# Lecture 3

# Protein Structure Prediction

Marc A. Marti-Renom

Assistant Adjunct Professor

Department of Biopharmaceutical Sciences

August 22, 2003

BayGenomics

# *Summary*

- Protein Structure Prediction and why is it useful?

- Methods in Protein Structure Prediction

- Comparative Modeling

  - ✓ Steps in CM (overview + resources)

  - ✓ Accuracy/Applications of comparative models

  - ✓ Case example in MODELLER

  - ✓ CM and Structural Genomics

BayGenomics

# *Summary*

- **Protein Structure Prediction and why is it useful?**

- Methods in Protein Structure Prediction

- Comparative Modeling

    - ✓ Steps in CM (overview + resources)

    - ✓ Accuracy/Applications of comparative models

    - ✓ Case example in MODELLER

    - ✓ CM and Structural Genomics

BayGenomics

# Why protein structure **prediction**?

| | Y 2003 | Y 2005 |
|---|---|---|
| Sequences | 1,000,000 | millions |
| Structures | 28,000 | 50,000 |

# Why protein structure prediction?

| | Y 2003 |
|---|---|
| Sequences | 1,000,000 |
| Structures | 300,000 |

**Theory**
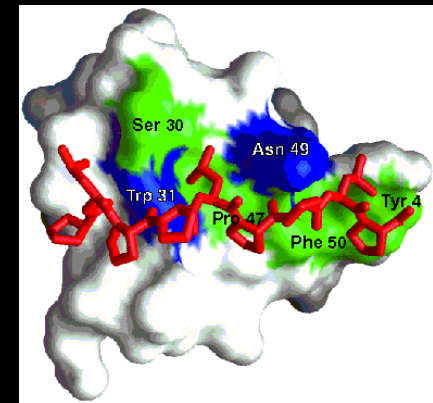
**Experiment**

http://salilab.org/modbase/
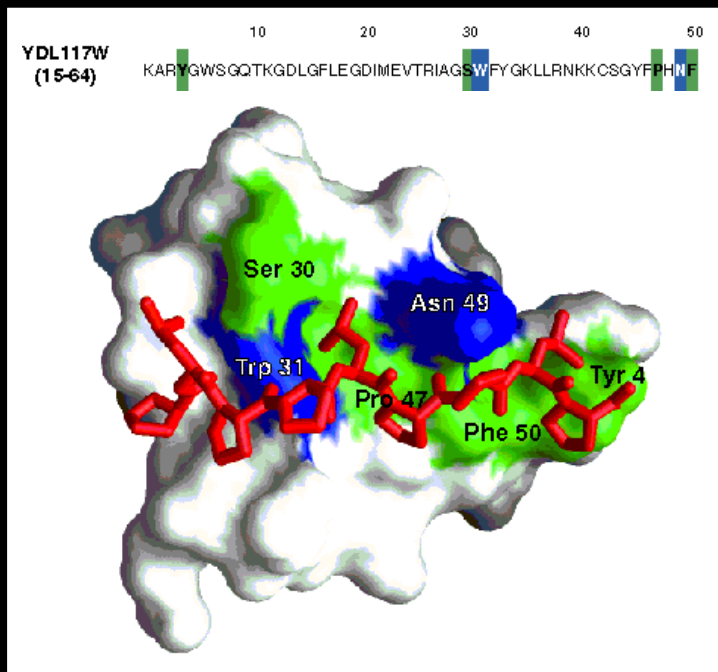
BayGenomics

# Function *via* Structure

Sequence  → **Structure**  → Function

ASILPKRLFGNC

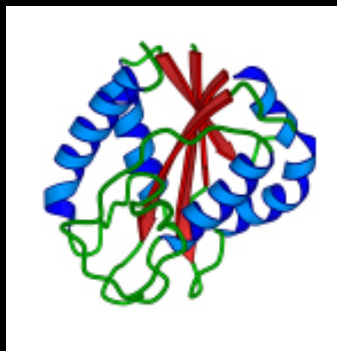# Why is it useful to know the structure of a protein, not only its sequence?

- The biochemical function (activity) of a protein is defined by its interactions with other molecules.

- The biological function is in large part a consequence of these interactions.

- 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.



In addition, since evolution tends to conserve function and function depends more directly on structure than on sequence, structure is more conserved in evolution than sequence
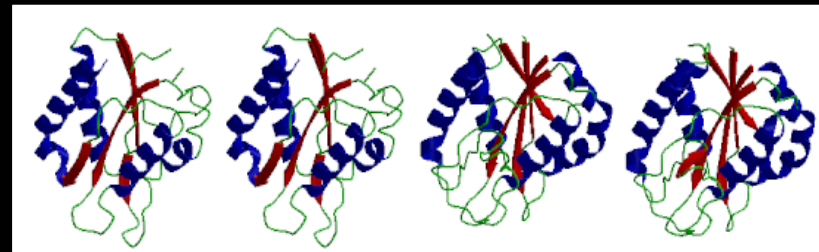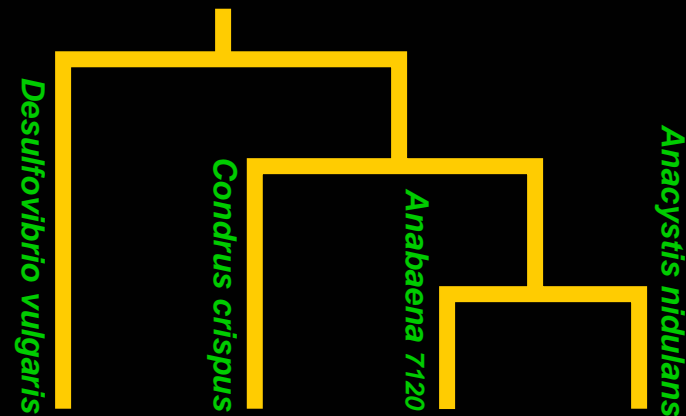
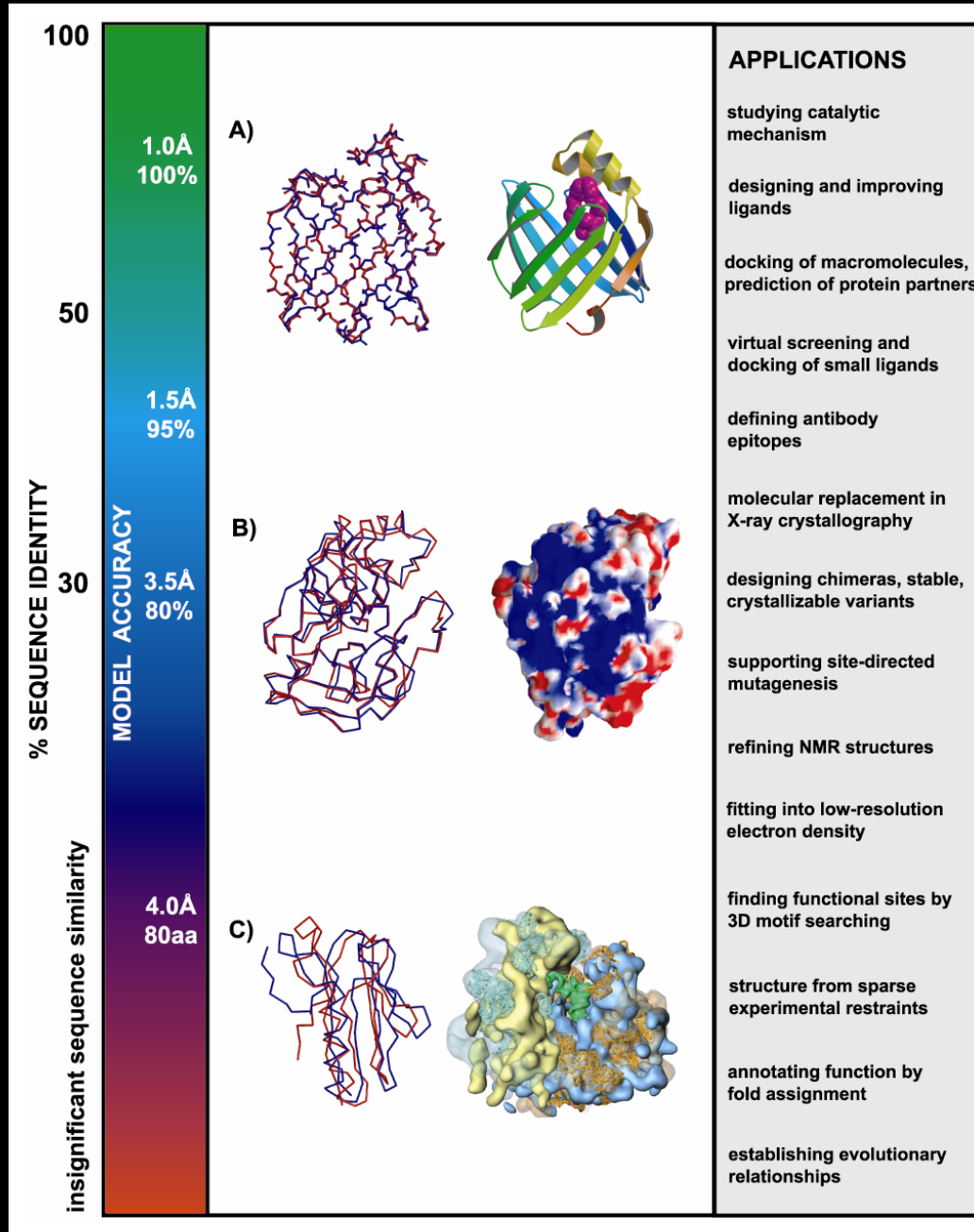The net result is that patterns in space are frequently more recognizable than patterns in sequence.

# *Summary*

- Protein Structure Prediction and why is it useful?

- **Methods in Protein Structure Prediction**

- Comparative Modeling

  - ✓ Steps in CM (overview + resources)

  - ✓ Accuracy/Applications of comparative models

  - ✓ Case example in MODELLER

  - ✓ CM and Structural Genomics

BayGenomics

# Resolution ⟵→ Methods



A. Šali & J. Kuriyan.
*TIBS* **22**, M20, 1999.

# **Methods** for Protein Structure Prediction

- *Ab Initio*
  - *ROSETTA*
    *[http://depts.washington.edu/bakerpg/]*

- *Threading – Fold assignment*
  - *THREADER*
    *[http://www.hgmp.mrc.ac.uk/Registered/Option/threader.html]*

- *Comparative Modeling*
  - *MODELLER*
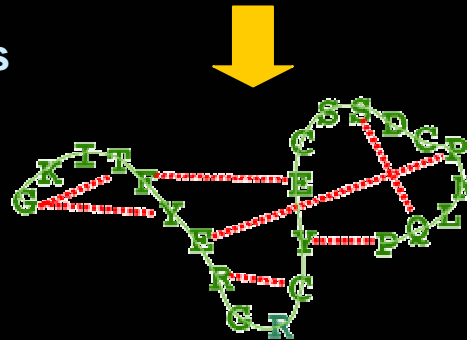    *[http://www.salilab.org/modeller]*

BayGenomics

# *Summary*

- Protein Structure Prediction and why is it useful?

- Methods in Protein Structure Prediction

- Comparative Modeling

  - ✓ Steps in CM (overview + resources)

  - ✓ Accuracy/Applications of comparative models

  - ✓ Case example in MODELLER

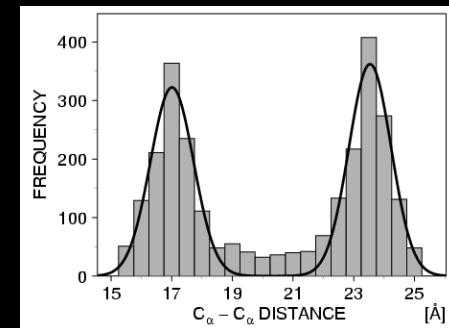  - ✓ CM and Structural Genomics

BayGenomics

# Comparative Modeling by Satisfaction of Spatial Restraints (MODELLER)

3D GKITFYERGFQGHCYESDC-NLQP...
SEQ GKITFYERG---RCYESDCPNLQP...

**1. Extract spatial restraints**



**2. Satisfy spatial restraints**
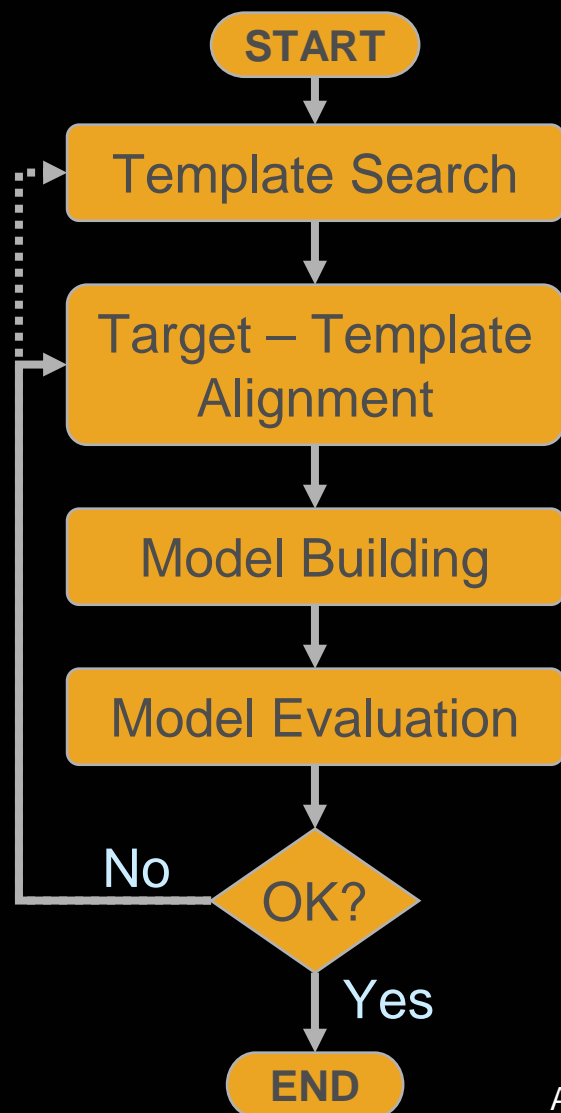
$$F(\mathbf{R}) = \prod_i p_i(f_i/l)$$

A. Šali & T. Blundell. *J. Mol. Biol.* **234**, 779, 1993.
J.P. Overington & A. Šali. *Prot. Sci.* **3**, 1582, 1994.
A. Fiser, R. Do & A. Šali, *Prot. Sci.*, 9, 1753, 2000.

http://www.salilab.org/

# *Summary*

- Protein Structure Prediction and why is it useful?

- Methods in Protein Structure Prediction

- Comparative Modeling

  - ✓ **Steps in CM (overview + resources)**

  - ✓ Accuracy/Applications of comparative models

  - ✓ Case example in MODELLER

  - ✓ CM and Structural Genomics

BayGenomics

# Steps in **Comparative** Protein Structure **Modeling**

**START**

**Template Search**

**Target – Template Alignment**

**Model Building**

**Model Evaluation**

**OK?**

No

Yes

**END**

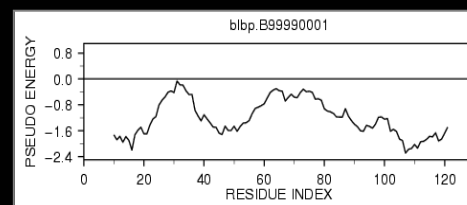**TARGET**

ASILPKRLFGNCEQTSDEGLK
IERTPLVPHISAQNVCLKIDD
VPERLIPERASFQWMNDK

ASILPKRLFGNCEQTSDEGLKIERTPLVPHISAQNVCLKIDDVPERLIPE
MSVIPKRLYGNCEQTSEEAIRIEDSPIV---TADLVCLKIDEIPERLVGE

**TEMPLATE**

blbp.B99990001

PSEUDO ENERGY

0.8
0.0
−0.8
−1.6
−2.4

0   20   40   60   80   100   120
RESIDUE INDEX

A. Šali, *Curr. Opin. Biotech.* 6, 437, 1995.
R. Sánchez & A. Šali, *Curr. Opin. Str. Biol.* 7, 206, 1997.
M. A. Martí-Renom *et al. Ann. Rev. Biophys. Biomolec. Struct.*, 29, 291, 2000.

BayGenomics

# Template Search Methods

- ## Sequence similarity searches
  - BLAST [http://www.ncbi.nlm.nih.gov/BLAST/]
  - FastA program [http://www.ebi.ac.uk/fasta33/]

- ## Profile and iterative methods
  - HMMs [http://www.cse.ucsc.edu/research/compbio/HMM-apps/]
  - PSI-BLAST [http://www.ncbi.nlm.nih.gov/BLAST/]

- ## Structure based threading
  - THREADER [http://bioinf.cs.ucl.ac.uk/threader/]
  - PROFIT [http://www.came.sbg.ac.at/]

BayGenomics

# Target – Template Alignment Methods

- Dynamic Programming Pairwise Alignment
    - ALIGN [http://www.salilab.org/modeller/]

- Multiple Alignments,
    - Psi-Blast [http://www.ncbi.nlm.nih.gov/BLAST/]
    - HMM [http://www.cse.ucsc.edu/research/compbio/HMM-apps/]
    - ALIGN4D [http://www.salilab.org/modeller/]
    - CLUSTALW [http://www.ebi.ac.uk/clustalw/]

- Structure based approaches
    - Threading [http://bioinf.cs.ucl.ac.uk/threader/]

BayGenomics

# Model Building Methods

- ## Rigid Body Assembly
    - COMPOSER [http://www-cryst.bioc.cam.ac.uk/]

- ## Segment Matching
    - SEGMOD

- ## Satisfaction of Spatial Restraints
    - MODELLER [http://www.salilab.org/modeller/]

BayGenomics

# Model Evaluation methods

- ## Stereochemistry

  - PROCHECK/ WHAT-IF
    [http://www.biochem.ucl.ac.uk/~roman/procheck/procheck.html]

- ## Environment

  - VERIFY3D [http://www.doe-mbi.ucla.edu/Services/Verify_3D/]

- ## Statistical potentials based methods

  - PROSAII [http://www.came.sbg.ac.at/]

  - ANOLEA [http://protein.bio.puc.cl/cardex/servers/index.html]

`http://www.salilab.org/bioinformatics_resources.shtml`

BayGenomics

# *Summary*

- Protein Structure Prediction and why is it useful?

- Methods in Protein Structure Prediction

- Comparative Modeling

  - ✓ Steps in CM (overview + some details)

  - ✓ Accuracy/Applications of comparative models

  - ✓ Case example in MODELLER

  - ✓ CM and Structural Genomics

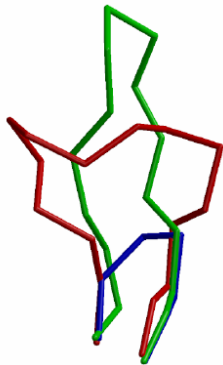BayGenomics

# Typical Errors in Comparative Models
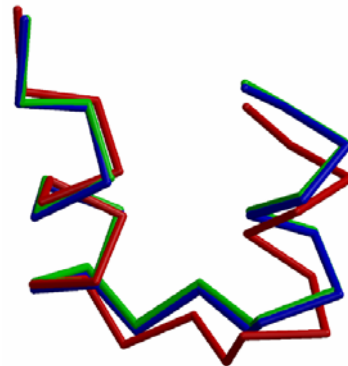


MODEL
X-RAY
TEMPLATE

**Incorrect template**

**Misalignment**

**Region without a template**

**Distortion in correctly aligned regions**

**Sidechain packing**

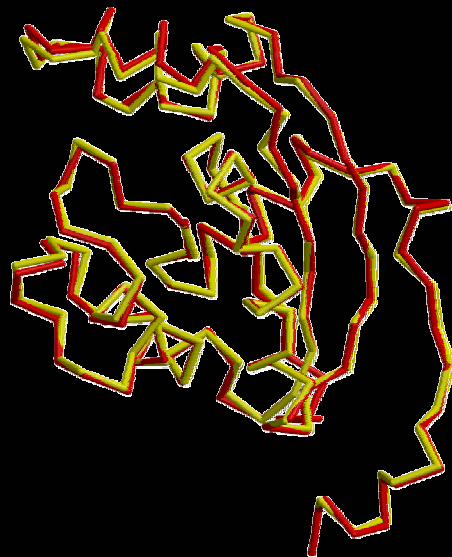# Model Accuracy as a Function of Target-Template Sequence Identity



Sánchez, R., Šali, A. *Proc Natl Acad Sci U S A.* 95 pp13597-602. (1998).

# Model Accuracy

Marti-Renom *et al.* Annu.Rev.Biophys.Biomol.Struct. **29**, 291-325, 2000.

**HIGH ACCURACY**      **MEDIUM ACCURACY**      **LOW ACCURACY**

NM23                        CRABP                      EDN
Seq id 77%                  Seq id 41%                 Seq id 33%
C$\alpha$ equiv 147/148     C$\alpha$ equiv 122/137    C$\alpha$ equiv 90/134
RMSD 0.41Å                  RMSD 1.34Å                 RMSD 1.17Å

Sidechains                  Sidechains                 Sidechains
                            Core backbone              Core backbone
                            Loops                      Loops
                                                       Alignment
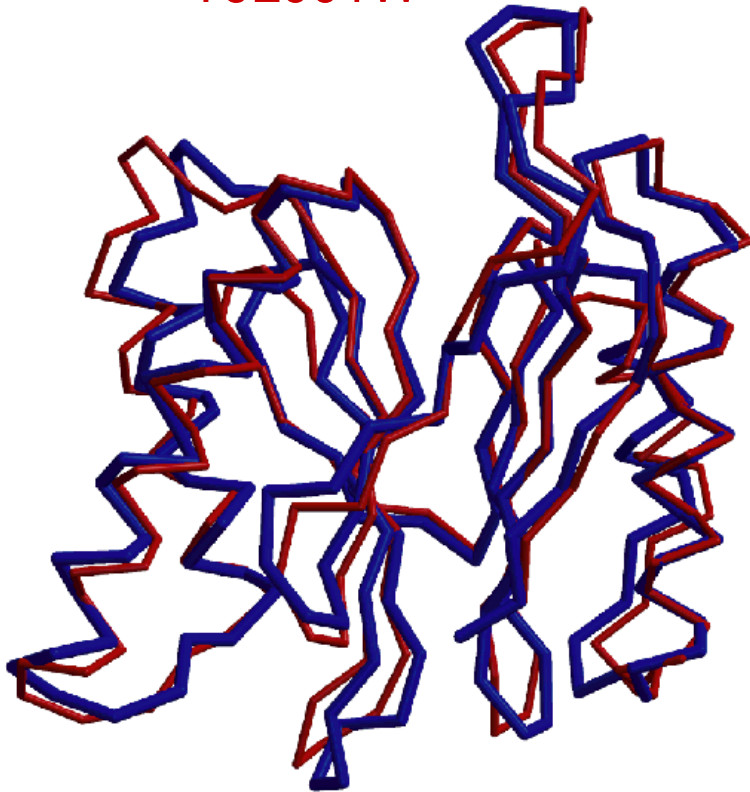X-RAY    MODEL                                         Fold assignment

BayGenomics

# Some Models Can Be Surprisingly Accurate
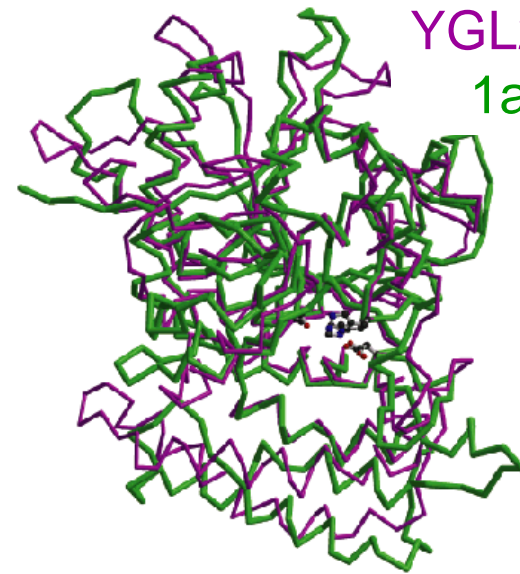## (in Some Core or Active Site Regions)



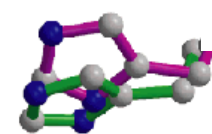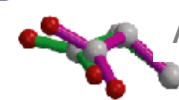24% sequence identity

YJL001W

25% sequence identity
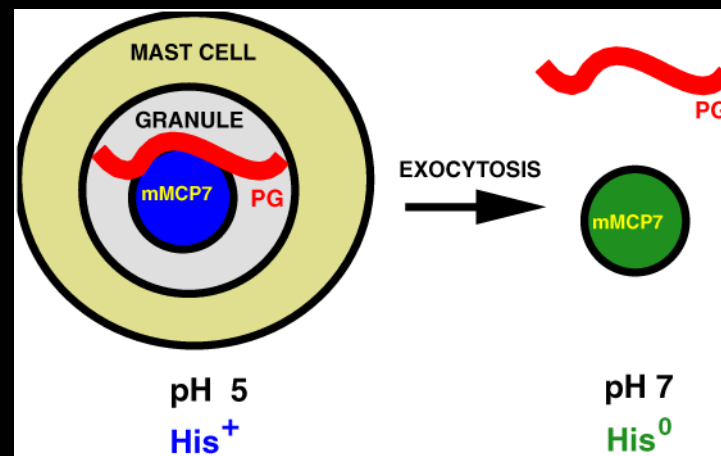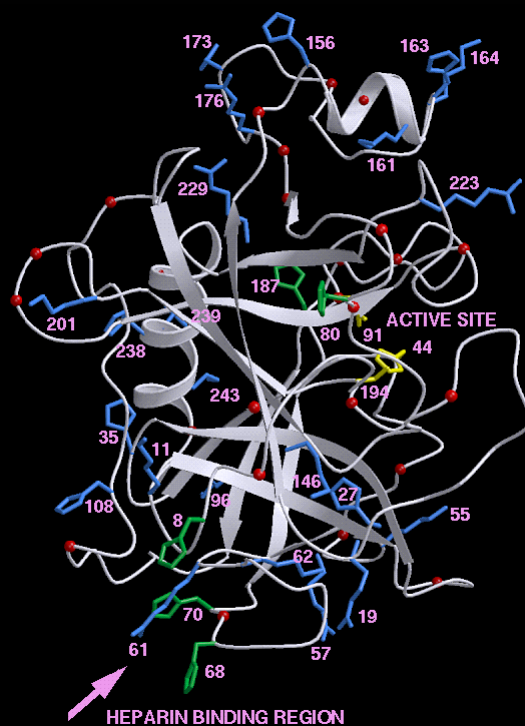
YGL203C
1ac5

His 488

Ser 176

Asp 383

BayGenomics

**Do mast cell proteases bind proteoglycans? Where? When?**

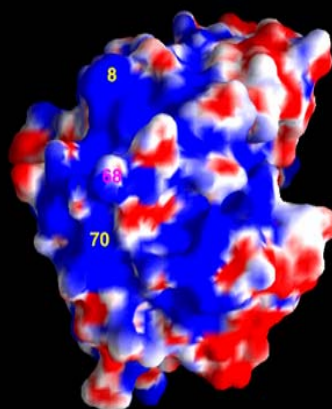**Predicting features of a model that are not present in the template**

1. mMCPs bind negatively charged proteoglycans through electrostatic interactions?
2. Comparative models used to find clusters of positively charged surface residues.
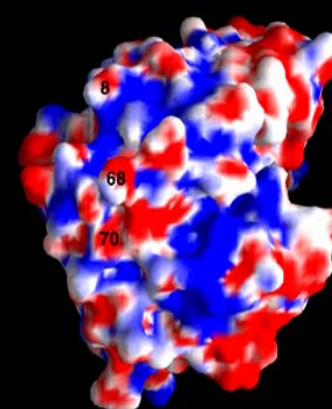3. Tested by site-directed mutagenesis.

Huang *et al. J. Clin. Immunol.* **18**,169,1998.
Matsumoto *et al. J.Biol.Chem.* **270**,19524,1995.
Šali *et al. J. Biol. Chem.* **268**, 9023, 1993.

# Some Models Can Be Used in Docking to Density Maps
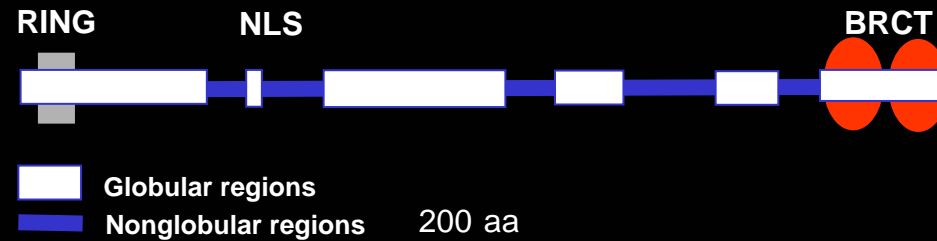## (Yeast Ribosomal 40S subunit)

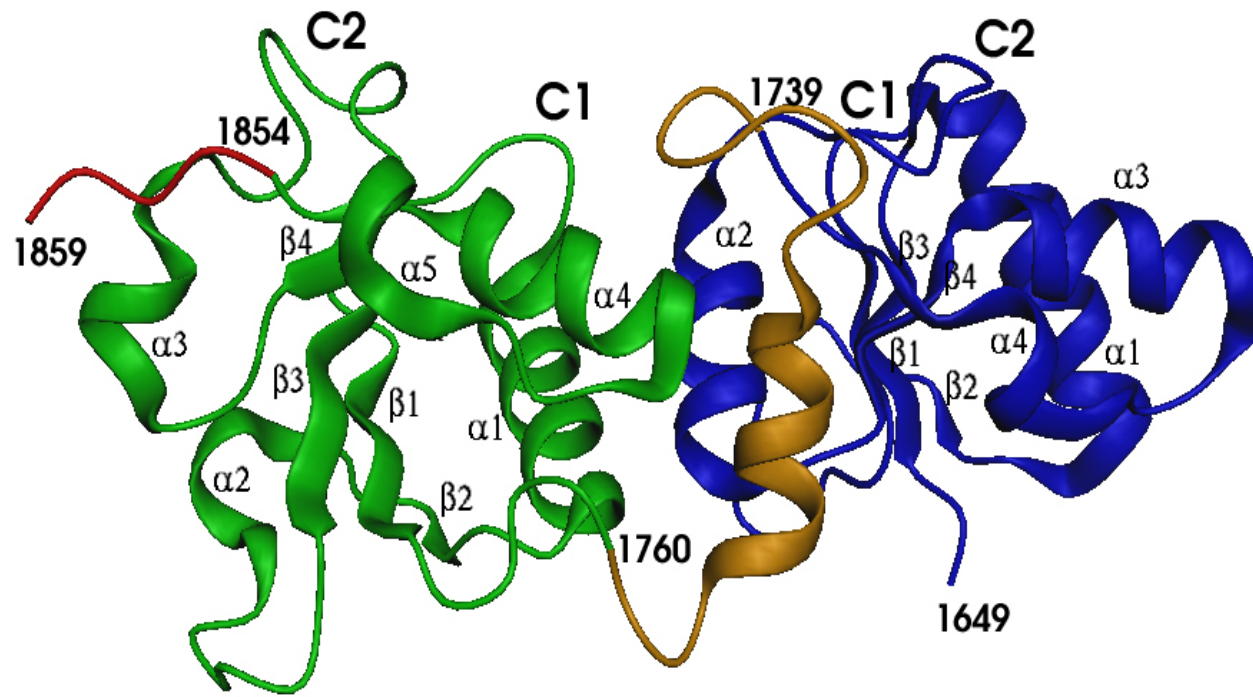

**Docking of comparative models into the cryo-EM map.**

Spahn *et al.* 2001 Cell **107**:373-386

Small 30S subunit from *Thermus thermophilus*
Large 50S subunit from *Haloarcula marismortui*

# Human BRCA1 and its two BRCT domains
## (structural analysis of missense mutations SNPs)



BRCA1 BRCT repeats, 1jnx

Williams, Green, Glover. *Nat.Struct.Biol.* 8, 838, 2001

CONFIDENTIAL

# MYRIAD

*BRACAnalysis*™
Comprehensive BRCA1-BRCA2 Gene Sequence Analysis Result

| | SPECIMEN | PATIENT |
|---|---|---|
| Niecee Singer, MS | Specimen Type: Blood | Name: |
| Strang Cancer Prevention Center | Draw Date: n/a | Date of Birth: Feb 02, 1953 |
| 428 E 72nd St | Accession Date: Oct 27, 2000 | Patient ID: |
| New York, NY 10021 | Report Date: Nov 17, 2000 | Gender: Female |
| | | Accession #: 00019998 |
| | | Requisition #: 56694 |

Physician: Fred Gilbert, MD

## Test Result

| Gene Analyzed | Specific Genetic Variant |
|---|---|
| BRCA2 | H2116R |
| BRCA1 | None Detected |

## Interpretation

### GENETIC VARIANT OF UNCERTAIN SIGNIFICANCE

The BRCA2 variant H2116R results in the substitution of arginine for histidine at amino acid position 2116 of the BRCA2 protein. Variants of this type may or may not affect BRCA2 protein function. Therefore, the contribution of this variant to the relative risk of breast or ovarian cancer cannot be established solely from this analysis. The observation by Myriad Genetic Laboratories of this particular variant in an individual with a deleterious truncating mutation in BRCA2, however, reduces the likelihood that H2116R is itself deleterious.

Authorized Signature:

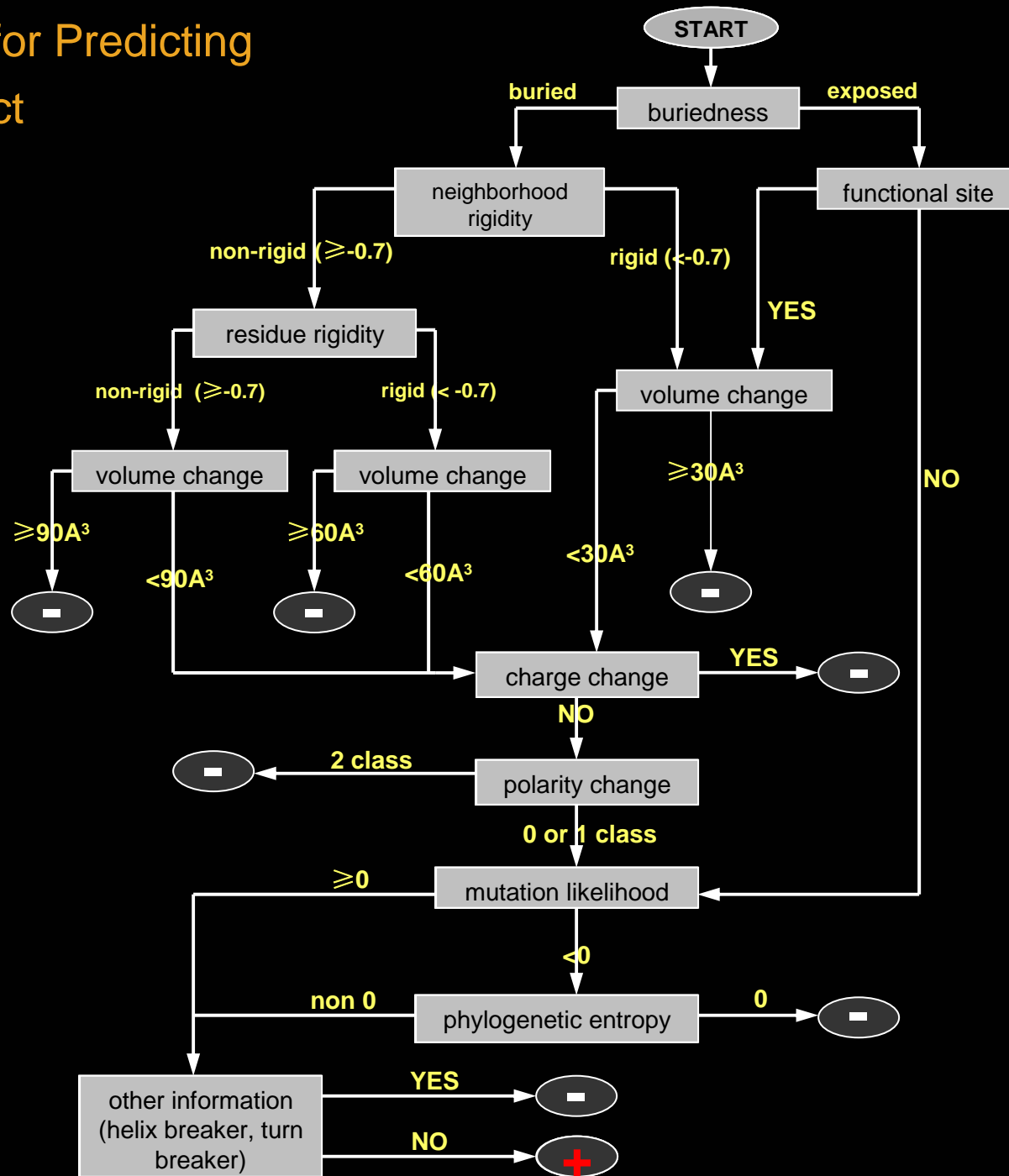Brian E. Ward, Ph.D.
Laboratory Director

Thomas S. Frank, M.D.
Medical Director

These test results should only be used in conjunction with the patient's clinical history and any previous analysis of appropriate family members. It is strongly recommended that these results be communicated to the patient in a setting that includes appropriate counseling. The accompanying Technical Specifications summary describes the analysis, method, performance characteristics, nomenclature, and interpretive criteria of this test. This test may be considered investigational by some states. This test was developed and its performance characteristics determined by Myriad Genetic Laboratories. It has not been reviewed by the U.S. Food and Drug Administration. The FDA has determined that such clearance or approval is not necessary.
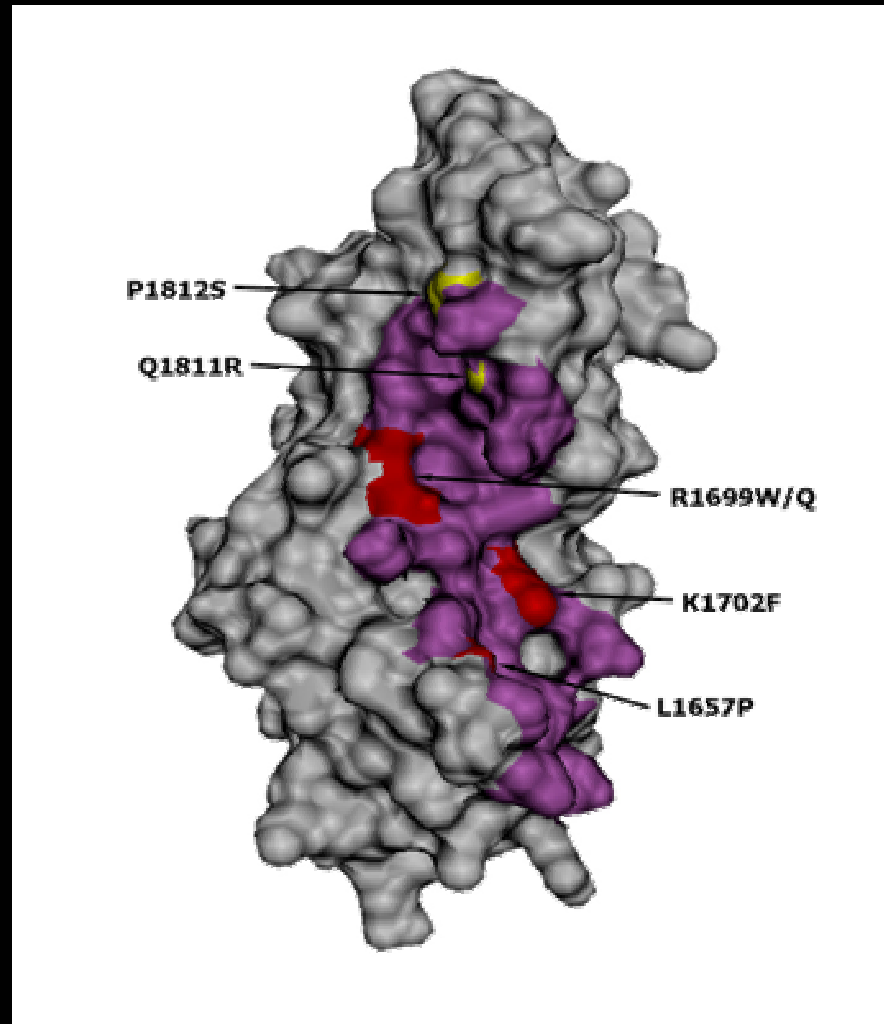
BayGenomics

# Missense Mutations in BRCT Domains by Function

|  | cancer associated | not cancer associated | ? |
|---|---|---|---|
| **no transcription activation** | C1697R<br>R1699W<br>A1708E<br>S1715R<br>P1749R<br>M1775R |  | M1652K L1705P F1761S<br>L1657P S1715N M1775E<br>E1660G S1722F M1775K<br>H1686Q F1734L L1780P<br>R1699Q G1738E I1807S<br>K1702E G1743R V1833E<br>Y1703H A1752P A1843T<br>F1704S F1761I |
| **transcription activation** |  | M1652I<br>A1669S | V1665M<br>D1692N<br>G1706A<br>D1733G<br>M1775V<br>P1806A |
| **?** |  |  | M1652T W1718S R1751P C1787S A1823T<br>V1653M T1720A R1751Q G1788D V1833M<br>L1664P W1730S R1758G G1788V W1837R<br>T1685A F1734S L1764P G1803A W1837G<br>T1685I E1735K I1766S V1804D S1841N<br>M1689R V1736A P1771L V1808A A1843P<br>D1692Y G1738R T1773S V1809A T1852S<br>F1695L D1739E P1776S V1809F P1856T<br>V1696L D1739G D1778N V1810G P1859R<br>R1699L D1739Y D1778G Q1811R<br>G1706E V1741G D1778H P1812S<br>W1718C H1746N M1783T N1819S |

"Decision" Tree for Predicting Functional Impact of Genetic Variants

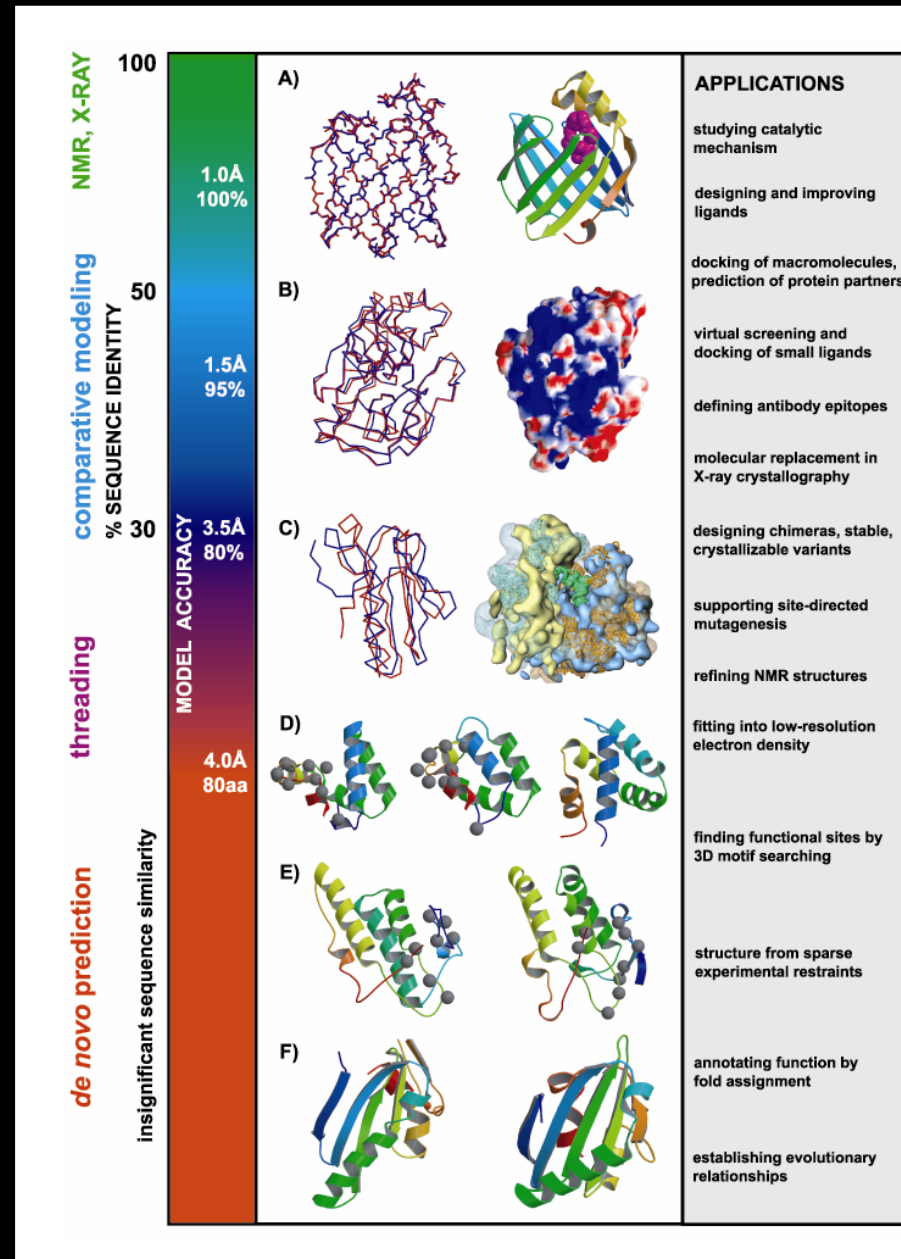# Putative Binding Site on BRCA1



RMSMV**VSG**L**TPEE**FM**LVY**KFARKHHIT**LT**N**LITEETTHVVMKTDAEFVC**E**RTL**K**Y**F**LGIAGGKWVVSYFWVTQSIKERKM
LNEHDFEVRGDVVNGRNHQGPKRARESQDRKIFRGLEICCYGPF**TNM**PTDQLEWMVQLCGASVVKELSSFTLGTGVHPIV
VV**QP**DAWTEDNGFHAIGQMCEAPVVT**RE**WV**L**DSVALYQCQELDTYLIPQIP
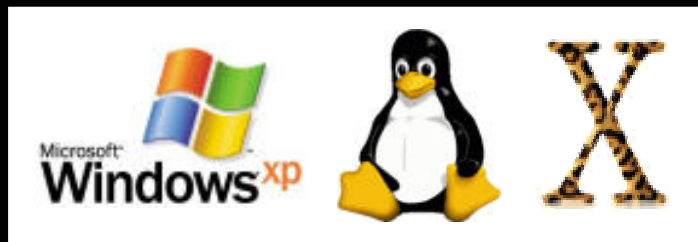
# **Applications** of Comparative Models

# *Summary*

- Protein Structure Prediction and why is it useful?

- Methods in Protein Structure Prediction

- Comparative Modeling

  - ✓ Steps in CM (overview + resources)

  - ✓ Accuracy/Applications of comparative models

  - ✓ Case example in MODELLER

  - ✓ CM and Structural Genomics

BayGenomics

# Obtaining MODELLER and related information

- MODELLER (6v2) web page
- `http://www.salilab.org/modeller/`

  - Download Software (Linux/Windows/Mac)
  - HTML Manual
  - Join Mailing List

# Using MODELLER

- No GUI! ☹
- Controlled by command file (script) ☹☹
- Script is written in TOP language ☹☹☹
- TOP language is simple ☺☺☺☺

BayGenomics

# Using MODELLER

- **INPUT:**
  - Target Sequence (FASTA/PIR format)
  - Template Structure (PDB format)
  - TOP command file

- **OUTPUT:**
  - Target-Template Alignment
  - Model in PDB format
  - Other data

BayGenomics

# Example 1: Modeling of BLBP
## Input

✓ Target: Brain lipid-binding protein (BLBP)

✓ BLBP sequence in PIR (MODELLER) format:

```
>P1;blbp

sequence:blbp:::::::::

VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSI
DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA*
```

- PSI-BLAST template search: Template: PDB file 1HMS:_

BayGenomics

# Example 1: Modeling of BLBP

## STEP 1: Align blbp and 1hms sequences
*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'

SEQUENCE_TO_ALI ALIGN_CODES = '1hms'

READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on

ALIGN

WRITE_ALIGNMENT FILE='blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'

WRITE_ALIGNMENT FILE='blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing `mod align.top` directory where you have the TOP file.
MODELLER will produce a `align.log` file

# Example 1: Modeling of BLBP

**STEP 1: Align blbp and 1hms sequences**

*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'

SEQUENCE_TO_ALI ALIGN_CODES = '1hms'

READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on

ALIGN

WRITE_ALIGNMENT FILE='blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'

WRITE_ALIGNMENT FILE='blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing `mod align.top` directory where you have the TOP file.
MODELLER will produce a `align.log` file

# Example 1: Modeling of BLBP

**STEP 1: Align blbp and 1hms sequences**

*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'

SEQUENCE_TO_ALI ALIGN_CODES = '1hms'

READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on

ALIGN

WRITE_ALIGNMENT FILE 'blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'

WRITE_ALIGNMENT FILE 'blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing `mod align.top` directory where you have the TOP file.
MODELLER will produce a `align.log` file

BayGenomics

# Example 1: Modeling of BLBP

## STEP 1: Align blbp and 1hms sequences
*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'

SEQUENCE_TO_ALI ALIGN_CODES = '1hms'

READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on

ALIGN

WRITE_ALIGNMENT FILE='blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'

WRITE_ALIGNMENT FILE='blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing `mod align.top` directory where you have the TOP file.
MODELLER will produce a `align.log` file

BayGenomics

# Example 1: Modeling of BLBP
## STEP 1: Align blbp and 1hms sequences
## *Output*

```
>P1;1hms

structureX:1hms:   1 : : 131 : :undefined:undefined:-1.00:-1.00

VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIIEKNGDILTLKTHSTFKNTEISFKLGVEFDETTA

DDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIDGKLILTLTHGTAVCTRTYEKE*

>P1;blbp

sequence:blbp:     : :      : : : : 0.00: 0.00

VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSI

DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA*
```

BayGenomics

# Example 1: Modeling of BLBP
## STEP 1: Align blbp and 1hms sequences
### *Output*

```
>P1;1hms

structureX:1hms:    1 : : 131 : :undefined:undefined:-1.00:-1.00

VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIIEKNGDILTLKTHSTFKNTEISFKLGVEFDETTA

DDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIDGKLILTLTHGTAVCTRTYEKE*

>P1;blbp

sequence:blbp:      : :     : : : : 0.00: 0.00

VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSI

DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA*
```

# Example 1: Modeling of BLBP
## STEP 1: Align blbp and 1hms sequences
### *Output*

```
_aln.pos          10        20        30        40        50        60
1hms      VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIIEKNGDILTLKTHSTFKNT
blbp      VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNT
_consrvd ****    ****  ** *** *** *********    **** **    *       *  *****


_aln.pos          70        80        90       100       110       120
1hms      EISFKLGVEFDETTADDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIDGKLILTLTHG
blbp      EINFQLGEEFEETSIDDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFG
_consrvd ** * ** ** **    ***   ** * *** ** * ***** **    **   ***   *** *


_aln.pos         130
1hms      TAVCTRTYEKE
blbp      DIVAVRCYEKA
_consrvd    *   * ***
```

# Example 1: Modeling of BLBP

**STEP 2: Model the blbp structure using the alignment from step 1.**
*TOP script for model building*

```
INCLUDE

SET ALNFILE = 'blbp-1hms.ali'

SET KNOWNS = '1hms'

SET SEQUENCE = 'blbp'

SET STARTING_MODEL = 1

SET ENDING_MODEL = 1

CALL ROUTINE = 'model'
```

Run by typing `mod model.top.`
Check file `model.log`

# Example 1: Modeling of BLBP

**STEP 2: Model the blbp structure using the alignment from step 1.**

*TOP script for model building*

```
INCLUDE

SET ALNFILE = 'blbp-1hms.ali'

SET KNOWNS = '1hms'

SET SEQUENCE = 'blbp'

SET STARTING_MODEL = 1

SET ENDING_MODEL = 1

CALL ROUTINE = 'model'
```

Run by typing `mod model.top`.
Check file `model.log`

BayGenomics

# Example 1: Modeling of BLBP

**STEP 2: Model the blbp structure using the alignment from step 1.**

*TOP script for model building*

```
INCLUDE

SET ALNFILE = 'blbp-1hms.ali'

SET KNOWNS = '1hms'

SET SEQUENCE = 'blbp'

SET STARTING_MODEL = 1

SET ENDING_MODEL = 1

CALL ROUTINE = 'model'
```

Run by typing `mod model.top`.
Check file `model.log`

# Example 1: Modeling of BLBP

**STEP 2: Model the blbp structure using the alignment from step 1.**
*Output coordinates file*

## Model file → `blbp.B99990001`

- PDB file

- Can be viewed with Chimera

[http://www.cgl.ucsf.edu/chimera/]

Rasmol

[http://www.bernstein-plus-sons.com/software/rasmol/]

**What is the physiological ligand of Brain Lipid-Binding Protein?**

**Predicting features of a model that are not present in the template**

BLBP/oleic acid

BLBP/Docosahexaenoic acid

*Cavity is **not** filled*

*Cavity **is** filled*



Ligand binding cavity

1. BLBP binds fatty acids.

2. Build a 3D model.

3. Find the fatty acid that fits most snuggly into the ligand binding cavity.

L. Xu, R. Sánchez, A. Šali, N. Heintz, *J. Biol. Chem.* **271**, 24711, 1996.

BayGenomics

# 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 ~16,000 sequence families (30% sequence id).

- **Solution:**
  - Determine many protein structures.
  - Increase modeling distance.

Šali. *Nat. Struct. Biol.* **5**, 1029, 1998.
Šali & Kuriyan. *TIBS* **22**, M20, 1999.

Burley *et al. Nat. Genet.* **23**, 151, 1999.
Sanchez *et al. Nat. Str. Biol.* 7, 986, 2000

BayGenomics

# How can **Comparative Modeling** be used in **Structural Genomics**?

- **Target Selection**

  How many structures need to be solved?
  Which structures should we solve first?

- **Target Amplification**
  How much of the sequence space is covered by:
  a new structure
  all structures

**MODPIPE: Large-Scale Comparative Protein Structure Modeling**

START

**PSI-BLAST**

Prepare PSI-BLAST PSSM by comparing the sequence against the NR database of sequences

Use the sequence PSSM to search against the representative set of PDB chains (F and no-F)

Use the PDB chain PSSMs to search against the sequence (F and no-F)

Select Templates using a permissive E-value cutoff

1

**MODELLER**

1

Align the matched part of the target sequence with the template structure

Build a model for the target segment by satisfaction of spatial restraints

Evaluate the model

For each template

For each sequence

END

BayGenomics

# Modeling Coverage Of The Sequence Space

Not Attempted
42.4%

Reliable Model +
Fold Assignment
43.9%

Attempted
1.0%

Reliable Model Only
0.2%

Fold Assignment
Only
12.6%

**Fold assignment:** **PSI-BLAST E-value $\leq$ 1e$^{-4}$**
**Reliable Model:** **Model Score $\geq$ 0.7**

BayGenomics

# Comparative Models for TrEMBL Sequences

**Target-Template Sequence Identity**

**Model Length**

**E-values of Alignment**

**Coverage of SCOP domains**

**ModBase**: A database of comparative protein structure models and properties.
`http://www.salilab.org/modbase`

# Mod Web

## A Server for Comparative Protein Structure Modeling
http://www.salilab.org/modweb

# Conclusions

✓ Comparative models help to understand protein's function:
  - ✓ Detecting remote structural (functional?) relationships.
  - ✓ Revealing features that are not present in the templates.
  - ✓ Revealing features that are not recognizable from the sequence.

✓ Currently, useful 3D models can be obtained for domains in approximately 50% of the proteins (30% of domains), because of the improved **methods** and because of the many **known protein structures** and **sequences**.

✓ We will be able to calculate useful models for most globular domains soon after the completion of the genome projects, because of **structural genomics**.

# Acknowledgments



Andrej Sali

Frank Alber
Fred Davis
Damien Devos
Narayanan Eswar
Rachel Karchin
Libusha Kelly
Michael F. Kim
Dmitry Korkin
M. S. Madhusudhan
Nebosja Mirkovic
Ursula Pieper
Andrea Rossi
Min-yi Shen
Maya Topf

`http://www.salilab.org`