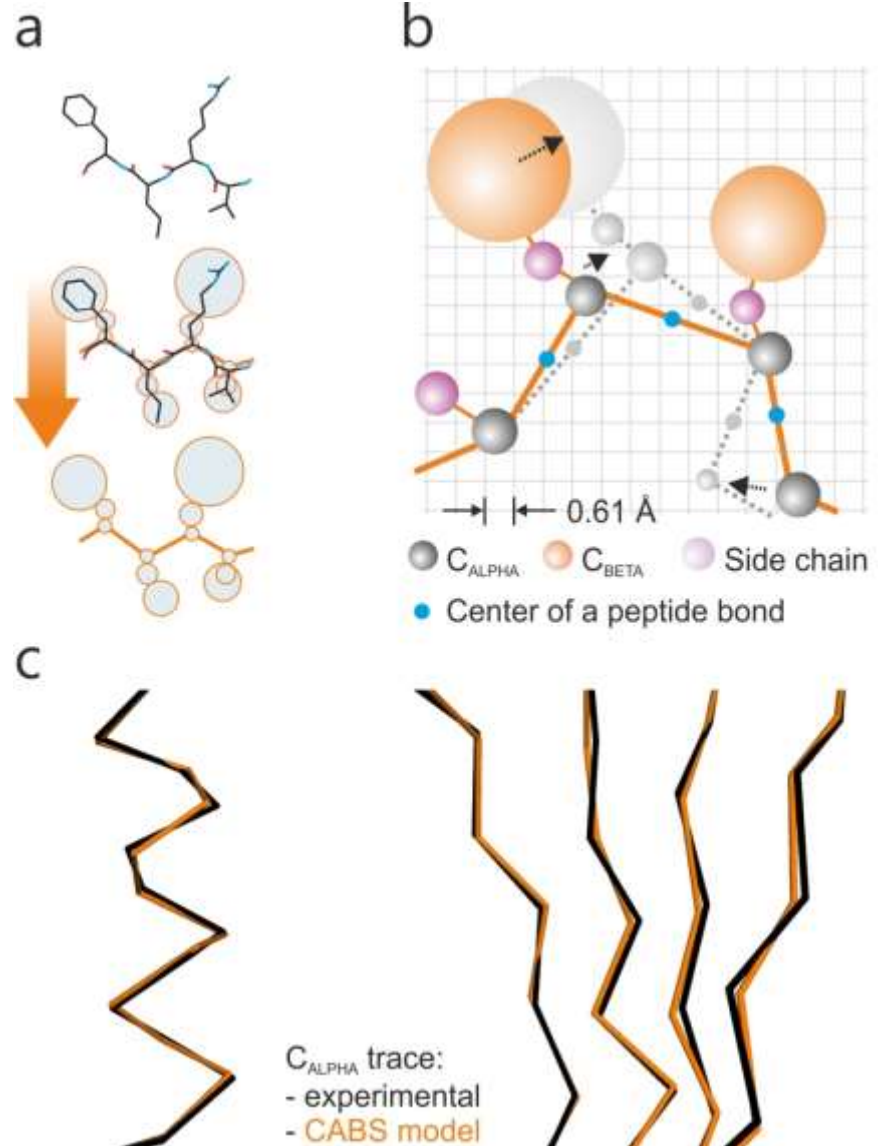
Monte Carlo dynamics

Figures:

a) Building reduced model

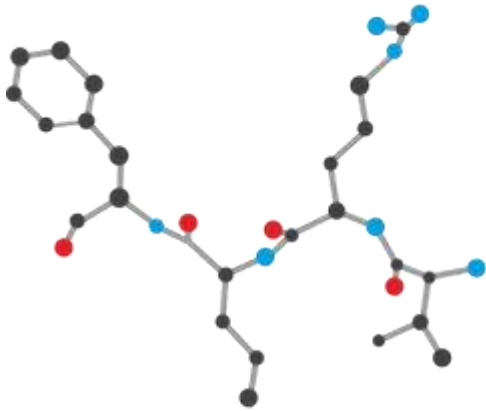
b) MC moves on the high-coordination lattice

c) Accuracy (C_{α} -traces)



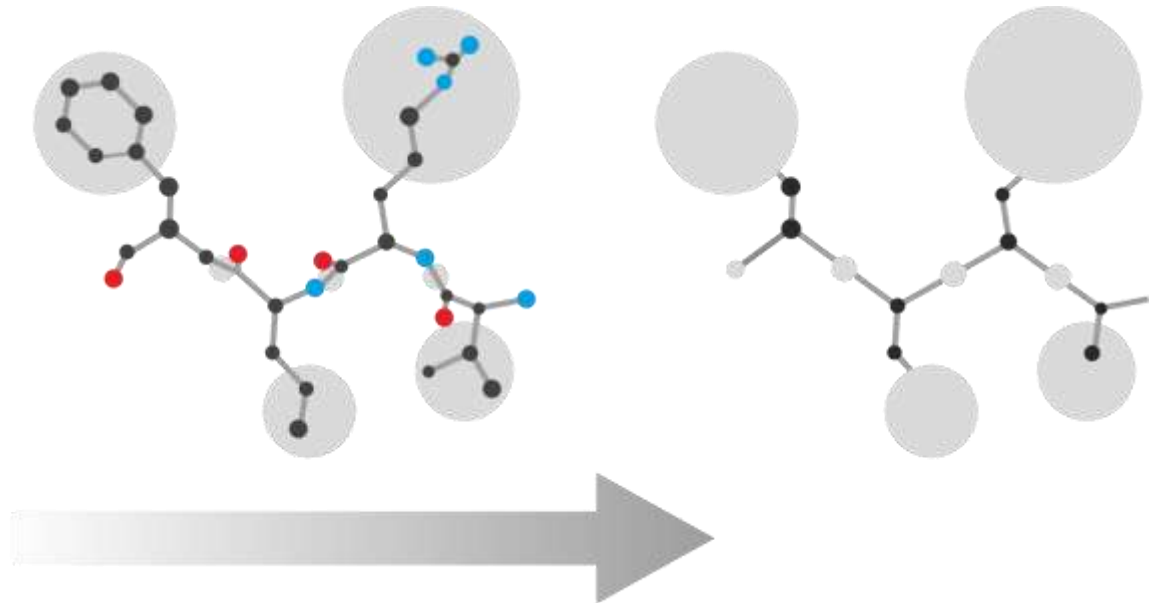
Time scales MD vs. CABS

All-atom molecular dynamics (MD)



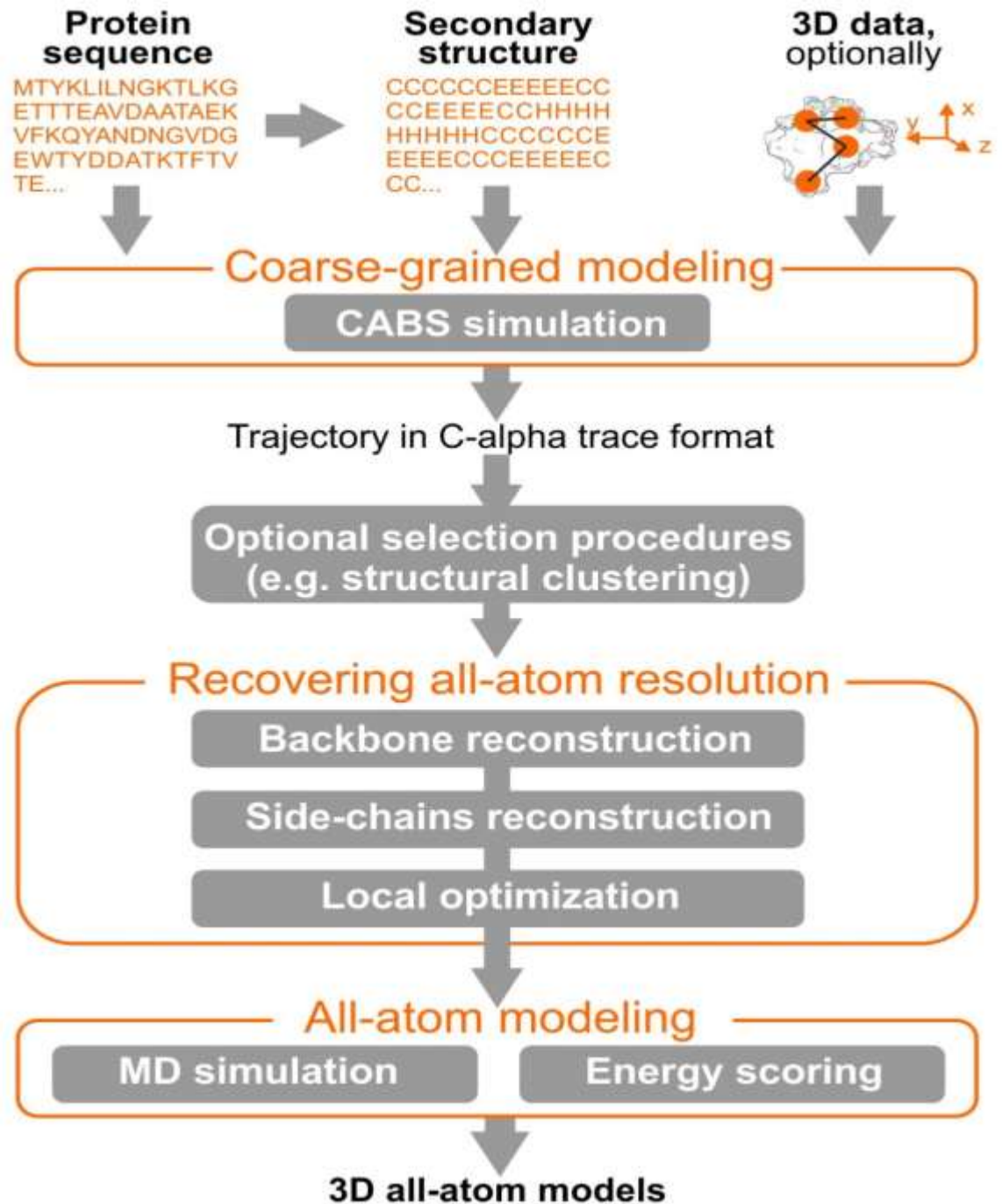
Max ~ 1 millisecond

CABS Monte Carlo dynamics



~ $10^3 / 10^4$ faster

Multiscale modeling with CABS



APPLICATINS

Structure prediction

Protein dynamics

Protein docking

CASP and CAPRI



CASP Competition

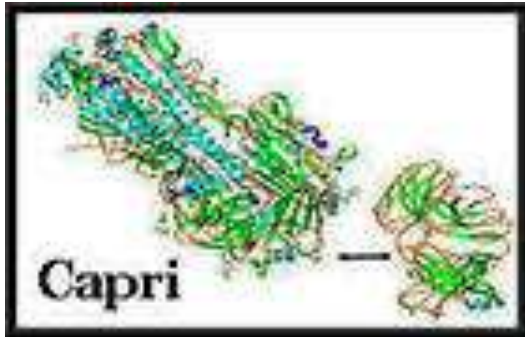
- CASP competition (Critical Assessment of Techniques for Protein Structure Prediction)
<http://predictioncenter.llnl.gov/>
- Their goal is to help advance the methods of identifying protein structure from sequence.

CASP Experiment

- Experimentalists are solicited to provide information about structures expected to be soon solved
- Predictors retrieve the sequence from prediction center (predictioncenter.llnl.gov)
- Deposit predictions throughout the season
- Meeting held to assess results

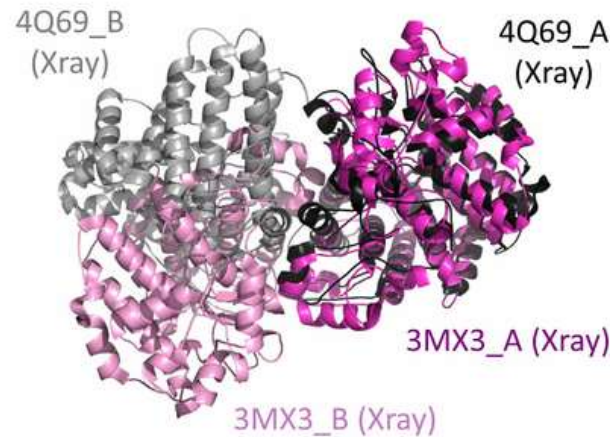
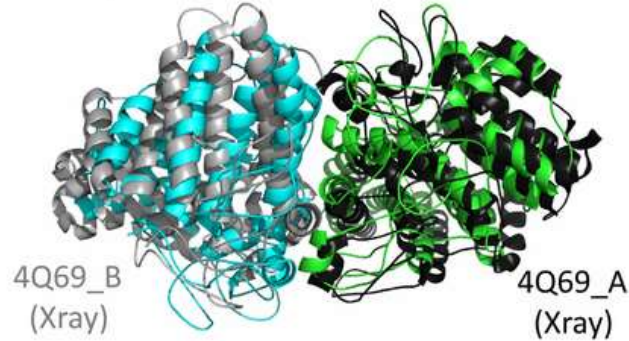
29

Polish scientists in CASP: Ginalski, Rychlewski, Bujnicki, Kolinski, Liwo, and others

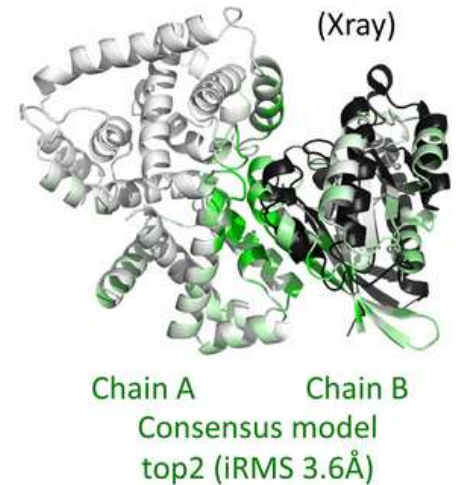


CAPRI: Critical Assessment of PRediction of Interactions

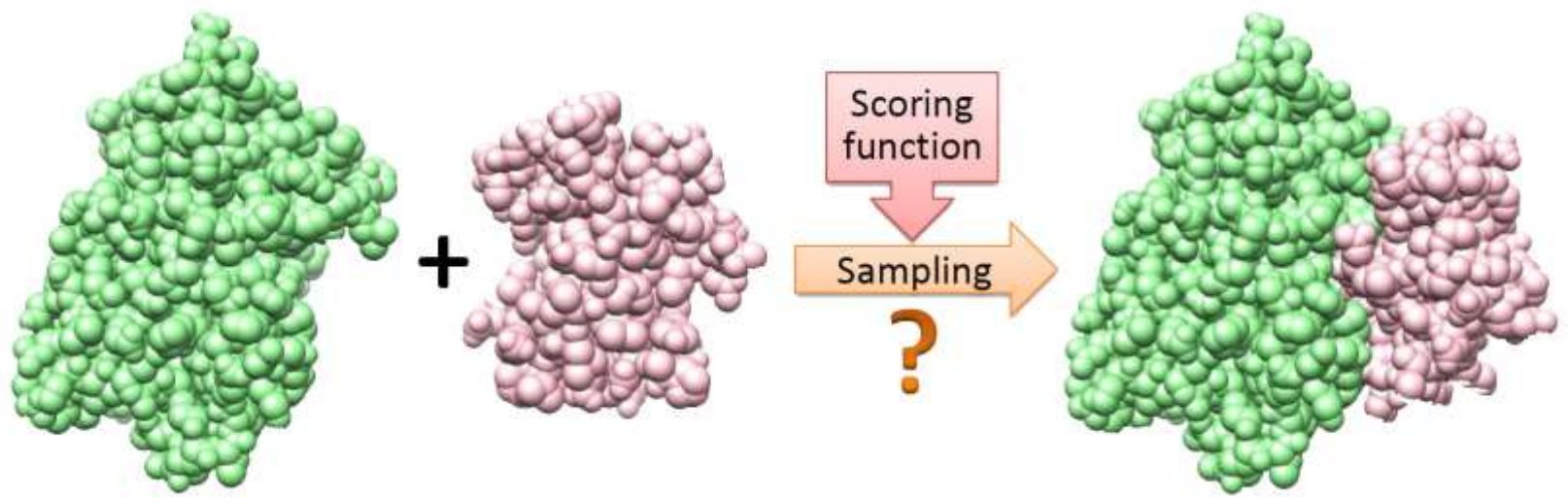
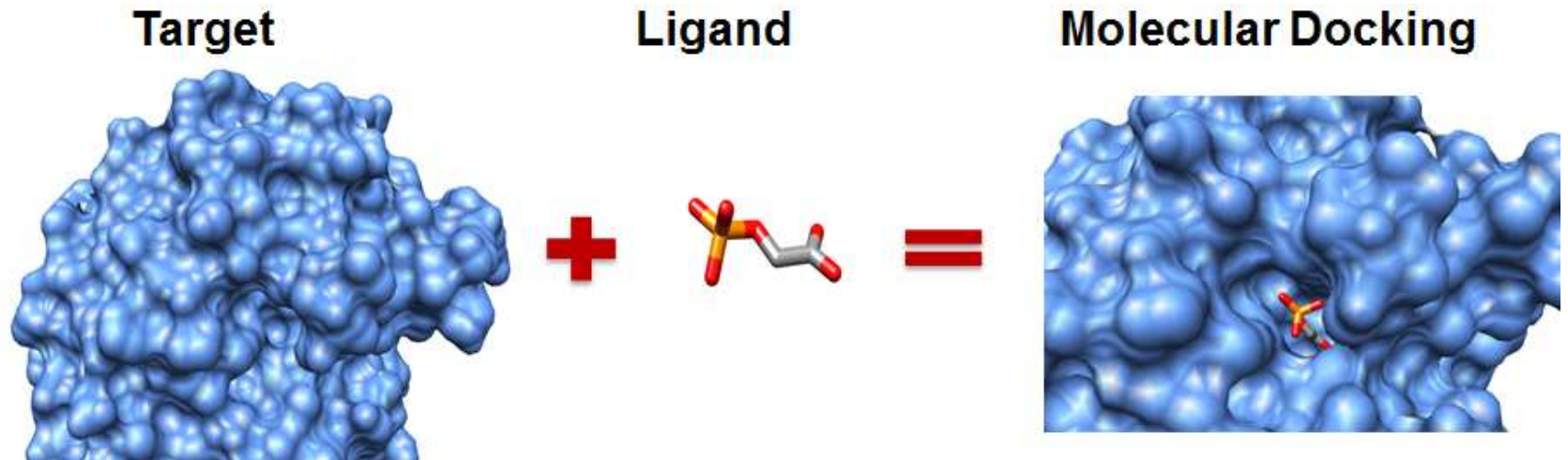
A CAPRI T72 Consensus Top5 Chain B CAPRI T72 Consensus Top5 Chain A



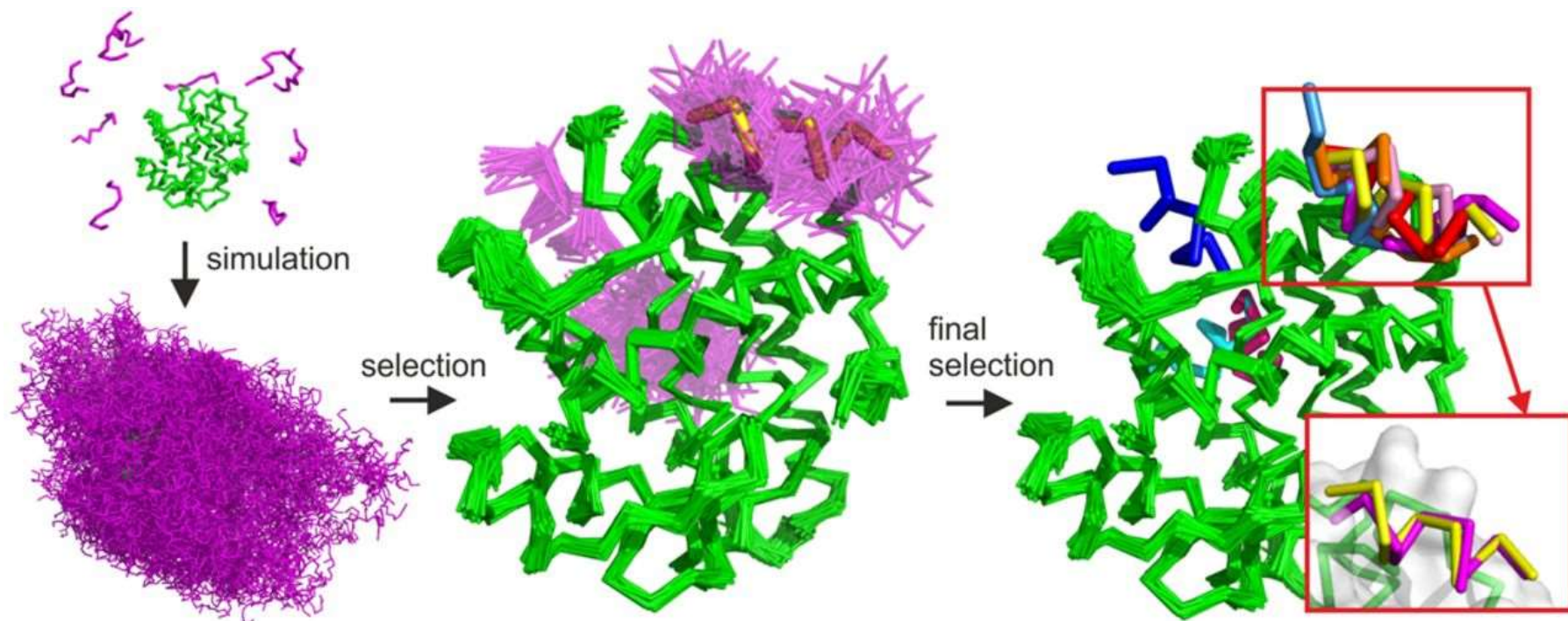
B 2G77 (Xray)



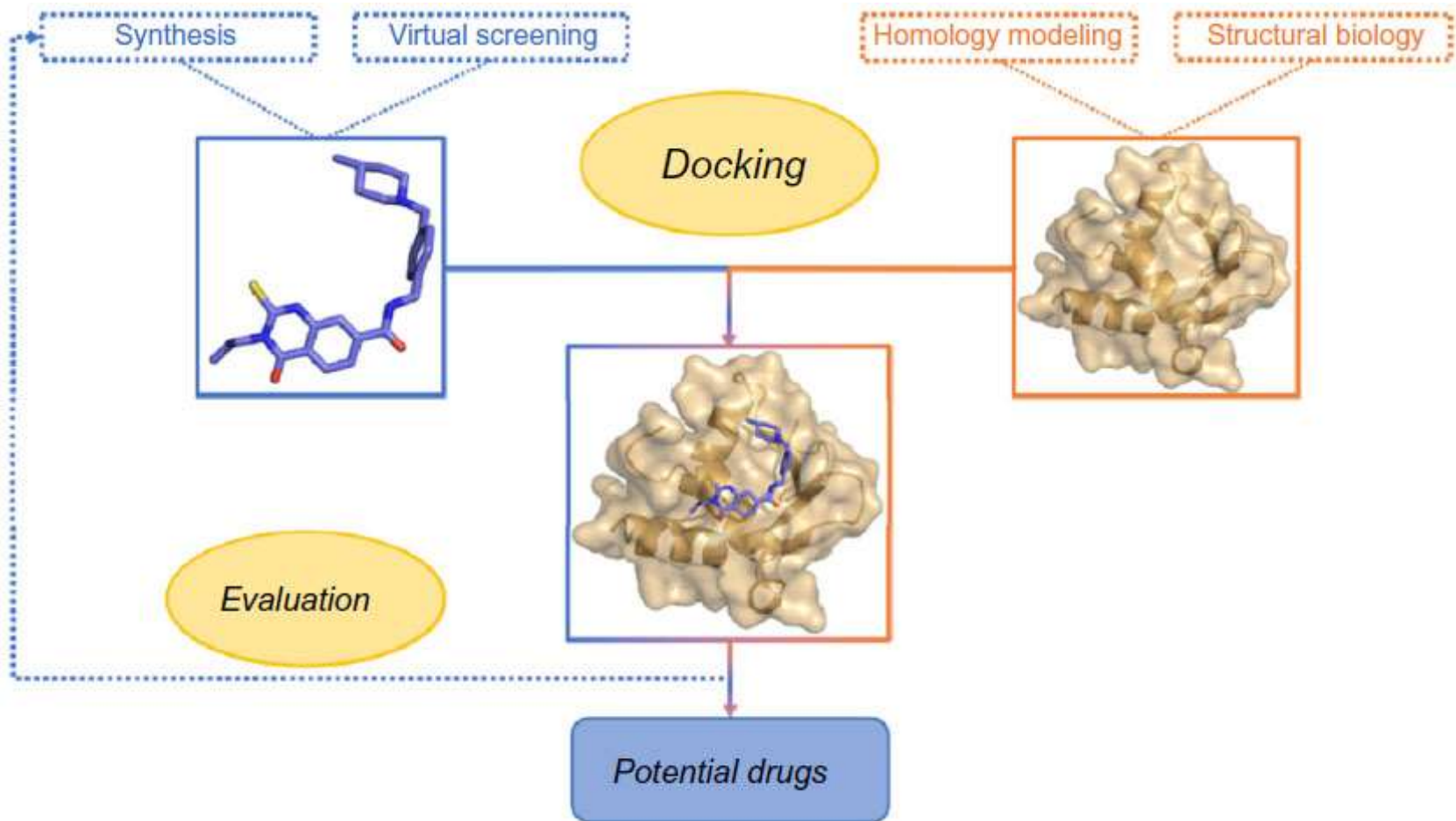
Molecular docking



Peptide docking with CABS model



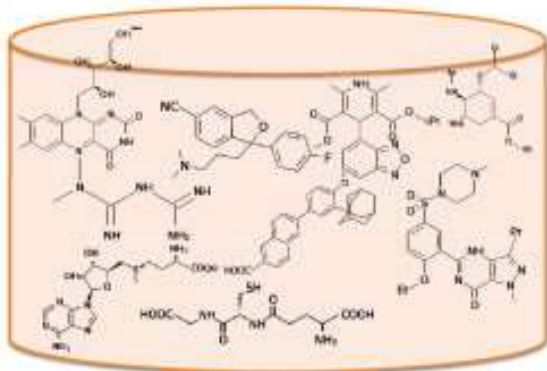
M. Kurcinski, M. Jamroz, M. Blaszczyk, A. Kolinski & S. Kmiecik, “CABS-dock web server for the flexible docking of peptides to proteins without prior knowledge of the binding site”, *Nucleic Acids Research*, 2015



(A) Docking



Protein of interest



Chemical database

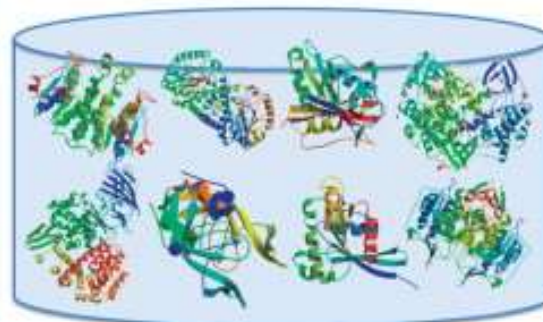


Possible binding ligand

(B) Reverse docking



Active compound or existing drug



Protein target database



Possible binding protein

Check our tools at:
<http://biocomp.chem.uw.edu.pl/tools>



LABORATORY
of THEORY of
BIOPOLYMERS

Andrzej Kolinski Research Group

RESEARCH

PEOPLE

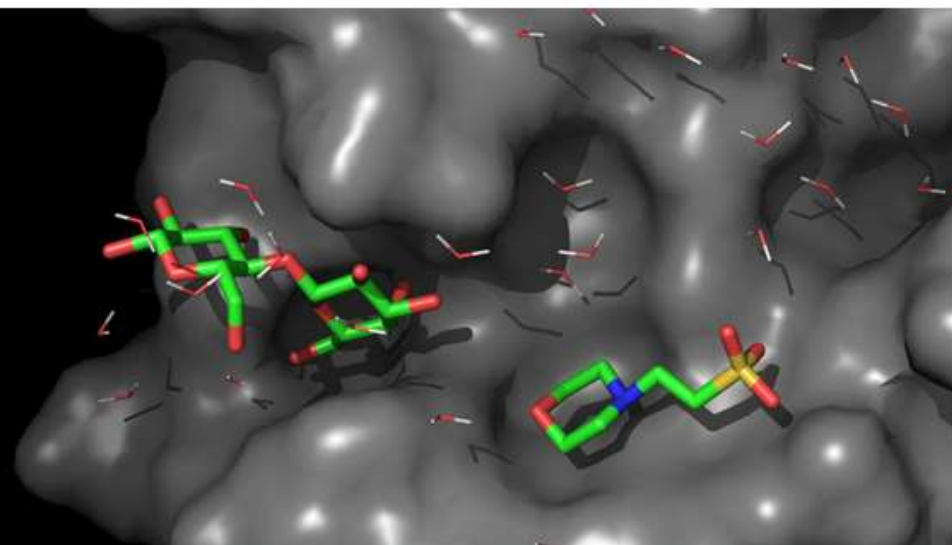
PUBLICATIONS

TOOLS

CONTACT

Modeling Software & Servers

[SEE OUR TOOLS](#)



NEWS

06.09.2013

We are pleased to announce an opening call for positions within the TEAM programme. We are looking for students to work in our project aimed at development of new modeling tools for structure and dynamics prediction of proteins and other biomolecules. >>>

PUBLICATIONS

CABS-flex: server for fast simulation of protein structure fluctuations

Authors: M. Jamroz, A. Kolinski, S. Kmiecik
Nucleic Acids Research, 41:W427-W431, 2013

ABSTRACT

NEWSLETTER

Subscribe to the LTB Newsletter - get news about our research and seminars!

E-mail

SUBSCRIBE

TOOLS

CABS-DOCK
CABS-FLEX
CABS-FOLD
AGGRESKAN3D
CABS
SURPASS
CABS-NMR
PYCABS
CLUSCO
BBQ
BIOSHELL
CCOMP
MSITE
BIODESIGNER
IMOL

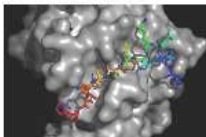
METHODS: The used methodologies were reviewed in our [review paper on Coarse-Grained Protein Models and Their Applications Tools](#) in *Chemical Reviews* journal.



FUNDING: CABS-dock, CABS-flex, CABS-fold and Aggreskan3D tools were funded by the Foundation for Polish Science TEAM project (TEAM/2011-7/5) co-financed by the European Regional Development Fund operated within the Innovative Economy Operational Program.

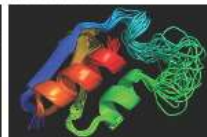
CABS-dock

server for protein-peptide docking and prediction of landing sites



CABS-flex

server for fast simulations of flexibility of protein structures



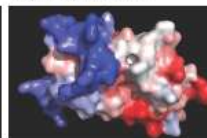
CABS-fold

server for de novo and consensus-based prediction of protein structure



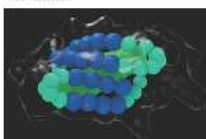
AGGRESKAN3D

server for prediction of aggregation properties of protein structures



SURPASS

SURPASS coarse-grained protein model of low-resolution



BBQ

program for protein backbone reconstruction from C-alpha coordinates



pyCABS

package for simulations of long-time protein dynamics using CABS reduced model



ClusCo

a software for GPU/CPU clustering and comparison of protein models

