# Algorithmic Foundations of Computational Biology
# CSCI 1820/2820
# An overview

- **Ch. 1 The BLAST Algorithm and Karlin-Altschul Statistics**
- **Ch. 2 Genome Assembly Algorithms and Haplotype Assembly Algorithms**
- **Ch. 3 The Protein Folding Problem: From HP-lattice models to AlphaFold**
- **Ch. 4 Recombination and Ancestral Recombination Graphs (ARGs) Algorithms**
- **Ch.5 Rigorous clustering: Spectral Graph Theory Algorithms**
- **Ch. 6 The Regulatory Genome and Gene Regulatory Networks**

# Ch. 1: BLAST Algorithm



**Given** a biomolecular query sequence Q
and a database DB of biomolecular sequences

**Find** all the biomolecular sequences in DB that have high alignment scores to the query

Biomolecular: DNA, RNA, protein

Problems we need to solve along the way:

Problem 1. General scoring schemes as hypotheses testing frameworks
          The Karlin-Altschul Statistics and the max scoring subsequence

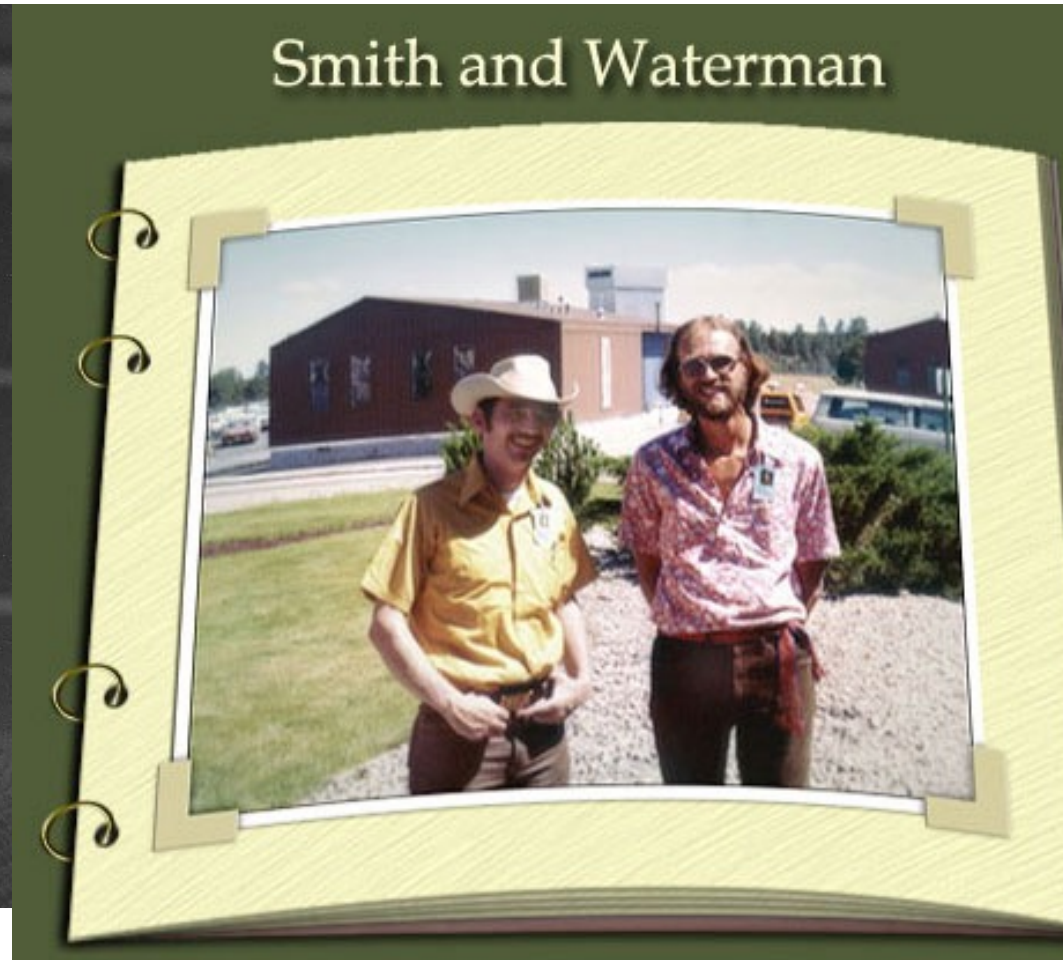Problem 2. Random Walks Theory and The Gambler's Ruin's Problem

Problem 3. De-noising: how long an alignment needs to be non-random?

Problem 4. Information Theory and the theory of scoring matrices for alignment

**Dr. Margaret Oakley Dayhoff**
**The Mother & Father of Bioinformatics**

**Temple Smith and Michael Waterman**
at Los Alamos, New Mexico
Photo by David Lipman, Taken Summer of 1980

# Sir Ronald Aylmer Fisher

The Lady Tasting Tea Problem

| | |
|---|---|
| | the<br>**Null Hypothesis** |
| Born | Ronald Aylmer Fisher<br><br>17 February 1890<br>London, England |
| Died | 29 July 1962 (aged 72)<br>Adelaide, South Australia, Australia |

•Linear discriminant analysis is a generalization of Fisher's li
discriminant[47][83]
•Fisher information, see also scoring algorithm also known a
scoring, and Minimum Fisher information, a variational prin
which, when applied with the proper constraints needed to
empirically known expectation values, determines the best
distribution that characterizes the system.[84]
•*F*-distribution, arises frequently as the null distribution of a
statistic, most notably in the analysis of variance
•Fisher–Tippett–Gnedenko theorem : Fisher's contribution t
made in 1927
•Fisher–Tippett distribution
•Fisher-Yates shuffle algorithm
•Von Mises–Fisher distribution[85]
•Inverse probability, a term Fisher used in 1922, referring to
fundamental paradox of inverse probability" as the source o
confusion between statistical terms which refer to the true v
estimated, with the actual value arrived at by estimation, w
subject to error.[86]
•Fisher's permutation test
•Fisher's inequality[87]
•Sufficient statistic, when a statistic is *sufficient* with respect
a statistical model and its associated unknown parameter if
statistic that can be calculated from the same sample provid
additional information as to the value of the parameter".[88]
•Fisher's noncentral hypergeometric distribution, a generaliz
the hypergeometric distribution, where sampling probabiliti
modified by weight factors.
•Student's *t*-distribution, widely used in statistics.[89][90]
•The concept of an ancillary statistic and the notion (the anc
principle) that one should condition on ancillary statistics.

# The BLAST algorithm **Professor Istrail**

☐ Detect all *word hits* (exact, or nearly identical matches) of a given length between the two sequences

- k=10 for nucleotide sequences (exact word matches)
- k=3 for protein sequences (nearly identical word matches)

☐ Extend the word hits in both directions to high-scoring *gap-free* segment pairs (HSPs)

- retain only HSPs that score above a threshold
- start from the center of the HSP (original BLAST, 1990), or from the center of a pair of HSPs located close to each other on the same diagonal (gapped BLAST, 1997)

☐ Extend the HSPs in both directions allowing for gaps

- use dynamic programming, and stop when the alignment score falls more than a threshold X below the best score yet seen

☐ Report all statistically significant local alignments

- E-value (starting with BLAST 2.0) is used to measure the statistical significance
- *E-value* = the number of alignments with score equal to or higher than *s* one would expect to find by chance when searching the database

# Ch. 2: Genome Assembly Algorithms

Questions: What algorithms to use to assemble DNA pieces into contigs and scafolds?

How long are the contigs?

How much the DNA target region is covered by the contigs?

How to measure the success of a genome assembly?

Problems we need to solve along the way

Problem 1. Genome Assembly Algorithms

Problem 2. Poisson statistics for DNA and Genome Assembly

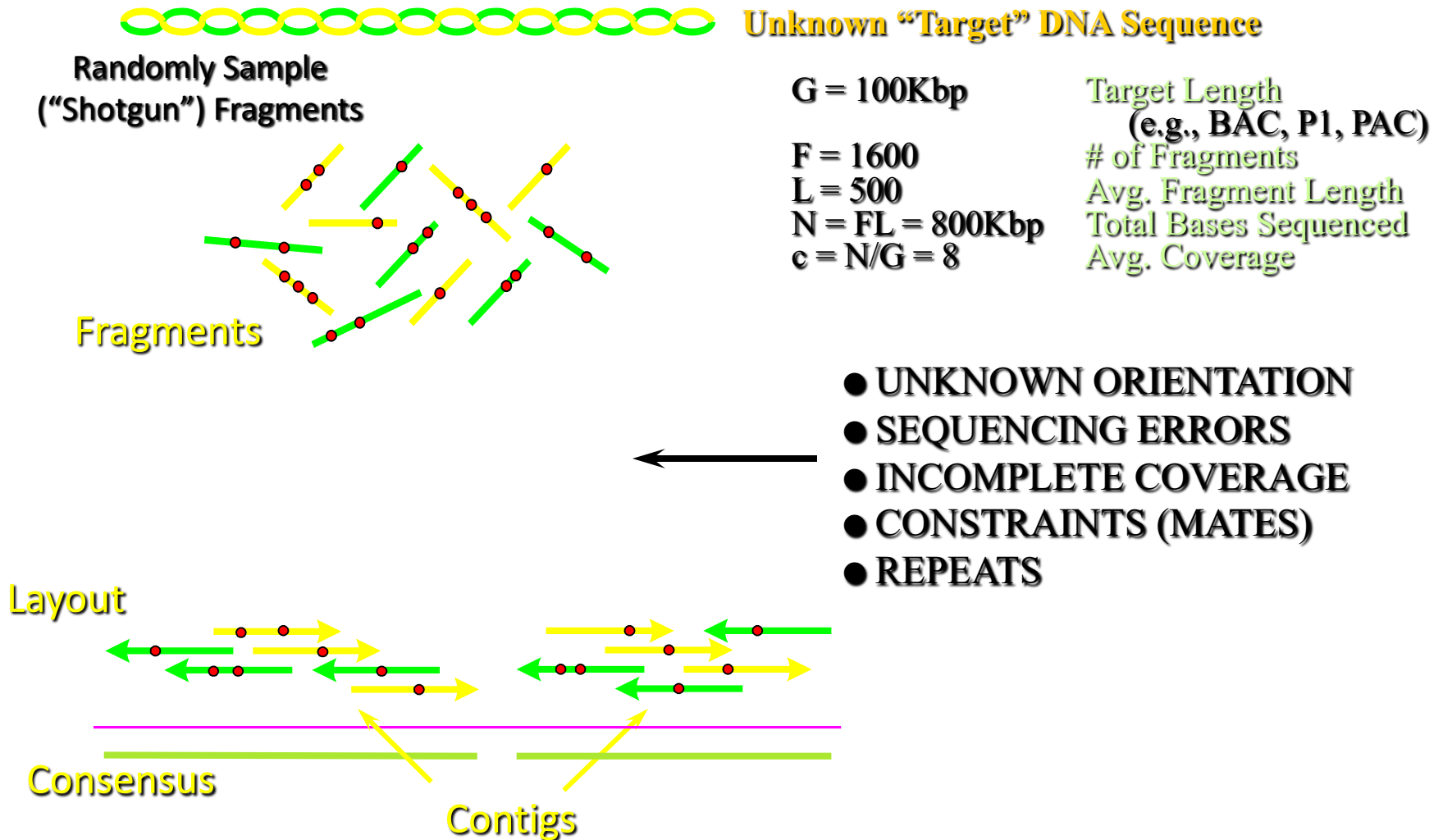Problem 3.  Ham Smith's DNA Lab with no windows
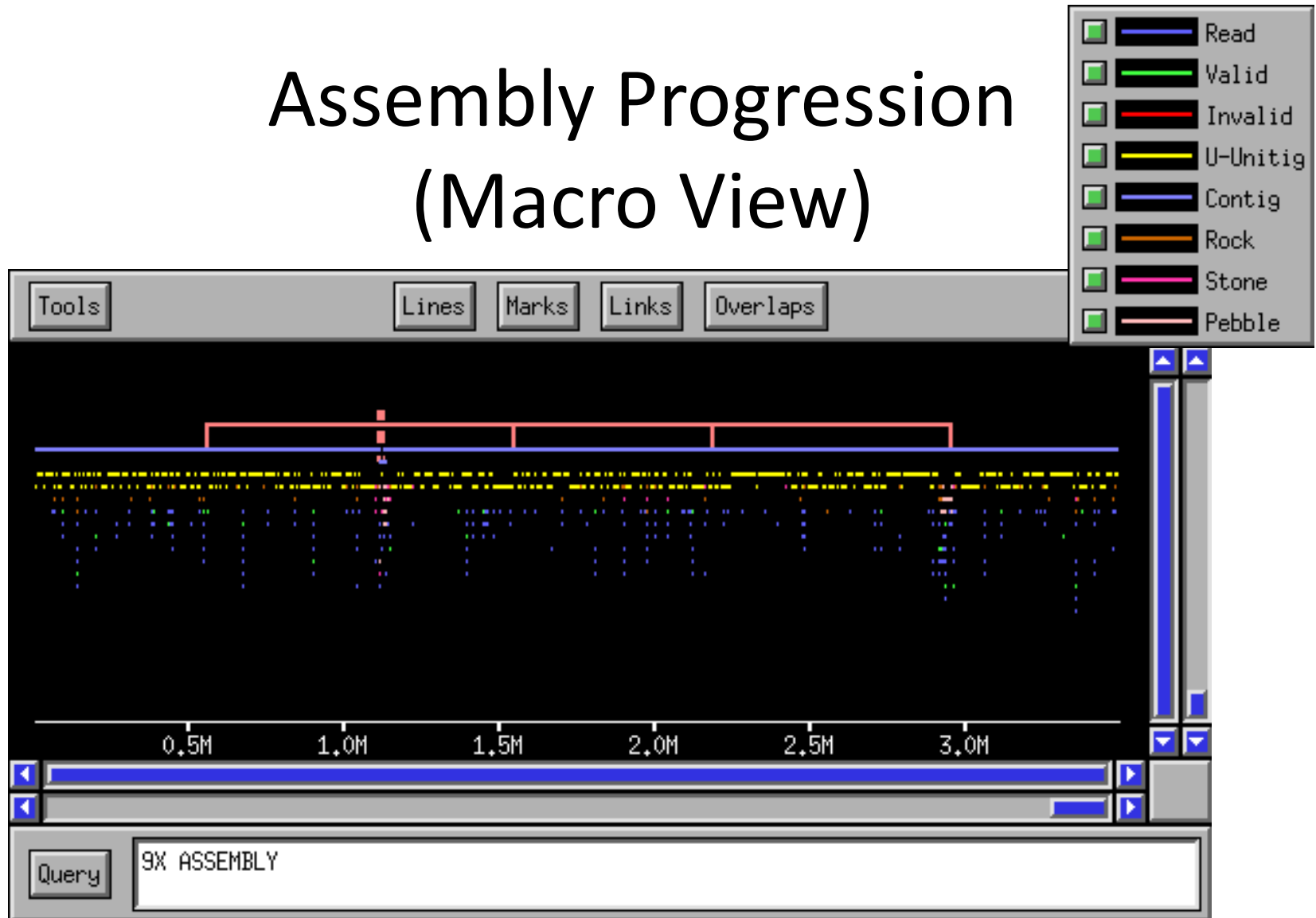
# Whole Genome Shotgun Sequencing

# Shotgun DNA Sequencing (Technology)

**DNA target sample**

**SHEAR**

**SIZE SELECT**

**End Reads (Mates)**

**Primer**

**SEQUENCE**

**LIGATE & CLONE**

**Vector**

# Shotgun DNA Sequencing (Computation)



Unknown "Target" DNA Sequence

**Randomly Sample ("Shotgun") Fragments**

$G = 100Kbp$ — Target Length (e.g., BAC, P1, PAC)

$F = 1600$ — # of Fragments
$L = 500$ — Avg. Fragment Length
$N = FL = 800Kbp$ — Total Bases Sequenced
$c = N/G = 8$ — Avg. Coverage

Fragments

- UNKNOWN ORIENTATION
- SEQUENCING ERRORS
- INCOMPLETE COVERAGE
- CONSTRAINTS (MATES)
- REPEATS

Layout

Consensus

Contigs

# Assembly Progression (Macro View)

# Siméon Denis Poisson

| | |
|---|---|
| Born | 21 June 1781<br>Pithiviers, Kingdom of France<br>(present-day Loiret) |
| Died | 25 April 1840 (aged 58)<br>Sceaux, Hauts-de-Seine, Kingdom of France |
| Alma mater | École Polytechnique |
| Known for | Poisson process<br>Poisson equation<br>Poisson kernel<br>Poisson distribution<br>Poisson limit theorem<br>Poisson bracket<br>Poisson algebra<br>Poisson regression<br>Poisson summation formula<br>Poisson's spot<br>Poisson's ratio<br>Poisson zeros<br>Conway–Maxwell–Poisson distribution<br>Euler–Poisson–Darboux equation |
| **Scientific career** | |
| Fields | Mathematics and physics |
| Institutions | École Polytechnique<br>Bureau des Longitudes<br>Faculté des sciences de Paris<br>École de Saint-Cyr |
| Academic advisors | Joseph-Louis Lagrange<br>Pierre-Simon Laplace |
| Doctoral students | Michel Chasles<br>Joseph Liouville |
| Other notable students | Nicolas Léonard Sadi Carnot<br>Peter Gustav Lejeune Dirichlet |

POISSON.

# de Bruijn Genome Assembly

# Ch. 3 Recombination and Ancestral Recombination Graphs (ARG) Algorithms



How do we reconstruct genealogies of a sample of individuals incorporating past mutations and recombinations?

Recombination + Phylogenetic Trees = ARG

# Ancestral Recombination Graph and Marginal Trees

# Ch. 5 Rigorous Clustering Algorithms Spectral Graph Theory Algorithms

**Algorithms and Statistical Theory**

- An introduction to Linear Algebra foundations for graph theory

- Principles of  Clustering Theory

- Graph Laplacians

- Graph cuts and random walks intuitions for Spectral Clustering

- Unnormalized Spectral Clustering Algorithms

- Normalized Spectral Clustering Algorithms

- Algorithmic Fairness and Clustering

# Pierre-Simon Laplace



| Born | 23 March 1749 Beaumont-en-Auge, Normandy, Kingdom of France |
|---|---|
| Died | 5 March 1827 (aged 77) Paris, Kingdom of France |
| Alma mater | University of Caen |
| Known for | show |
| **Scientific career** ||
| Fields | Astronomy and Mathematics |

| Notable students | Siméon Denis Poisson Napoleon Bonaparte |
|---|---|

Keep You From Forgetting To Mail Your Wife's Letter  RUBE GOLDBERG (tm)  RGI 049

Mixed character of the problem :

continuous  mathematics
discrete      mathematics

GENOMIC
REGULATORY
SYSTEMS

# A Tale of Two Networks



## Sea Urchin



## Drosophila

# One gene, 30 years of study, 300 docs and postdocs
# A Proposal for Nobel Prize
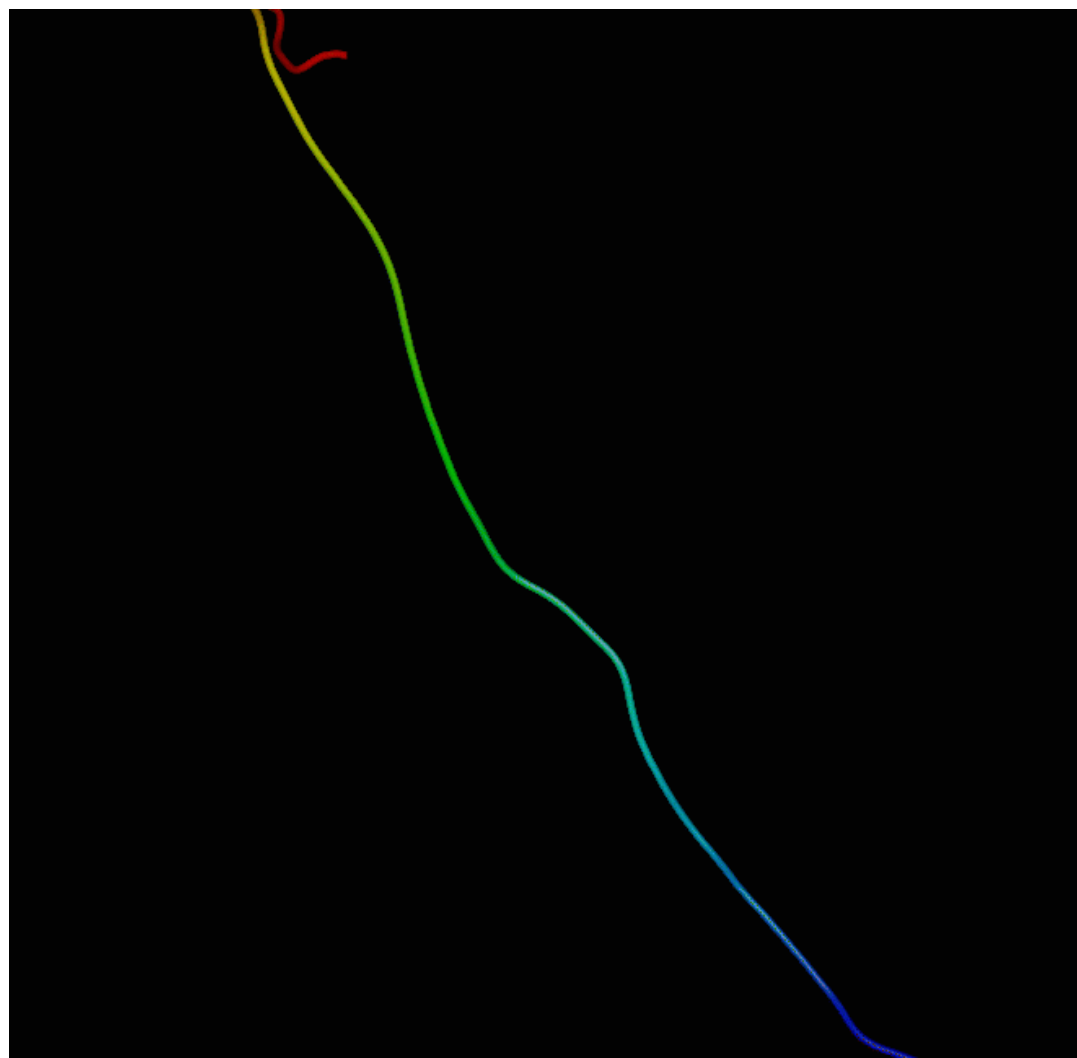


"Programs built into the DNA of every animal."
Eric H. Davidson

**Genomic Regulatory Systems**

Prof. Istrail

# The Protein Folding Problem

## Statistical Mechanics models

Mixed character of the problem :

continuous   mathematics  --  geometry of surfaces &
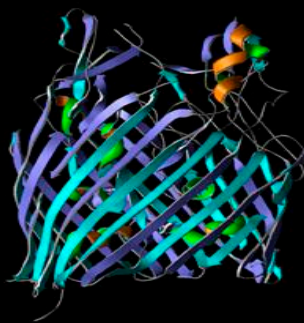discrete       mathematics  --  combinatorics of folds

*"The protein folding problem is three different problems: the folding code – the thermodynamic question of how a native structure results from the interatomic forces acting on an amino acid sequence; protein structure prediction – the computational problem of how to predict the native structure of a protein from its amino acid; and the folding speed (Levinthal's paradox) – the kinetic question of how a protein can fold so fast… Current knowledge of the folding code is sufficient to guide the successful design of new proteins and new materials. Current computer algorithms are now predicting the native structures of small simple proteins remarkable accurately, contributing to drug discovery and proteomics. Even once intractable Levinthal 'spuzzle now seems to have a very simple answer…"*

**Ken Dill**

K. A. Dill, S. Banu Ozkan, T. R. Weikl, J. D. Chodera and V. A. Voelz. The protein folding problem: when will it be solved? *Current Opinion in Structural Biology*, 17:2342-346, 2007.
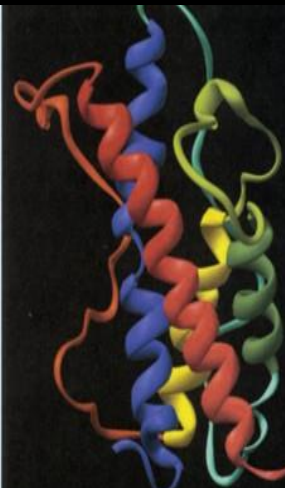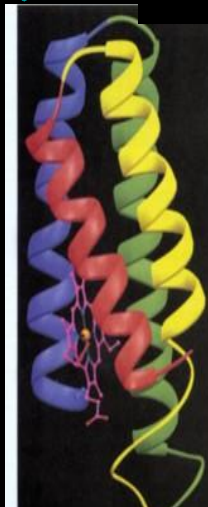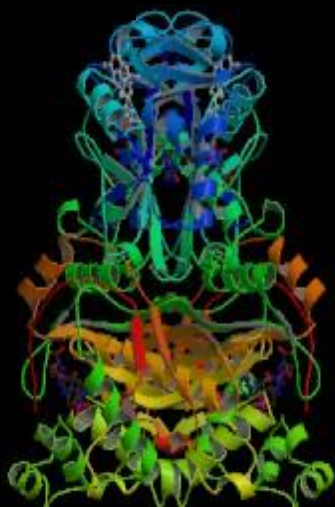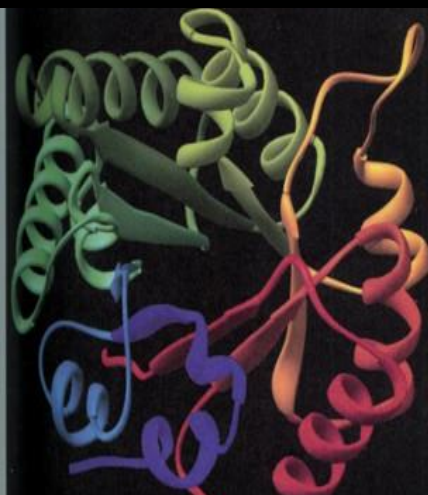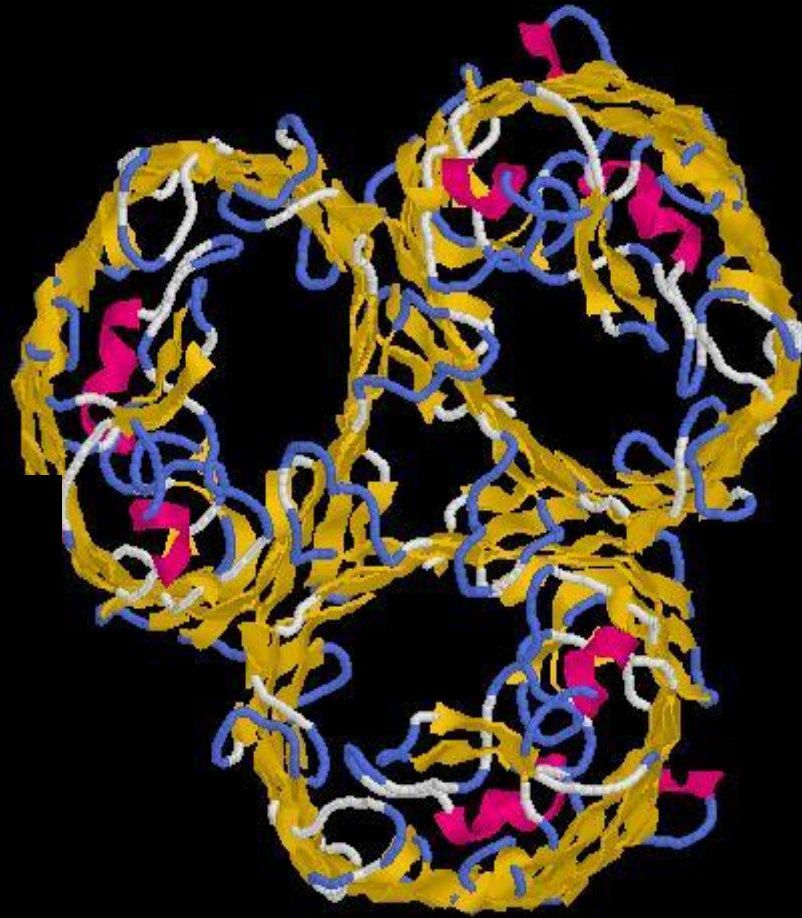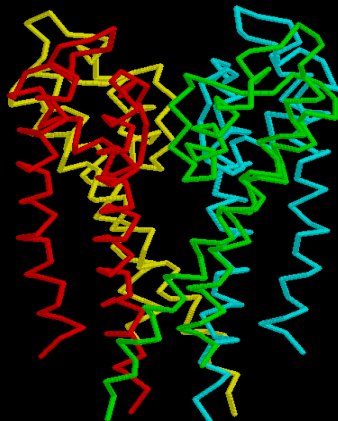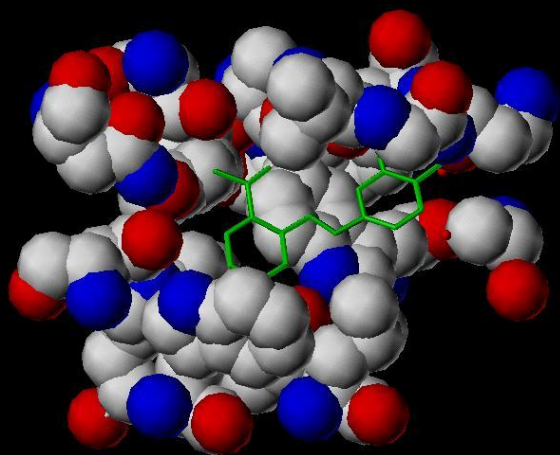
FhuA
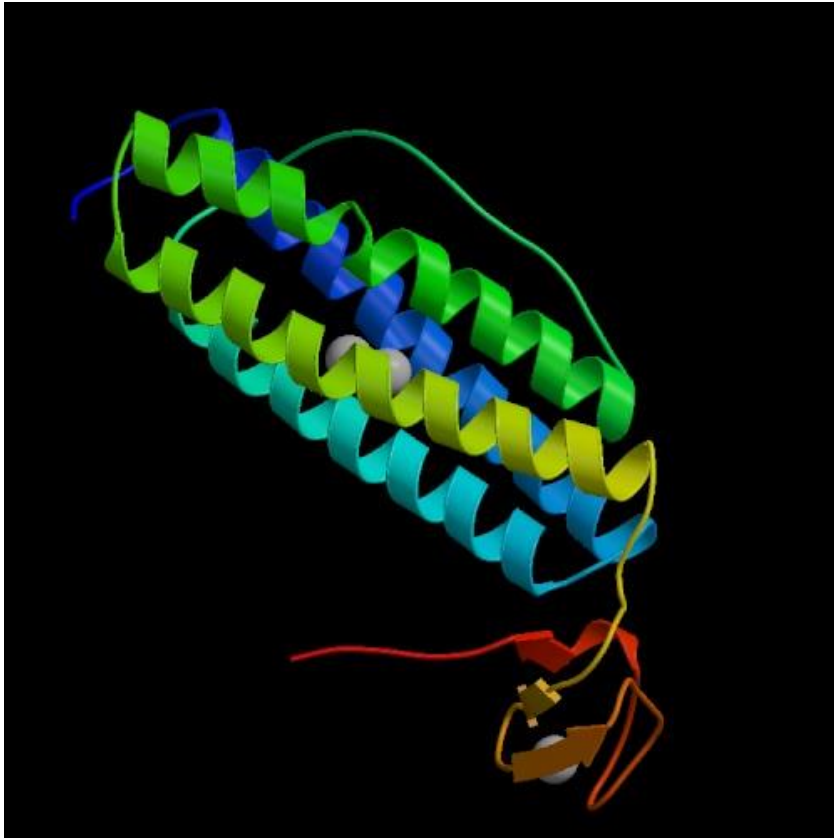(Ferguson et al., 1998)

FepA
(Buchanan et al., 1999)

OmpA
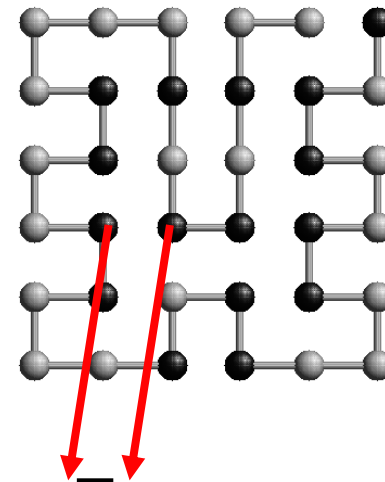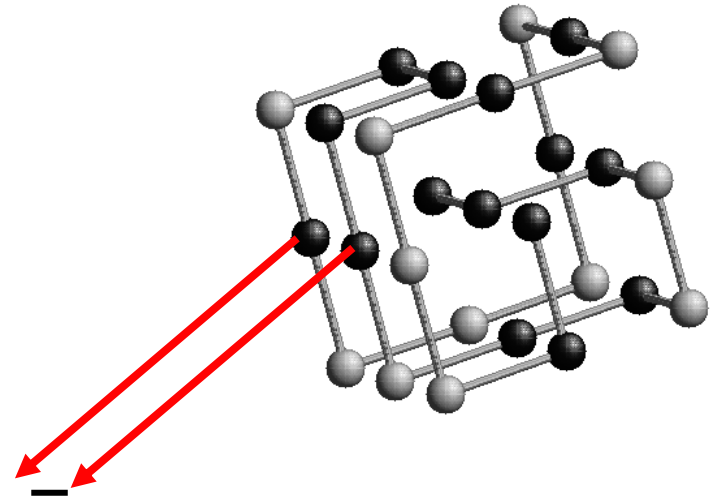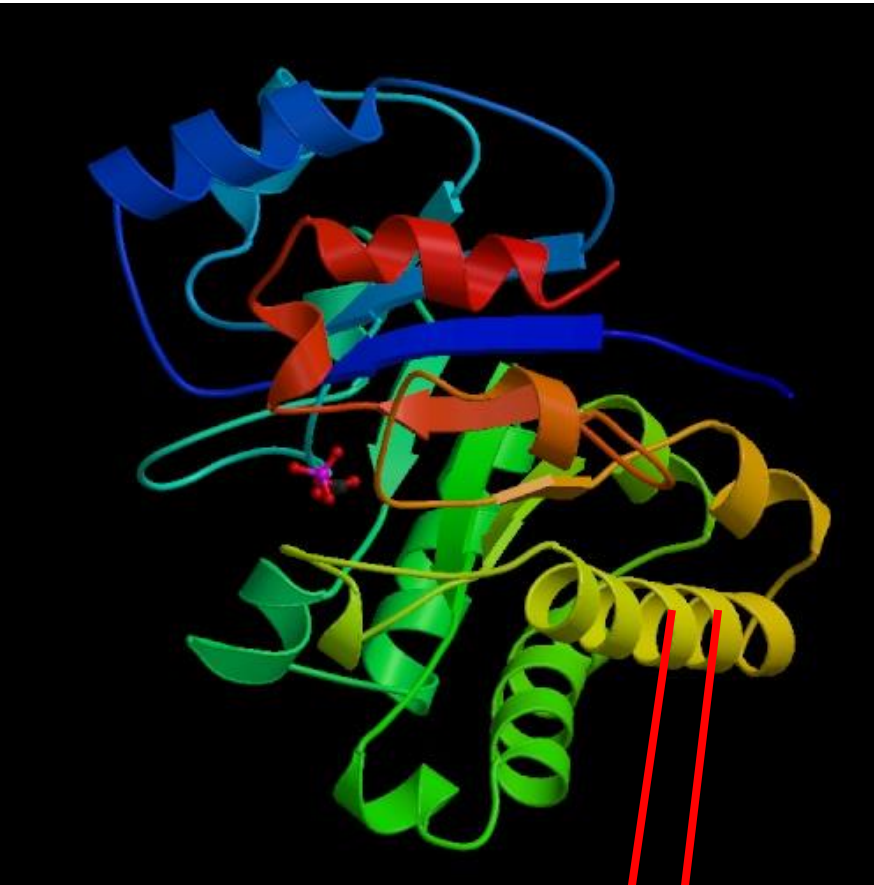(Pautsch & Schulz, 1998)

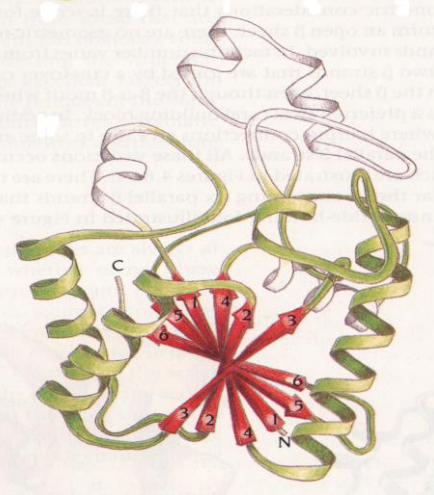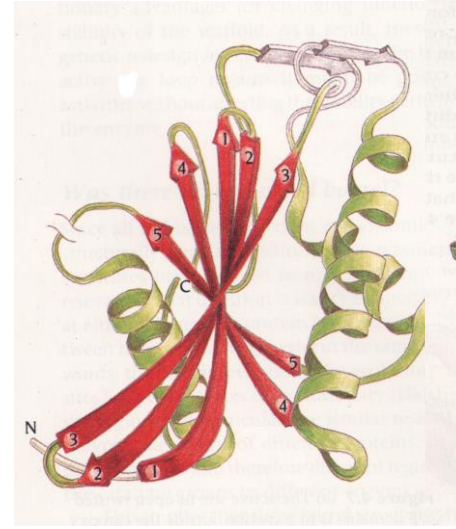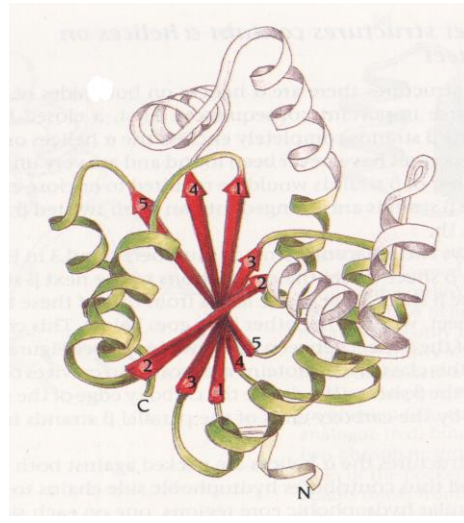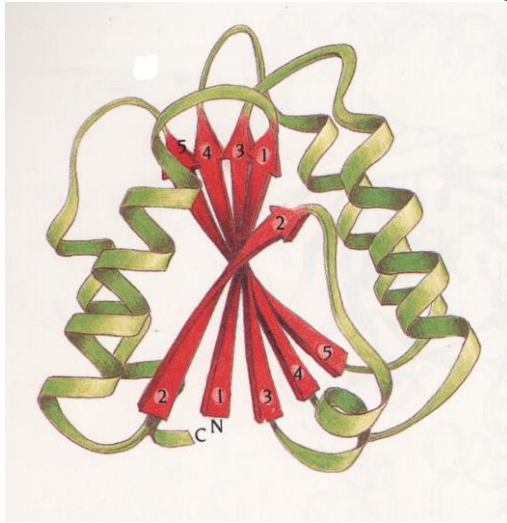Illustration © 1999 JHK

# SELF-AVOIDING WALK

# CONTACT



<= 4 Amstrongs

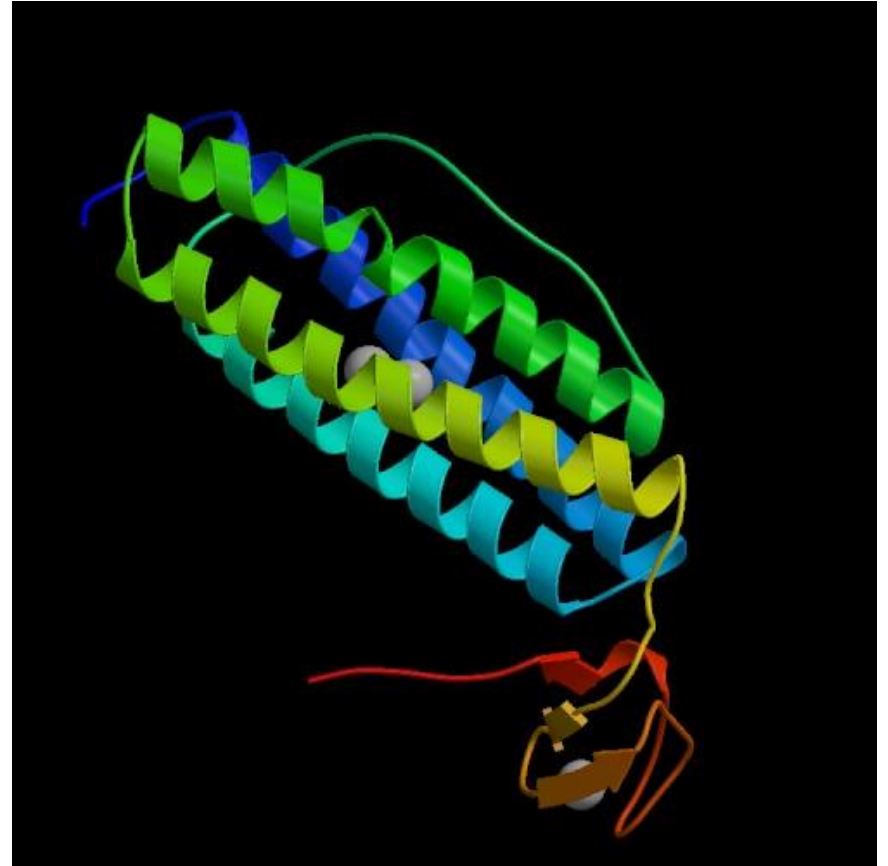# ASILOMAR: Protein Folding -- A Paradigm-Shift



**Structure Similarity**

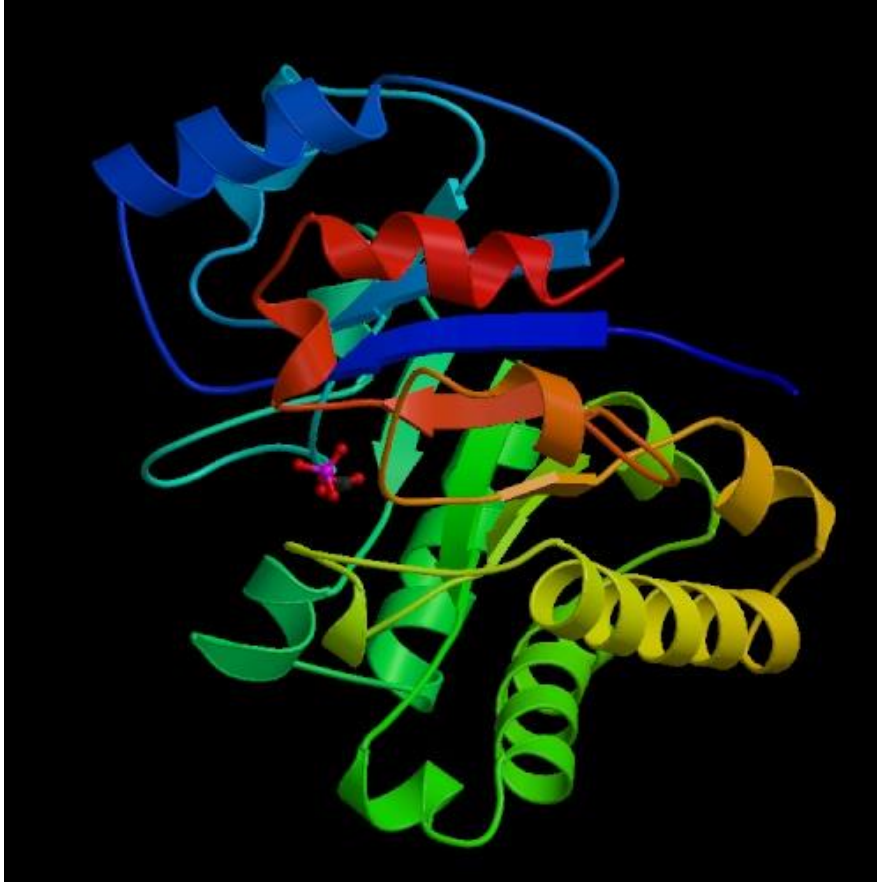**Structure Alignment**

**Fold Recognition**

**Fold Assignment**

# FOLD RECOGNITION

# Measure of Protein Structure Similarity

- Root-mean-square distance (RMSD)

- Difference of the distance matrices (DDM)

- Contact Map Overlap (CMO)

- Various more or less ad hoc scoring schemes based on local secondary structure, hydrogen bonding pattern, burial status, interaction pattern
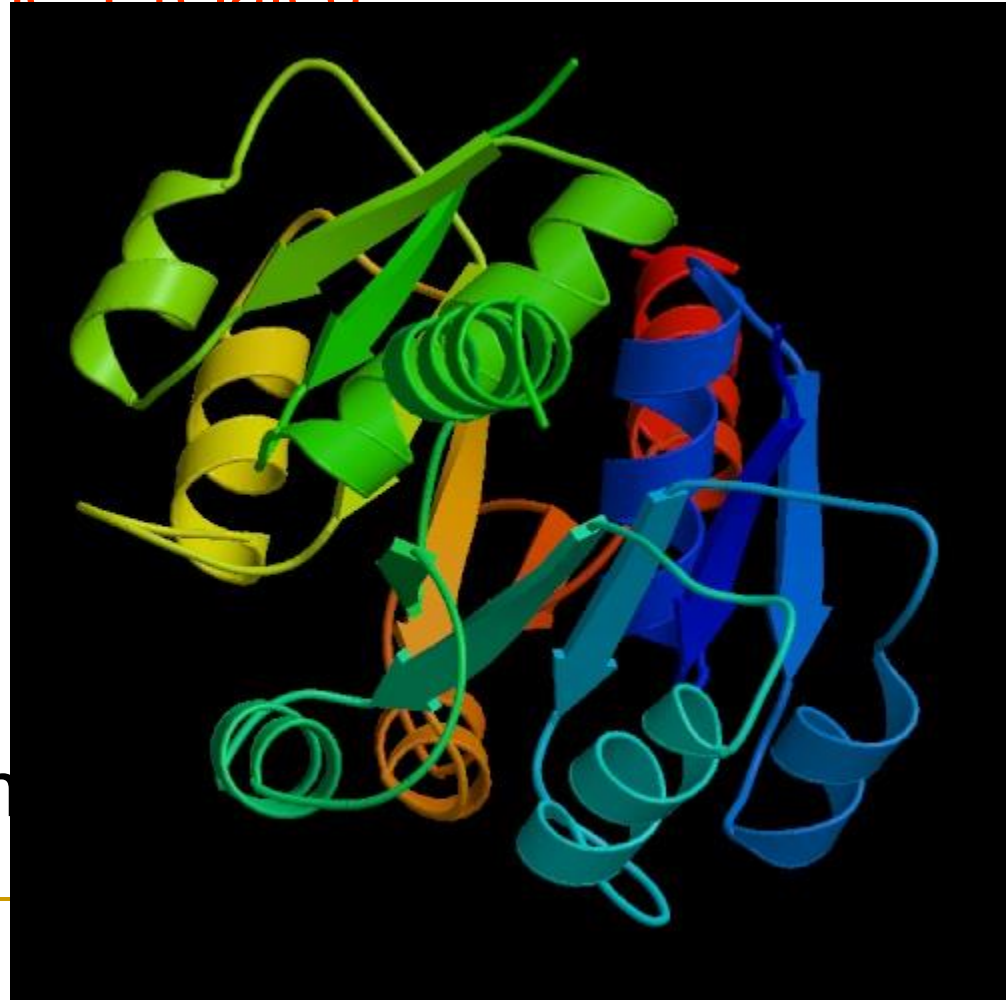
{1RCD} IRON STORAGE


{1FHA} IRON STORAGE

Objinserted 1PSP into 1FHA with 87 shared contacts

87 shared contacts

Save PostScript

# Skolnick Test

- Four Families
  1. Flavodoxin-like fold Che-Y related
  2. Plastocyanin
  3. TIM Barrel
  4. Ferratin

- alpha-beta
- 8 structures
- up to 124 residues
- 15-30% sequence sir
- < 3Å RMSD

# Skolnick Test

- Four Families
  1. Flavodoxin-like fold Che-Y related
  2. Plastocyanin
  3. TIM Barrel
  4. Ferratin
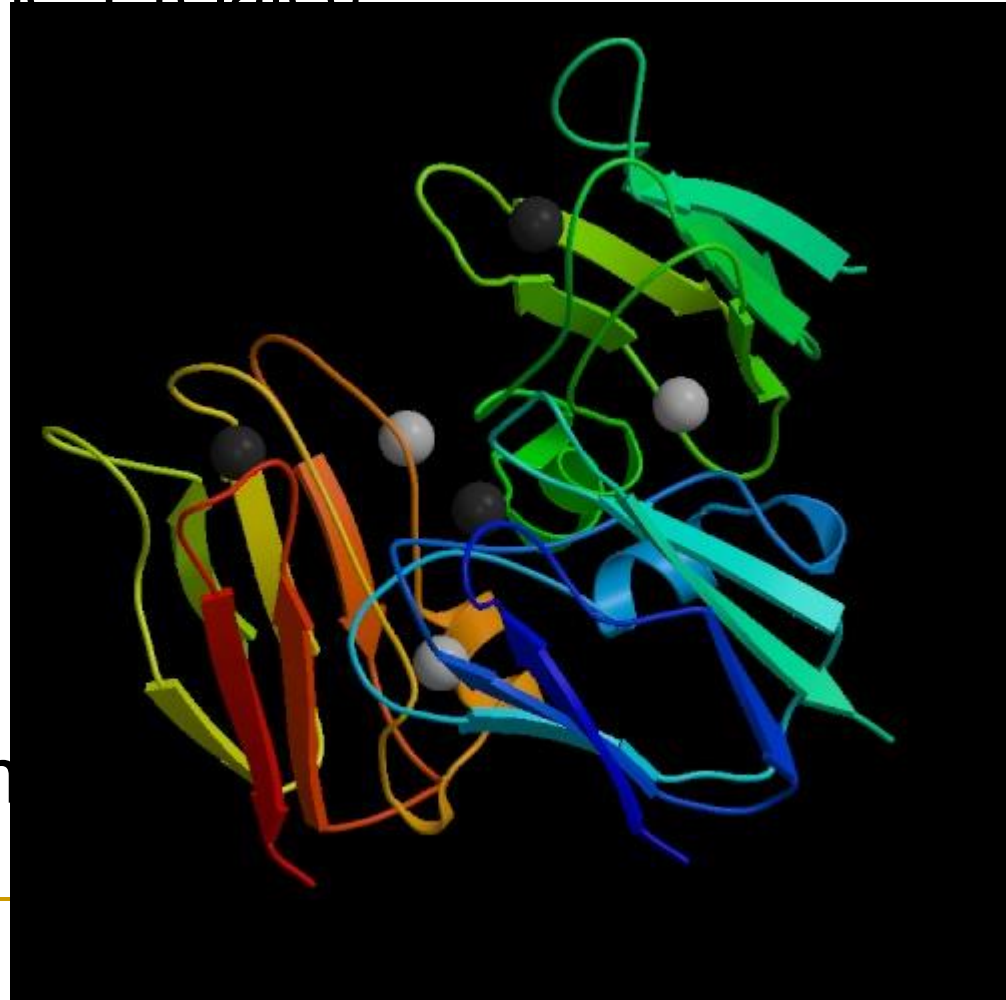
- beta
- 8 structures
- up to 99 residues
- 35-90% sequence sin
- < 2Å RMSD

# Skolnick Test

- ## Four Families
  1. Flavodoxin-like fold Che-Y related
  2. Plastocyanin
  3. TIM Barrel
  4. Ferratin

- ## alpha-beta
- ## 11 structures
- ## up to 250 residues
- ## 30-90% sequence sin
- ## < 2Å RMSD

# Skolnick Test

- Four Families
  1. Flavodoxin-like fold Che-Y related
  2. Plastocyanin
  3. TIM Barrel
  4. Ferratin

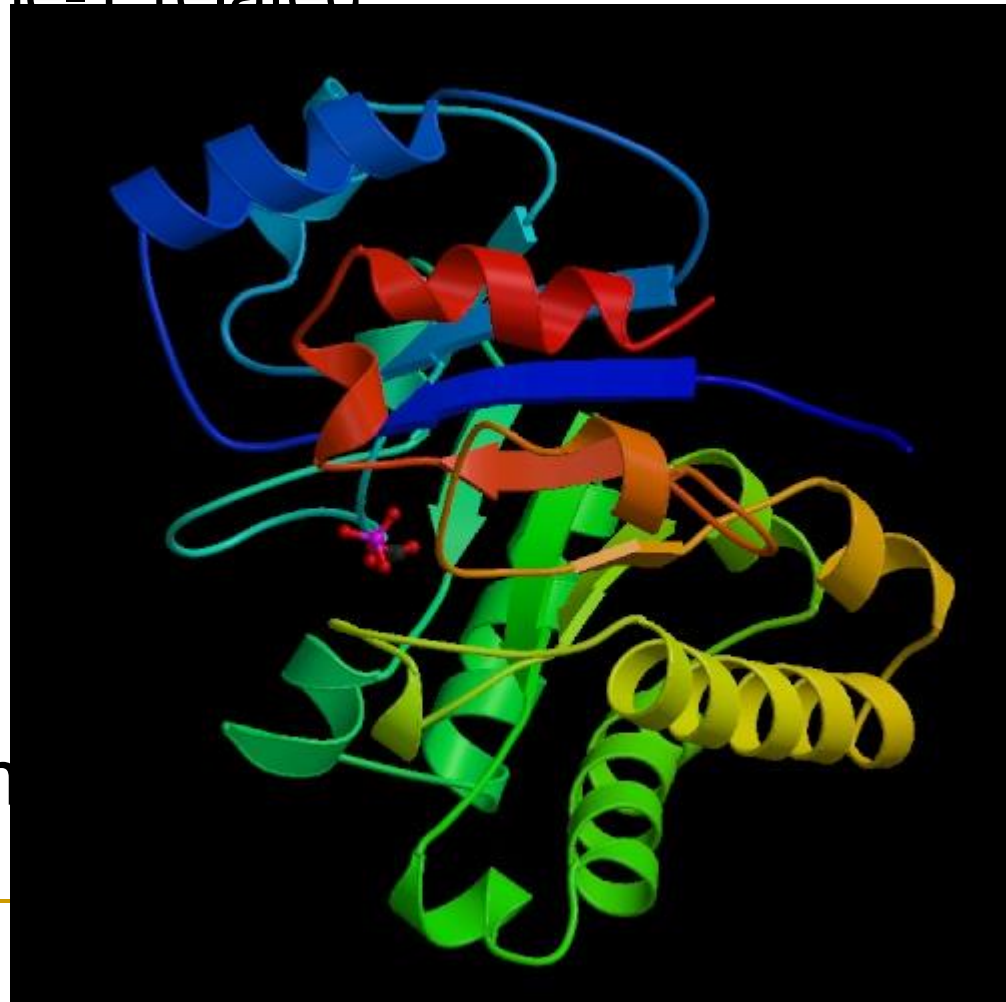- alpha

- 6 structures

- up to 170 residues

- 7-70% sequence simi

- < 4Å RMSD