# 3D Structure
## *Prediction and Assessment*



triose phosphate isomerase barrel



orthogonal views of Rop



Orthogonal beta sandwich fold of intestinal fatty acid–binding protein
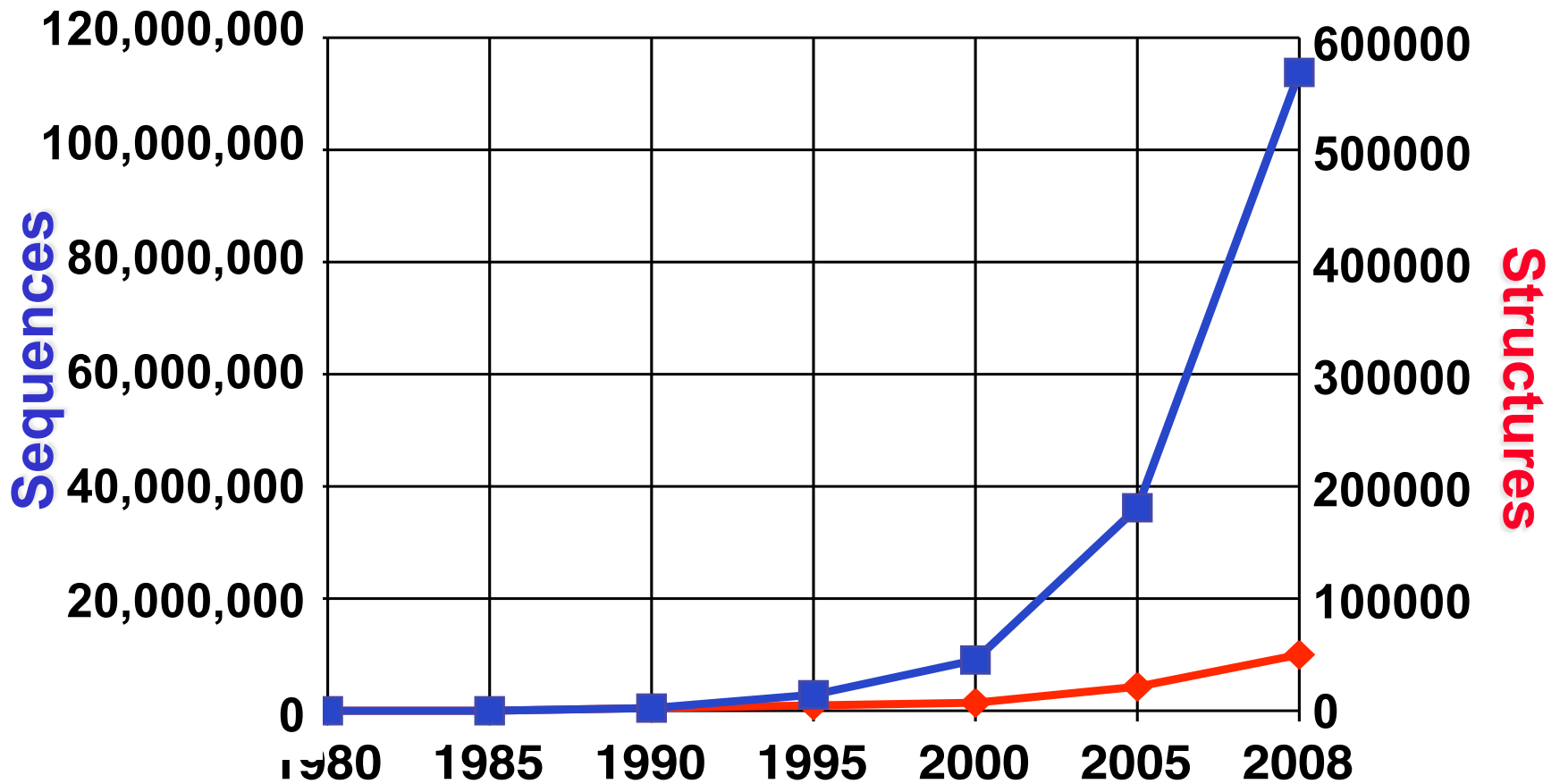
**David Wishart**

**Athabasca 3-41**

**david.wishart@ualberta.ca**

# Outline & Objectives*

- **Become familiar with the Protein Universe and the Protein Structure Initiative**

- **Learn principles of how to do homology (comparative) modelling of 3D protein structures**

- **Learn how to do homology modelling on the Web**

- **Learn how to assess 3D structures (modelled and experimental)**

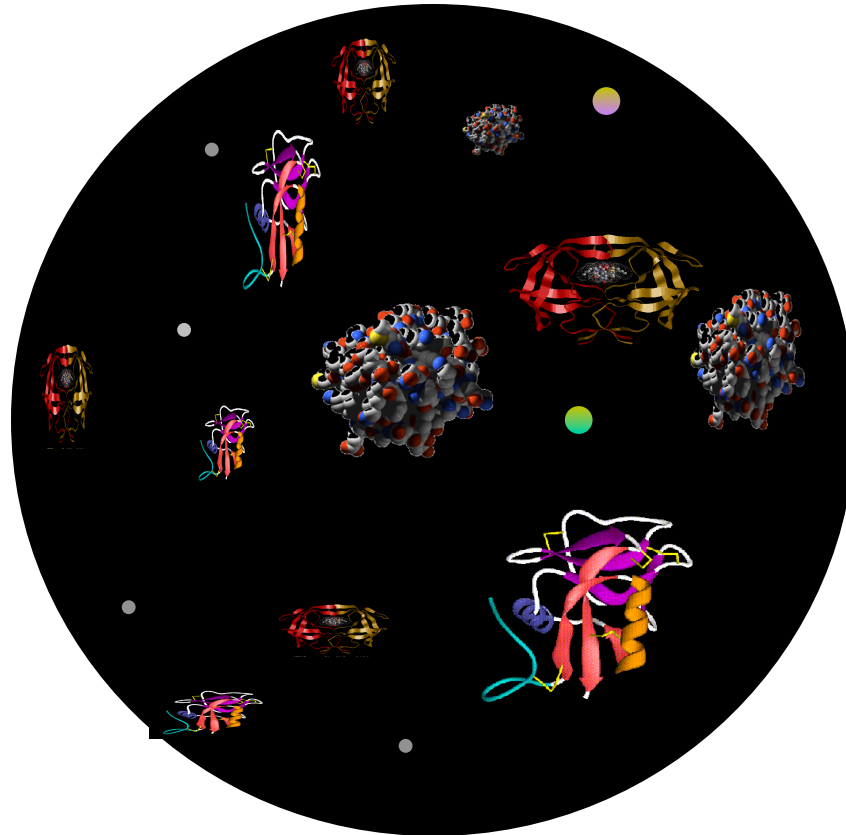# Structural Proteomics: The Motivation

# Protein Structure Initiative*

- **Organize all known protein sequences into sequence families**

- **Select family representatives as targets**

- **Solve the 3D structures of these targets by X-ray or NMR**

- **Build models for the remaining proteins via comparative (homology) modeling**

# Protein Structure Initiative*

- **Organize and recruit interested structural biologists and structure biology centres from around the world**

- **Coordinate target selection**

- **Develop new kinds of high throughput techniques**

- **Solve, solve, solve, solve….**
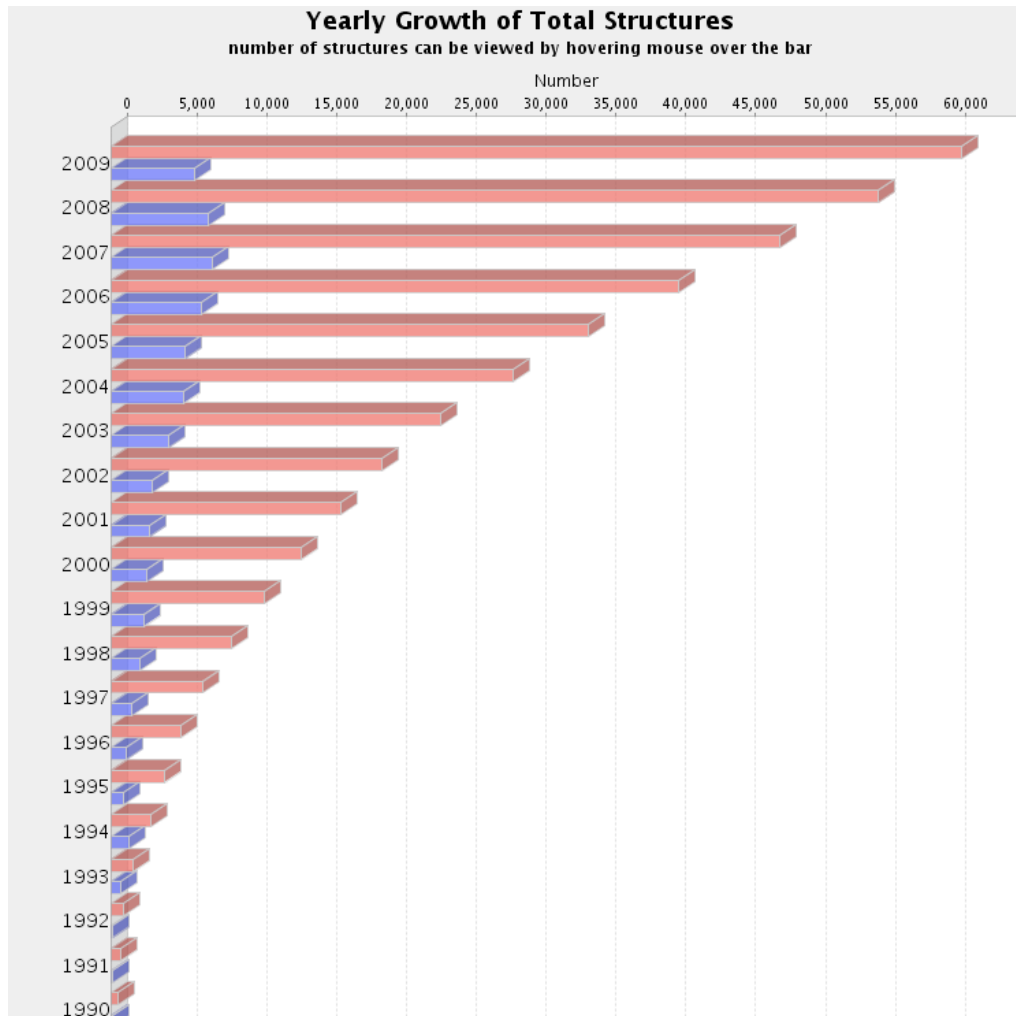
# The Protein Fold Universe

How Big Is It???

500?
 2000?
10000?
∞ ?

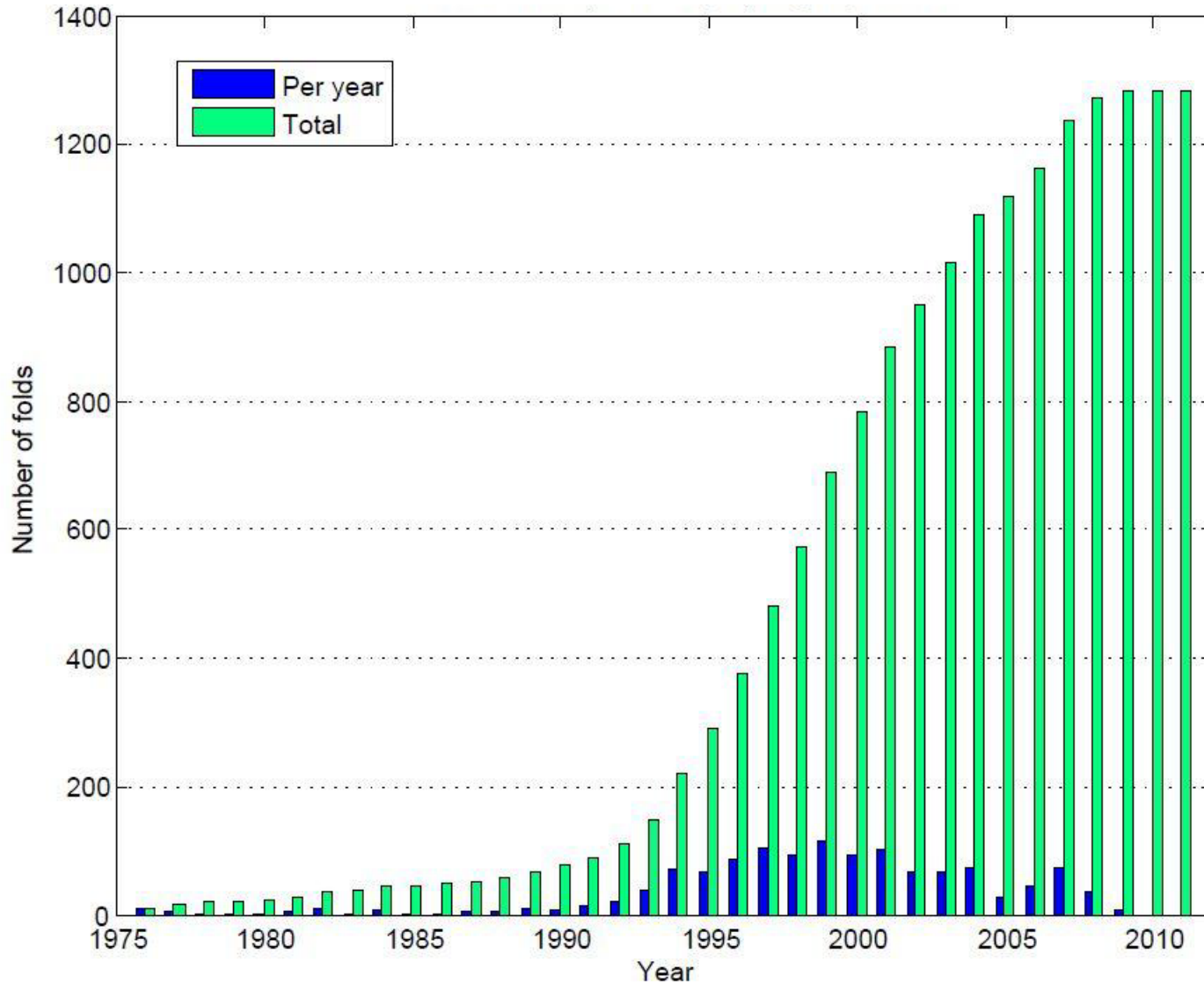*Human Genome Codes for ~21,000 Proteins*

# Structure Deposition Rate



**Yearly Growth of Total Structures**
number of structures can be viewed by hovering mouse over the bar

- **Growth has been exponential for the past 10 years**
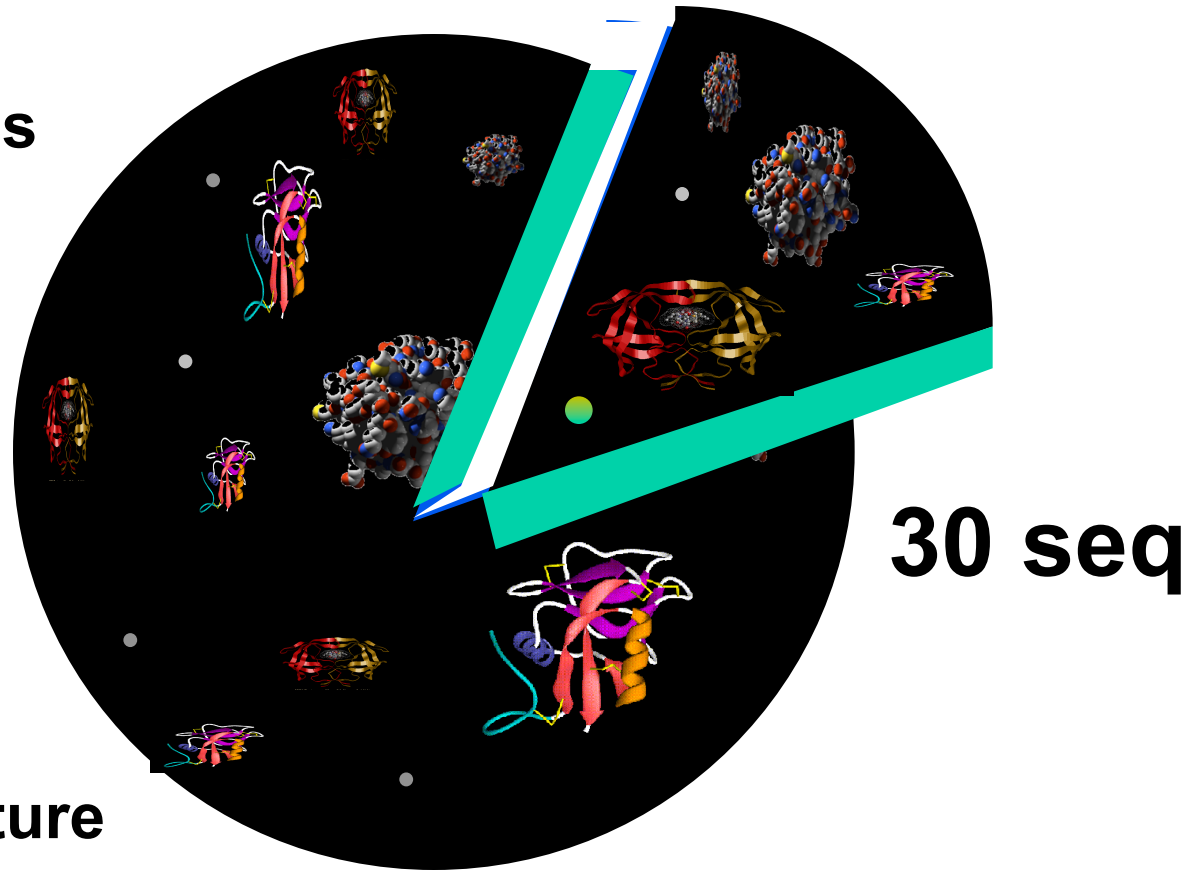- **Approximately 8000 new structures being added each year**
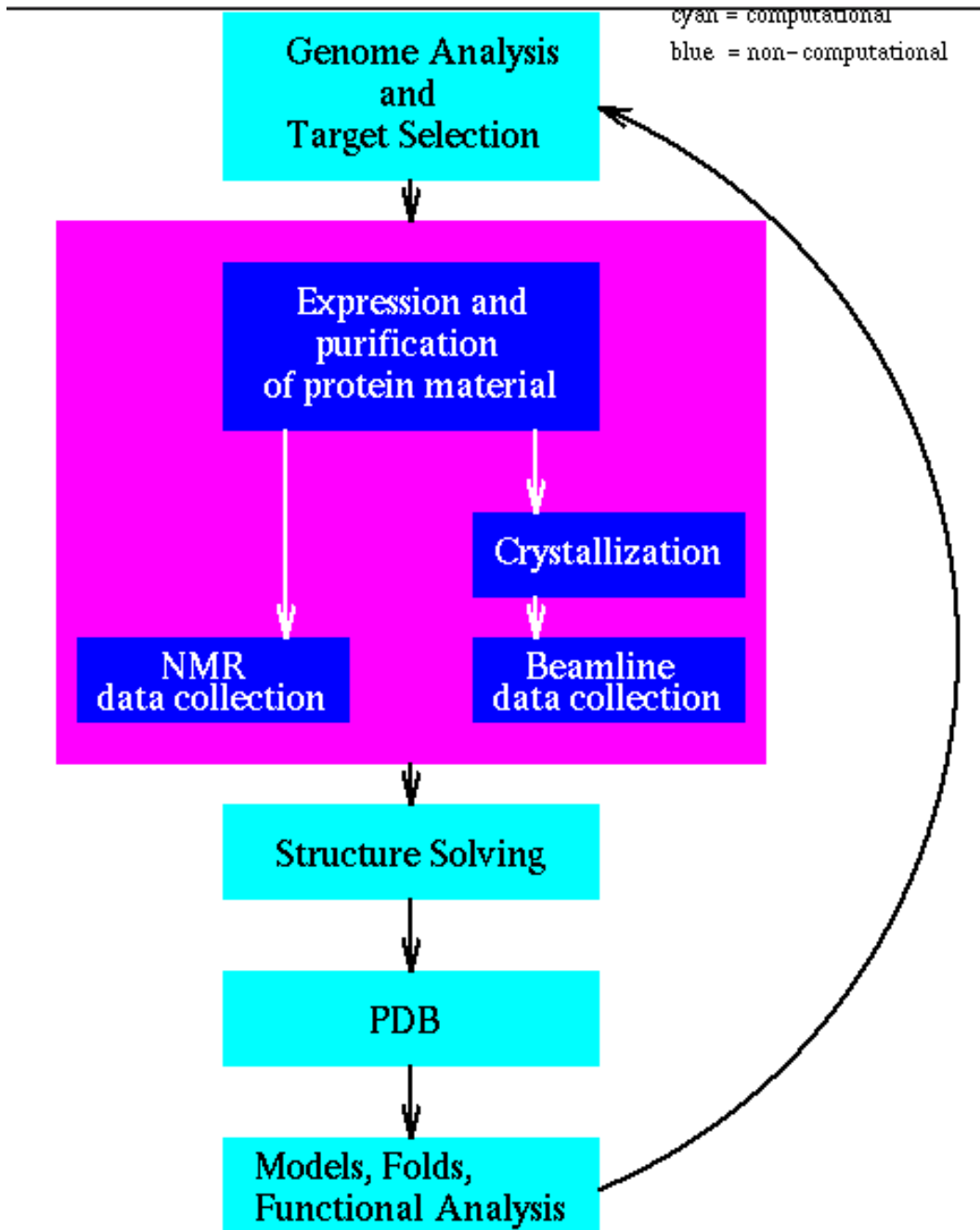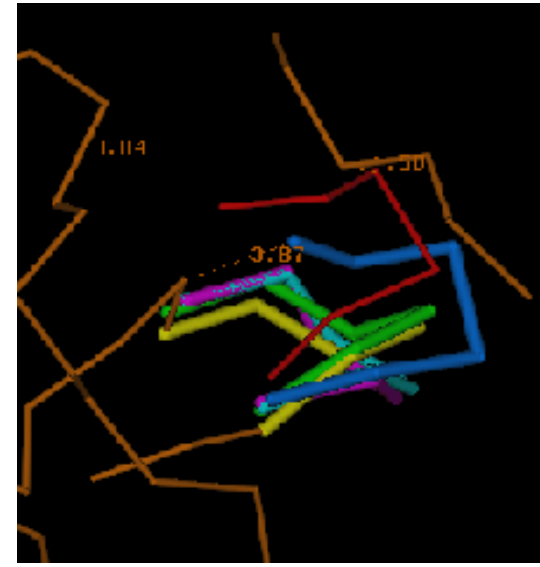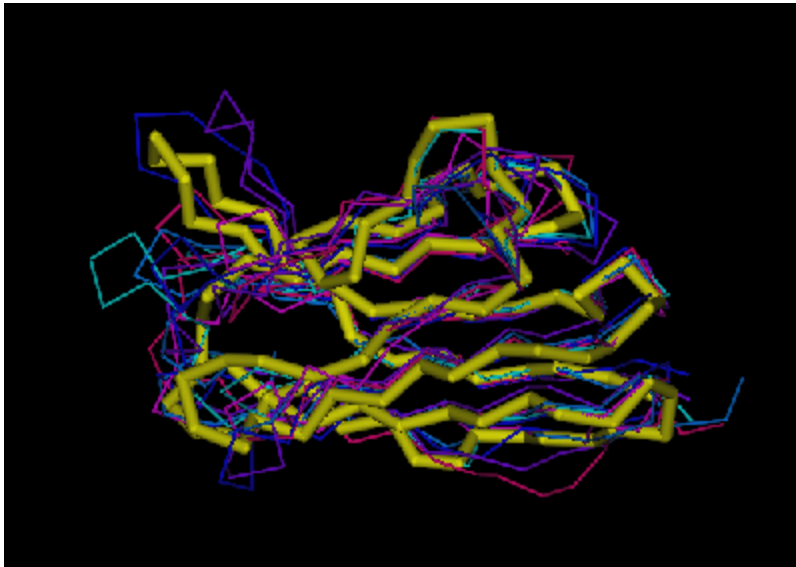
# Number of New Folds in The PDB*

# Protein Structure Initiative

- 25,000 proteins
- 10,000 subset
- 30% ID or
- 30 seq
- Solve by 2010
- $20,000/Structure

30 seq

cyan = computational
blue = non-computational

Genome Analysis
and
Target Selection

Expression and
purification
of protein material

Crystallization

NMR
data collection

Beamline
data collection

Structure Solving

PDB

Models, Folds,
Functional Analysis

# Comparative (Homology) Modelling



ACDEFGHIKLMNPQRST--FGHQWERT------TYREWYEGHADS
ASDEYAHLRILDPQRSTVAYAYE--KSFAPPGSFKWEYEAHADS
MCDEYAHIRLMNPERSTVAGGHQWERT----GSFKEWYAAHADD

# Homology Modelling*

- **Based on the observation that "Similar sequences exhibit similar structures"**

- **Known structure is used as a template to model an unknown (but likely similar) structure with known sequence**

- **First applied in late 1970's using early computer imaging methods (Tom Blundell)**
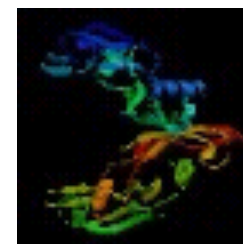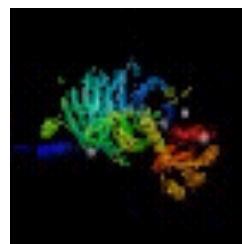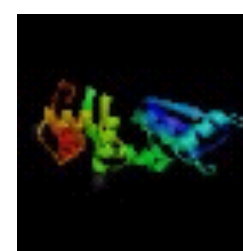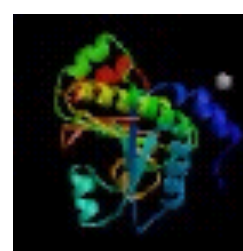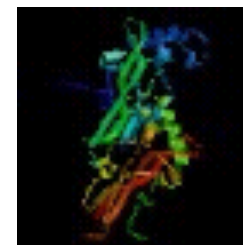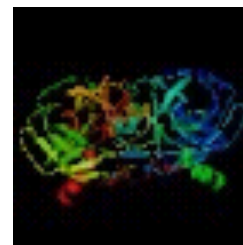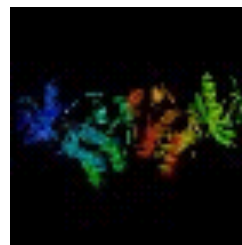
# Homology Modelling*

- **Offers a method to "Predict" the 3D structure of proteins for which it is not possible to obtain X-ray or NMR data**

- **Can be used in understanding function, activity, specificity, etc.**

- **Of interest to drug companies wishing to do structure-aided drug design**

- **A keystone of Structural Proteomics**

# Homology Modelling*

- **Identify homologous sequences in PDB**
- **Align query sequence with homologues**
- **Find Structurally Conserved Regions (SCRs)**
- **Identify Structurally Variable Regions (SVRs)**
- **Generate coordinates for core region**
- **Generate coordinates for loops**
- **Add side chains (Check rotamer library)**
- **Refine structure using energy minimization**
- **Validate structure**

# Step 1: ID Homologues in PDB



```
PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQNCEQWERYTRASDFHG
TREWQIYPASDFGHKLMCNASQERWW
PRETWQLKHGFDSADAMNCVCNQWER
GFDHSDASFWERQWK
```

**Query Sequence**                                          **PDB**

# Step 1: ID Homologues in PDB

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQNCEQWERYTRASDFHG
TREWQIYPASDFGHKLMCNASQERWW
PRETWQLKHGFDSADAMNCVCNQWER
GFDHSDASFWERQWK

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQQWEWEWQWEWEQWEWEWQ
RYEYEWQWNCEQWERYTRASDFHG
TREWQIYPASDWERWEREWRFDSFG

## Hit #1

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQNCEQWERYTRASDFHG
TREWQIYPASDFGHKLMCNASQERWW
PRETWQLKHGFDSADAMNCVCNQWER
GFDHSDASFWERQWK

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQQWEWEWQWEWEQWEWEWQ
RYEYEWQWNCEQWERYTRASDFHG
TR

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQNCEQWERYTRASDFHG
TREWQIYPASDFG

## Hit #2

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQNCEQWERYTRASDFHG
TREWQIYPASDFGPRTEINSEQENCEPR
TEINSEQUENCEPRTEINSEQNCEQWER
YTRASDFHGTREWQIYPASDFG
TREWQIYPASDFGPRTEINSEQENCEPR
TEINSEQUENCEPRTEINSEQNCEQWER

YTRASDFHGTREWQ

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQNCEQWERYTRASDFHG
TREWQIYPASDFG

PRTEINSEQENCEPRTEINSEQUENC
EPRTEINSEQNCEQWERYTRASDFHG
TREWQIYPASDFGPRTEINSEQENC

**Query Sequence**                **PDB**

# Step 2: Align Sequences
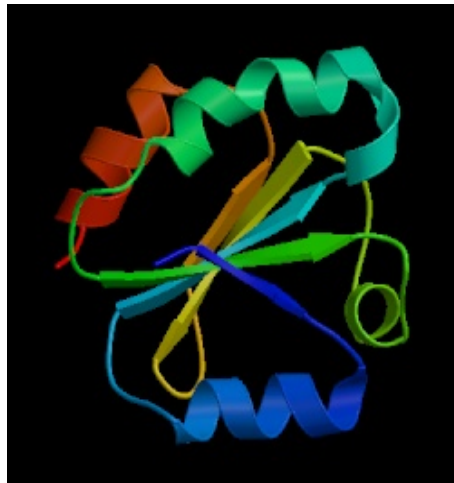
|       | G  | E  | N  | E  | T  | I  | C  | S  |
|-------|----|----|----|----|----|----|----|----|
| **G** | 10 | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| **E** | 0  | 10 | 0  | 10 | 0  | 0  | 0  | 0  |
| **N** | 0  | 0  | 10 | 0  | 0  | 0  | 0  | 0  |
| **E** | 0  | 0  | 0  | 10 | 0  | 0  | 0  | 0  |
| **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 |

**Dynamic Programming**

# Step 2: Align Sequences

```
Query  ACDEFGHIKLMNPQRST--FGHQWERT------TYREWYEG
Hit #1  ASDEYAHLRILDPQRSTVAYAYE--KSFAPPGSFKWEYEA
Hit #2  MCDEYAHIRLMNPERSTVAGGHQWERT----GSFKEWYAA
```



**Hit #1**



**Hit #2**

# Alignment*

- **Key step in Homology Modelling**
- **Global (Needleman-Wunsch) alignment is absolutely required**
- **Small error in alignment can lead to big error in structural model**
- **Multiple alignments are usually better than pairwise alignments**

# Alignment Thresholds*



Threshold for structural homology

# Step 3: Find SCR's

```
Query  ACDEFGHIKLMNPQRST--FGHQWERT------TYREWYEG
Hit #1  ASDEYAHLRILDPQRSTVAYAYE--KSFAPGSFKWEYEA
Hit #2  MCDEYAHIRLMNPERSTVAGGHQWERT----GSFKEWYAA
        HHHHHHHHHHHHHCCCCCCCCCCCCCCCCCBBBBBBBB
```

*SCR #1*                                    *SCR #2*



**Hit #1**                                   **Hit #2**

# Structurally Conserved Regions (SCR's)*

- **Corresponds to the most stable structures or regions (usually interior) of protein**

- **Corresponds to sequence regions with lowest level of gapping, highest level of sequence conservation**

- **Usually corresponds to secondary structures**

# Step 4: Find SVR's

|         |                    |          |            |
|---------|--------------------|----------|------------|
| Query   | ACDEFGHIKLMNPQRST   | --FGHQWERT----- | -TYREWYEG |
| Hit #1  | ASDEYAHLRILDPQRSTV  | AYAYE--KSFAPR    | GSFKWEYEA |
| Hit #2  | MCDEYAHIRLMNPERSTV  | AGGHQWERT-----   | GSFKEWYAA |
|         | HHHHHHHHHHHHH       | CCCCCCCCCCCCCCCCC | BBBBBBBBB |

*SVR (loop)*



**Hit #1**     **Hit #2**

# Structurally Variable Regions (SVR's)*

- **Corresponds to the least stable or most flexible regions (usually exterior) of protein**

- **Corresponds to sequence regions with highest level of gapping, lowest level of sequence conservation**

- **Usually corresponds to loops and turns**

# Step 5: Generate Coordinates

ALA

| ATOM | 1 | N | SER A | 1 | 21.389 | 25.406 | −4.628 | 1.00 | 23.22 | 2TRX | 152 |
| ATOM | 2 | CA | SER A | 1 | 21.628 | 26.691 | −3.983 | 1.00 | 24.42 | 2TRX | 153 |
| ATOM | 3 | C | SER A | 1 | 20.937 | 26.944 | −2.679 | 1.00 | 24.21 | 2TRX | 154 |
| ATOM | 4 | O | SER A | 1 | 21.072 | 28.079 | −2.093 | 1.00 | 24.97 | 2TRX | 155 |
| ATOM | 5 | CB | SER A | 1 | 21.117 | 27.770 | −5.002 | 1.00 | 28.27 | 2TRX | 156 |
| ATOM | 6 | OG | SER A | 1 | 22.276 | 27.925 | −5.861 | 1.00 | 32.61 | 2TRX | 157 |
| ATOM | 7 | N | ASP A | 2 | 20.173 | 26.028 | −2.163 | 1.00 | 21.39 | 2TRX | 158 |
| ATOM | 8 | CA | ASP A | 2 | 19.395 | 26.125 | −0.949 | 1.00 | 21.57 | 2TRX | 159 |
| ATOM | 9 | C | ASP A | 2 | 20.264 | 26.214 | 0.297 | 1.00 | 20.89 | 2TRX | 160 |
| ATOM | 10 | O | ASP A | 2 | 19.760 | 26.575 | 1.371 | 1.00 | 21.49 | 2TRX | 161 |

| ATOM | 1 | N | ALA A | 1 | 21.389 | 25.406 | −4.628 | 1.00 | 23.22 | 2TRX | 152 |
| ATOM | 2 | CA | ALA A | 1 | 21.628 | 26.691 | −3.983 | 1.00 | 24.42 | 2TRX | 153 |
| ATOM | 3 | C | ALA A | 1 | 20.937 | 26.944 | −2.679 | 1.00 | 24.21 | 2TRX | 154 |
| ATOM | 4 | O | ALA A | 1 | 21.072 | 28.079 | −2.093 | 1.00 | 24.97 | 2TRX | 155 |
| ATOM | 5 | CB | ALA A | 1 | 21.117 | 27.770 | −5.002 | 1.00 | 28.27 | 2TRX | 156 |
| ATOM | 6 | OG | SER A | 1 | 22.276 | 27.925 | −5.861 | 1.00 | 32.61 | 2TRX | 157 |
| ATOM | 7 | N | GLU A | 2 | 20.173 | 26.028 | −2.163 | 1.00 | 21.39 | 2TRX | 158 |
| ATOM | 8 | CA | GLU A | 2 | 19.395 | 26.125 | −0.949 | 1.00 | 21.57 | 2TRX | 159 |
| ATOM | 9 | C | GLU A | 2 | 20.264 | 26.214 | 0.297 | 1.00 | 20.89 | 2TRX | 160 |
| ATOM | 10 | O | GLU A | 2 | 19.760 | 26.575 | 1.371 | 1.00 | 21.49 | 2TRX | 161 |

# Step 5: Generate Core Coordinates*

- **For identical amino acids, transfer all atom coordinates (XYZ) to query protein**

- **For similar amino acids, transfer backbone coordinates & replace side chain atoms while respecting $\chi$ angles**

- **For different amino acids, transfer only the backbone coordinates (XYZ) to query sequence**

# Step 6: Replace SVRs (loops)



Query FGHQWERT

Hit #1 YAYE--KS

# Loop Library*

- **Loops extracted from PDB using high resolution (<2 Å) X-ray structures**

- **Typically thousands of loops in DB**

- **Includes loop coordinates, sequence, # residues in loop, Ca-Ca distance, preceding $2^o$ structure and following $2^o$ structure (or their Ca coordinates)**

# Step 6: Replace SVRs (loops)*

- **Must match desired # residues**

- **Must match Ca-Ca distance (<0.5 Å)**

- **Must not bump into other parts of protein (no Ca-Ca distance <3.0 Å)**

- **Preceding and following Ca's (3 residues) from loop should match well with corresponding Ca coordinates in template structure**

# Step 6: Replace SVRs (loops)

- **Loop placement and positioning is done using superposition algorithm**

- **Loop fits are evaluated using RMSD calculations and standard "bump checking"**

- **If no "good" loop is found, some algorithms create loops using randomly generated $\phi/\psi$ angles**

# Step 7: Add Side Chains

# Amino Acid Side Chains*



Two amino acid side chains to indicate the atom naming convention. Hydrogens are not shown.

lysine

tyrosine

# Newman Projections

# Newman Projections*

# Preferred Side Chain χ Angles*

Some  combinations are  **BAD.**

Some are **OK.**

# Relation Between $\chi$ and $\phi/\psi$*

Some $\phi, \chi_1$ combinations are **BAD.**

Some $\psi, \chi_1$ combinations are **BAD.**

The rest are **OK.**

$\phi = -180°, 0°$
$\chi_1 = g-$

$\psi = 0°, 180°$
$\chi_1 = t$

$\phi = -120°, 60°$
$\chi_1 = g+$

$\psi = -60°, 120°$
$\chi_1 = g+$

# Relation Between χ and φ/ψ



Histidine

# Relation Between $\chi$ and $\phi/\psi$



The Ramachandran Plot.

# Relation Between χ and φ/ψ*



SER r1=g+       SER r1=t       SER r1=g-

Box widths proportional to r1 populations in each phi,psi range

**g+**       **t**       **g-**

**Serine**

# Relation Between χ and φ/ψ*



VAL r1=g+

VAL r1=t

VAL r1=g-

Box widths proportional to r1 populations in each phi,psi range

**g+**

**t**

**g-**

**Valine**

# Step 7: Add Side Chains*

- **Done primarily for SVRs (not SCRs)**
- **Rotamer placement and positioning is done via a superposition algorithm using rotamers taken from a standardized library (Trial & Error)**
- **Rotamer fits are evaluated using simple "bump checking" methods**

# Step 8: Energy Minimization*

# Energy Minimization*

- **Efficient way of "polishing and shining" your protein model**

- **Removes atomic overlaps and unnatural strains in the structure**

- **Stabilizes or reinforces strong hydrogen bonds, breaks weak ones**

- **Brings protein to lowest energy in about 1-2 minutes CPU time**

# Energy Minimization (Theory)

- **Treat Protein molecule as a set of balls (with mass) connected by rigid rods and springs**

- **Rods and springs have empirically determined force constants**

- **Allows one to treat atomic-scale motions in proteins as classical physics problems (OK approximation)**

# Standard Energy Function*

$$E = K_r(r_i - r_j)^2 + \qquad \text{Bond length}$$
$$K_\theta(\theta_i - \theta_j)^2 + \qquad \text{Bond bending}$$
$$K_\phi(1-\cos(n\phi_j))^2 + \quad \text{Bond torsion}$$
$$q_i q_j/4\pi\varepsilon r_{ij} + \qquad \text{Coulomb}$$
$$A_{ij}/r^6 - B_{ij}/r^{12} + \qquad \text{van der Waals}$$
$$C_{ij}/r^{10} - D_{ij}/r^{12} \qquad \text{H-bond}$$

# Energy Terms*



$$K_r(r_i - r_j)^2$$

**Stretching**

$$K_\theta(\theta_i - \theta_j)^2$$

**Bending**

$$K_\phi(1-\cos(n\phi_j))^2$$

**Torsional**

# Energy Terms*

$$q_i q_j / 4\pi\varepsilon r_{ij}$$

**Coulomb**

$$A_{ij}/r^6 - B_{ij}/r^{12}$$

**van der Waals**

$$C_{ij}/r^{10} - D_{ij}/r^{12}$$

**H-bond**

# An Energy Surface

**High Energy**

**Low Energy**

**Overhead View**

**Side View**

# Minimization Methods*

- **Energy surfaces for proteins are complex hyperdimensional spaces**

- **Biggest problem is overcoming local minimum problem**

- **Simple methods (slow) to complex methods (fast)**
  - **Monte Carlo Method**
  - **Steepest Descent**
  - **Conjugate Gradient**

# Monte Carlo Algorithm

- **Generate a conformation or alignment (a state)**
- **Calculate that state's energy or "score"**
- **If that state's energy is less than the previous state accept that state and go back to step 1**
- **If that state's energy is greater than the previous state accept it if a randomly chosen number is $< e^{-E/kT}$ where E is the state energy otherwise reject it**
- **Go back to step 1 and repeat until done**

# Conformational Sampling



**Mid-energy**     **lower energy**     **lowest energy**     **highest energy**

# Monte Carlo Minimization



**Performs a progressive or directed random search**

# Steepest Descent & Conjugate Gradients

- **Frequently used for energy minimization of large (and small) molecules**

- **Ideal for calculating minima for complex (I.e. non-linear) surfaces or functions**

- **Both use derivatives to calculate the slope and direction of the optimization path**

- **Both require that the scoring or energy function be differentiable (smooth)**

# Steepest Descent Minimization



**Makes small locally steep moves down gradient**

# Conjugate Gradient Minimization



**Includes information about the prior history of path**

# Energy Minimization*

- **Very complex programs that have taken years to develop and refine**
- **Several freeware options to choose**
  - **XPLOR (Axel Brunger, Yale)**
  - **GROMACS (Gronnigen, The Netherlands)**
  - **AMBER (Peter Kollman, UCSF)**
  - **CHARMM (Martin Karplus, Harvard)**
  - **TINKER (Jay Ponder, Wash U))**

# The Final Result



Modelled

Actual

# Summary*

- **Identify homologous sequences in PDB**
- **Align query sequence with homologues**
- **Find Structurally Conserved Regions (SCRs)**
- **Identify Structurally Variable Regions (SVRs)**
- **Generate coordinates for core region**
- **Generate coordinates for loops**
- **Add side chains (Check rotamer library)**
- **Refine structure using energy minimization**
- **Validate structure**

# How Good are Homology Models?

# Outline

- **The Protein Universe and the Protein Structure Initiative**

- **Homology (Comparative) Modelling of 3D Protein Structures**

- <span style="color:red">**Homology Modelling on the Web**</span>

- **Assessing 3D Structures (modelled and experimental)**

# Modelling on the Web

- **Prior to 1998 homology modelling could only be done with commercial software or command-line freeware**

- **The process was time-consuming and labor-intensive**

- **The past few years has seen an explosion in automated web-based homology modelling servers**

- **Now anyone can homology model!**

# Swiss-Model*



http://swissmodel.expasy.org//SWISS-MODEL.html

# 3D-Jigsaw



http://bmm.cancerresearchuk.org/~3djigsaw/

# Proteus2*



http://www.proteus2.ca/proteus2/

# Modelled Protein Databases

- **Databases containing 3D structural models of 100,000's of proteins and protein domains**

- **Idea is to generate a 3D equivalent of GenBank (saves on everyone having to model everytime they want to look at a structure)**

- **Helps in Proteomics Target Selection**

ModBase Search Page

◄ ► ⌂ 🖶 A A + 🔖 http://modbase.compbio.ucsf.edu/modbase-cgi/search_form.cgi   ↻   Q▾ Google

📖 ▦ Department o...ell Biology   Login– Depar... of Alberta   Audiobaba Music Search   Bioinformati... the U of A!   Coilgun Basics 2   Pathguide: t...esource list   »

UCSF University of California, San Francisco | About UCSF | UCSF Medical Center

| Home | User Login | ModBase Search Page | ModWeb Modelling Server | Help | Current Logins |

# Database of Comparative Protein Structure Models

Welcome to ModBase, a database of three-dimensional protein models calculated by comparative modeling.

**General Information**

**Statistics**

**News**

**Project Pages**

**Documentation**

**Authors and Acknowledgements**

**Publications**

**Todo List**

**Related Resources**

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

**ModBase search form**                                                    [Search]

Search type ❓ [Model(Default) ▾]        Display type ❓ [Model Detail (graphical) ▾]

**All available datasets are selected** ❓           **Select specific dataset(s)**

To include the academic (comprehensive) dataset, go to 'User Login'!

**Search by properties**

Property ❓ [Database Accession Number ▾]        [_____]

Organism ❓ [ALL ▾]    or    [_____]

**Advanced search**

# Outline

- **The Protein Universe and the Protein Structure Initiative**

- **Homology (Comparative) Modelling of 3D Protein Structures**

- **Homology Modelling on the Web**

- **Assessing 3D Structures (modelled and experimental)**

# Why Assess Structure?

- **A structure can (and often does) have mistakes**

- **A poor structure will lead to poor models of mechanism or relationship**

- **Unusual parts of a structure may indicate something important (or an error)**

# Famous "bad" structures*

- **Azobacter ferredoxin** (wrong space group)

- **Zn-metallothionein** (mistraced chain)

- **Alpha bungarotoxin** (poor stereochemistry)

- **Yeast enolase** (mistraced chain)

- **Ras P21 oncogene** (mistraced chain)

- **Gene V protein** (poor stereochemistry)

# How to Assess Structure?*

- **Assess experimental fit (look at R factor or rmsd)**

- **Assess correctness of overall fold (look at disposition of hydrophobes)**

- **Assess structure quality (packing, stereochemistry, bad contacts, etc.)**

# A Good Protein Structure..*

## X-ray structure

- R = 0.59 random chain
- R = 0.45 initial structure
- R = 0.35 getting there
- R = 0.25 typical protein
- R = 0.15 best case
- R = 0.05 small molecule

## NMR structure

- rmsd = 4 Å random
- rmsd = 2 Å initial fit
- rmsd = 1.5 Å OK
- rmsd = 0.8 Å typical
- rmsd = 0.4 Å best case
- rmsd = 0.2 Å dream on

# A Good Protein Structure..*

- **Minimizes disallowed torsion angles**
- **Maximizes number of hydrogen bonds**
- **Maximizes buried hydrophobic ASA**
- **Maximizes exposed hydrophilic ASA**
- **Minimizes interstitial cavities or spaces**

# A Good Protein Structure..*

- **Minimizes number of "bad" contacts**
- **Minimizes number of buried charges**
- **Minimizes radius of gyration**
- **Minimizes covalent and noncovalent (van der Waals and coulombic) energies**

# Radius & Radius of Gyration

- **RAD = 3.95 x NUMRES$^{0.6}$ + 7.25**        **(Folded)**

- **RADG = 0.41 x (110 x NUMRES) $^{0.5}$**     **(Unfolded)**



*Radius*

*Radius of Gyration*

# Packing Volume



Loose Packing          Dense Packing          Protein

*Proteins are Densely Packed*

# Accessible Surface Area

# Accessible Surface Area*

# Accessible Surface Area*

- **Solvation free energy is related to ASA**
  - ◆ $\Delta G = \Sigma \Delta \sigma_i A_i$
- **Proteins typically have 60% of their ASA comprised of polar atoms or residues**
- **Proteins typically have 40% of their ASA comprised of nonpolar atoms or residues**
- $\Delta_{ASA}$ **(obs - exp.) reveals shape/roughness**

# Structure Validation Servers

- **WhatIf Web Server -** http://swift.cmbi.ru.nl/servers/html/index.html

- **Protein Structure Validation Suite -** http://psvs-1_3.nesg.org/

- **Verify3D -** http://nihserver.mbi.ucla.edu/Verify_3D/

- **Molprobity -** http://molprobity.biochem.duke.edu/

- **PROSESS -** http://www.prosess.ca/

- **VADAR -** http://vadar.wishartlab.com/

## Verify3D Structure Evaluation Server

[Servers Home]

**People** ♦ **Seminars**
**Lectures** ♦ **Webmail**
**Links** ♦ **Facilities**
**Software** ♦ **Home**

**UCLA**

The UCLA-DOE Structure Evaluation server is a tool designed to help in the refinement of crystallographic structures. It will provide you with a visual analysis of the quality of a putative crystal structure for a protein. Verify3D expects this crystal structure to be submitted in PDB format. Please note that Verify3D works best on proteins with at least 100 residues. To submit a crystal structure for analysis, simply select it with the file dialog which is activated by clicking on the Browse button below, then click the Send File button.

Form Based PDB File Upload:
[ Choose File ] no file selected

[ Send File ]  [ Clear Form ]  [ Refresh ]

Verify3D analyzes the compatibility of an atomic model (3D) with its own amino acid sequence (1D). Each residue is assigned a structural class based on its location and environment (alpha, beta, loop, polar, nonpolar, etc). A collection of good structures is used as a reference to obtain a score for each of the 20 amino acids in this structural class. The scores of a sliding 21-residue window (from -10 to +10) are added and plotted for individual residues.

## Obtain your own standalone copy of Profile Search/Environments program/Verify 3D

References: [Bowie *et al.*, 1991; Luethy *et al.*, 1992]. end_a_page_with_links();

**High scores = good   Low scores = bad**

# VADAR*



http://vadar.wishartlab.com/

# VADAR

# Structure Validation Programs

- **PROCHECK -** http://www.biochem.ucl.ac.uk/~roman/procheck/procheck.html

- **PROSA II -** http://lore.came.sbg.ac.at/People/mo/Prosa/prosa.html

- **WhatCheck -** http://swift.cmbi.kun.nl/gv/whatcheck/

- **PDB Validation Suite -** http://sw-tools.pdb.org/apps/VAL/index.html

- **DSSP -** http://swift.cmbi.kun.nl/gv/dssp/

# Procheck*

# Summary

- **Homology modeling is the most accurate method known for predicting 3D protein structures**

- **Recent advances have made homology modeling trivial to do over the web**

- **There are many different ways of evaluating and validating the quality of 3D structure models**

- *Homework: spend 15-20 minutes visiting the websites mentioned today*

# How To Do Your Assignment

- **Follow the instructions carefully**

- **Each of the programs or websites you need to use has been mentioned in the last 3 lectures, if you're smart you may only need to use 3 (local) tools**

- **This assignment will take 4-5 hours to complete and should be 6-8 pages long**