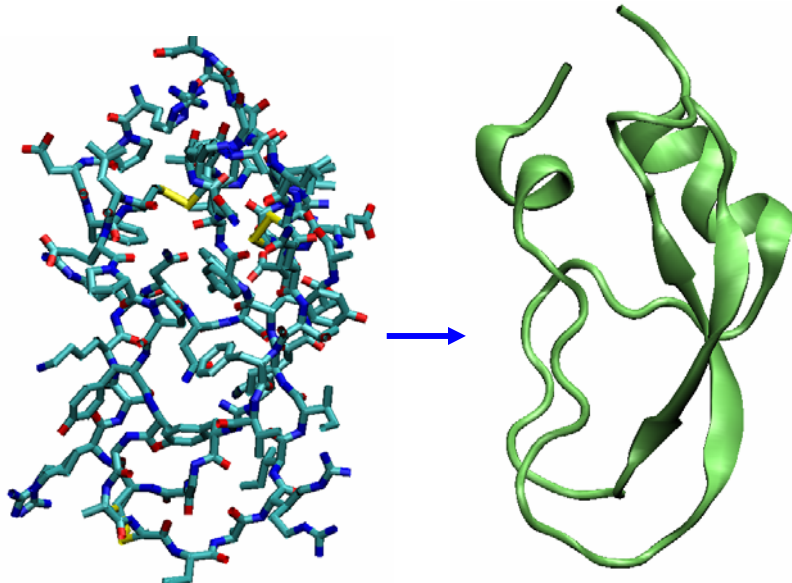# Bioinformatics

Use of computer to analyze and archive biological data (sequence and structural information) on a large scale

– includes development of analysis algorithms, visualization software, database design

- Secondary structure assignment
- Secondary structure prediction
- Sequence alignment
- Structural alignment
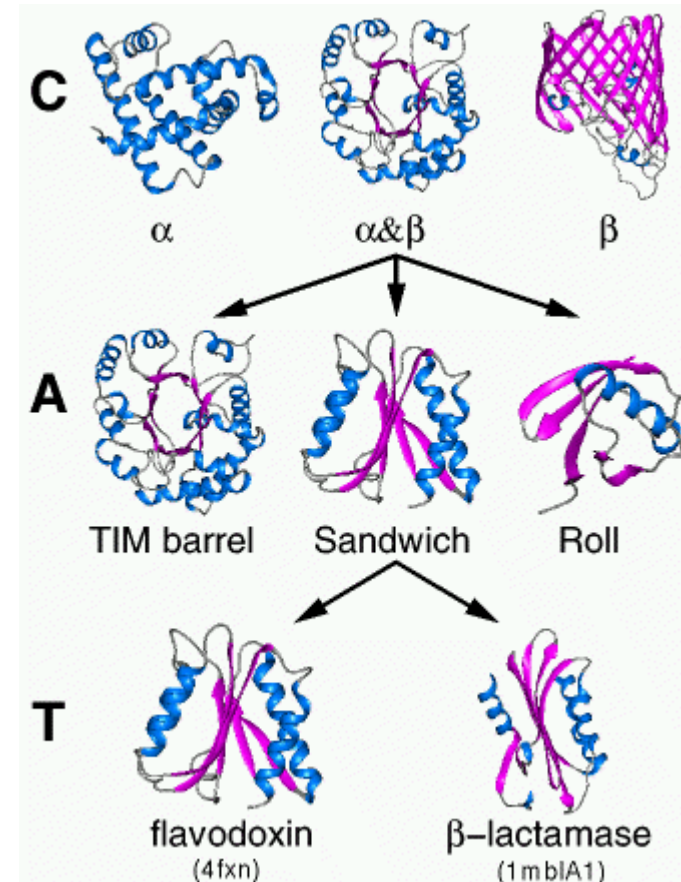- Tertiary structure prediction

# Secondary structure assignment

Easy visualization

Structural classification



Detection of structural motifs and improved sequence-structure searches

Structural alignment

# Given a structure, identify the regions of secondary structure

– DSSP, Stride, DEFINE

– implementation dependent

# Secondary structure prediction

Tertiary structure prediction from the amino acid sequence is very difficult

Prediction of secondary structure is feasible and more reliable

In some models of protein folding, secondary structural elements form first before a tertiary structure is formed



A-C-H-Y-T-T-E-K-R-G-G-S-G-T-K-K-R-E-A

H-H-H-H-H-H-H-H-C-C-C-C-C-S-S-S-S-S-S

Knowing the region of secondary structure is critical for some applications
– transmembrane domain of the membrane protein GPCR
– secondary structural info may be sufficient for some studies

# Prediction methods

Use known secondary structure propensities of individual amino acids—either statistical or experimental
- – helix former, helix breaker, helix neutral, sheet former, sheet breaker, etc
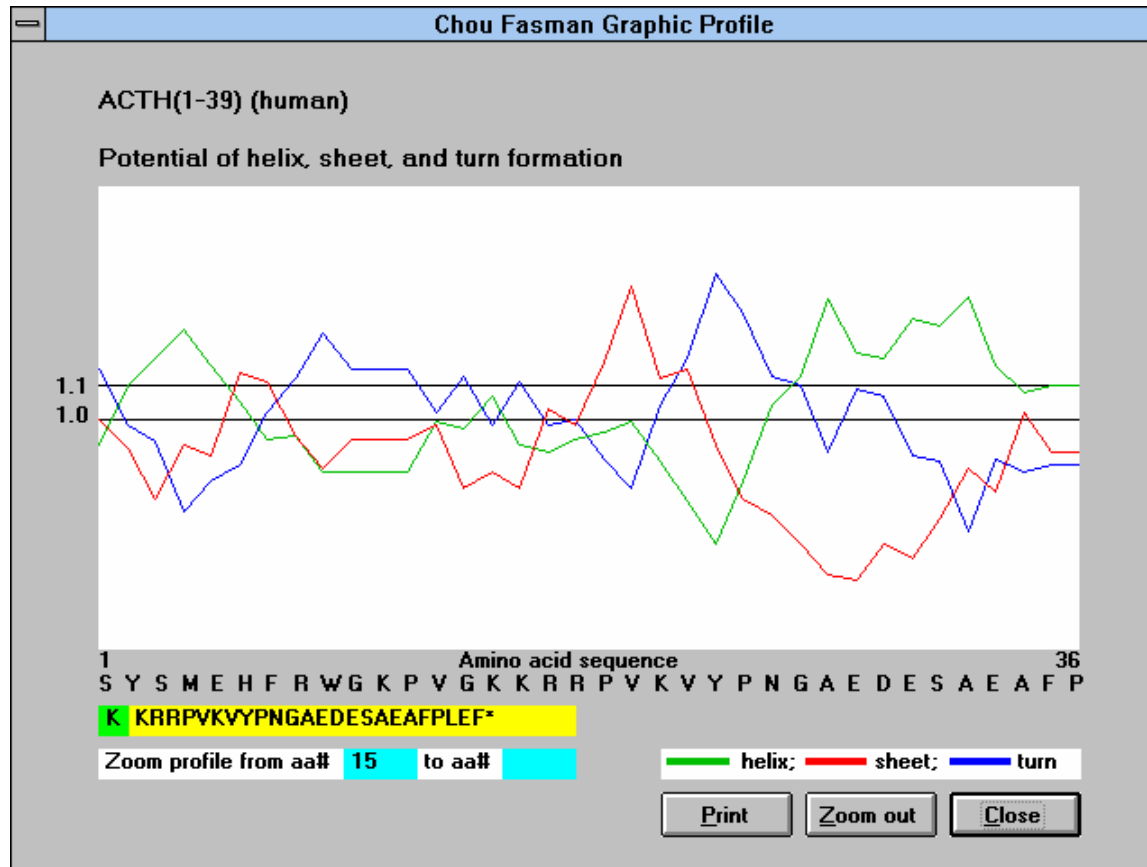- – develop heuristic rules for identifying and extending a helix or a sheet

## Chou Fasman Graphic Profile

ACTH(1-39) (human)

Potential of helix, sheet, and turn formation

1.1
1.0

1                          Amino acid sequence                          36
S Y S M E H F R W G K P V G K K R R P V K V Y P N G A E D E S A E A F P

K KRRPVKVYPNGAEDESAEAFPLEF*

Zoom profile from aa#  15  to aa#

━━━ helix;  ━━━ sheet;  ━━━ turn

Print    Zoom out    Close

Examine sets of adjacent amino acids (e.g. windows of 11-21 amino acids) rather than individual amino acids
- probability of an amino acid to be in a particular secondary structure considering the nearby residues
- HHSLCSHSHHSC less likely than HHHHHCCCSSSS
- local context is important

Secondary structure prediction services
- PredictProtein
- PHD
- JPRED
- PSSP

Limitations
- overall prediction: 60%
- beta-strands prediction: ~35%
- predictions include small secondary elements that cannot be easily integrated into longer structures

# Sequence alignment

Process of comparing two or more sequences by looking for a series of individual characters or character patterns (similar vs. identical) that are in the same order in the sequences

Sequence alignment lies at the heart of bioinformatics
– newly discovered sequence may be related to known sequence
– models evolutionary relationship
– assist in engineering and 3D prediction
– basis to functional genomics
– population genomics—genetic variations in an isolated group (DeCode).

Identity vs. similarity – definition of similarity

http://www.ncbi.nlm.nih.gov/

http://www.ebi.ac.uk/

http://www.expasy.ch/

**National Center for**
**Biotechnology**
**Information**
**National Library of**
**Medicine**
**National Institute of**
**Healtlh**
USA

**Eouropean**
**Bioinformatics Institute**
**European Molecular**
**Biology Laboratory**
UK

**Ex**pert **P**rotein **A**nalysis **Sy**stem proteomics server of the Swiss Institute of Bioinformatics

# Multiple sequence alignment (MSA)

Incorporate evolutionary information through multiple sequence alignment

- information on sequence conservation, substitution, and potential interaction
- ClustalW
- T-Coffee



| | Global multiple alignment | Local multiple alignment |

# Structural alignment

Structures are more conserved than sequences

In the "twilight zone" of sequence similarity, structural alignment might help to correctly determine the relations between two proteins

Structural alignment is more predictive of function than sequence alignment

sequence 1

sequence 2

similar local structure

# Alignment v. superposition

Superposition assumes the two are related—translate and rotate one of them to minimize the total rmsd

Alignment is a means of determining if two are structurally related by mapping stretches of atoms from one protein to another
- integral to structural classification

- Distance alignment matrix (DALI)
- Combinatorial extension (CE)
- Sequential Structure Alignment Program (SSAP)
- Spatial Arrangements of Backbone Fragments (SARF2)
- Structural Alignment of Multiple Proteins (STAMP)
- Structure based Alignment Program (STRAP)

Many are available as web services

(a)

C

helix

N

helix

loop

(b)

N

helix

helix

loop

C

©1999 GARLAND PUBLISHING INC.
A member of the Taylor & Francis Group

SARF: pair of secondary structure

3D-BLAST: unusual definitions

CE: longest path of aligned fragment pairs

Protein B

$D_{ij}$

$p_i^B$

$p_j^B$

$d_{p_i^B+k, p_j^B+l}^B$

AFP $i$

AFP $j$

$p_i^A$

$d_{p_i^A+k, p_j^A+l}^A$

Protein A

A   3D Structure Database

c

b

a

1brb_I

c'

b'

a'

1bf0

· · ·

B

$C_\alpha^{i-1}$

$C_\alpha^{i-2}$

$C_\alpha^i$

$\kappa$

$C_\alpha^{i+1}$

$\alpha$

$C_\alpha^{i+2}$

C

1brb_I:

T K T E N

C B Y Y Y

180

120

Kappa

60

0

-180   -90   0   90   180

Alpha

D   Structural Alphabet Database

1brb_I          a              b

TKTENKQPXTN KHKHKKEMD SQTNKHKNF KMMPFK
IQPKVP CBYYYY

1bf0        c              a'          b'

PXIGDQTKHEHKXTGPF KEKEKKF GLSQ TNKXNFF
KQVQNKQRPFVT IGGGGD DLLS
· · ·          c'

E          a              b                  c

SCORE: 6646261121116262661218666601166122161371821111112
1brb_I: TKTENKQPXTN KHKHKKEMD SQTNKHKNF KMMPFKIQPKVP CBYYYY
        TK E+K + + K K+KK+++SQTNK + FK + K +P+V+++ +
1bf0: TKHEHKXTGPF KEKEKKFGL SQTNKXNFF KQVQNKQRPFVT IGGGGD
        a'              b'                  c'

F

c'

c

b

a
a'

# Tertiary structure prediction

- Detailed structural information is essential to model function and to design methods to modulate function
- Experimentally determined structures are used as templates during structural prediction



Stevens et al, Science 294, 89-93 (2001)

Stevens & Wilson, Science 291, 519 (2001)

Baker & Sali, Science 294, 93 (2001)

# CASP



Critical Assessment of (Protein) Structure Prediction

Bi-annual competition for testing the current state of structure prediction capabilities

Contestants are given protein sequence and need to submit model structures to be compared against experimental structures

No limit on the technique

---

Judging the success of a prediction -- Global and local rmsd

---

Would like it to be high throughput to cover the full genome

A lot of experimental information cannot be modeled in high throughput, e.g. thermostability and functional site residues

Lack of resolution prevents mutagenesis data, information regarding solvent accessibility (e.g. H/D exchange, fluorescence) to be properly modeled

Domain arrangements (quaternary structure) are also difficult to model

# *Protein Structure Prediction Center*

## Genome Center
## University of California, Davis

## Welcome to the Protein Structure Prediction Center!

Our goal is to help advance the methods of identifying protein structure from sequence. The Center has been organized to provide the means of objective testing of these methods via the process of blind prediction. In addition to support of the CASP meetings our goal is to promote an evaluation of prediction methods on a continuing basis.

CASP experiments aim at establishing the current state of the art in protein structure prediction, identifying what progress has been made, and highlighting where future effort may be most productively focused. The organizers are thankful to CASP assessors for their valuable contribution to this field.

There have been seven previous CASP experiments.

**CASP1 (1994)** | **CASP2 (1996)** | **CASP3 (1998)** | **CASP4 (2000)** | **CASP5 (2002)** | **CASP6 (2004)** | **CASP7 (2006)**

**Proceedings**

**Click on the logo below to proceed to the main page of the latest CASP experiment.**

C
A
S
P
7

*7th Community Wide Experiment on the*

# Critical Assessment of Techniques for Protein Structure Prediction

*Asilomar Conference Center, Pacific Grove, CA*
*November 26-30, 2006*

Sponsored by the US National Library of Medicine (NIH/NLM), National Institute of General Medical Sciences (NIH/NIGMS)

Co-sponsored by: BioSapiens Network of Excellence, *hp invent*

| Targets | Predictions | Meeting | 3D Evaluation Results |
|---|---|---|---|
| CASP7 Target List<br>Refinement Target List<br>Domain definition<br>Domain classification<br>Prediction success charts | Categories of predictions<br>Server Predictions | Abstracts<br>Meeting Program<br>Meeting participants<br>GROUPS: by name  by number | **Target Perspective View**<br>Group Perspective View<br>Table Browser<br>Refinement Results |
| Thank you, experimentalists | CASP7 in numbers | Presentations | Results Page Description |

**For CASP7 raw data archives go to the Downloads Area of our main page or click here**

## Description of the experiment

Goals   Scope   Related   Timetable   Participation   Targets   Format   Assessment   Results   Meeting   Organizers

Targets: ranking by Target Number - Mozilla Firefox

File  Edit  View  Go  Bookmarks  Tools  Help

http://predictioncenter.gc.ucdavis.edu/casp7/targets/cgi/casp7-view.cgi?loc=predictioncenter.org;page=casp7/    Go  G

Buffalo  Finance  NCBI  ISI  Dictionary  Excite  Science  Journals  Temp  Vendors  UIUC

Search  NCBI    PubMed    for

You have been registered — The BALL Web...    BALLView — The BALL Website    Targets: ranking by Target Number

# Seventh Community Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction

# Targets: ranking by Target Number

Targets expire on specified date at noon (12:00) local time in California (GMT - 7 hours).
If information leak occurs after the three weeks since target release, evaluation will be limited to the models submitted within the initial 3 weeks only.

Order by Target Number    Reorder

Click here for help.

| Tar-id | Name | Nres | Method | Entry-date | Expiry-date/cancel | Description |
|--------|------|------|--------|------------|--------------------|-------------|
| T0283 | BH3980 | 112 | X-ray | 10 May | 21 Jun | JCSG target 10176605, Bacillus halodurans  PDB code 2HH6 |
| T0284 | PA4872 | 287 | X-ray | 11 May | 1 Jun | Hypothetical protein, Pseudomonas aeruginosa PAO1 |
| T0285 | AbfS | 125 | X-ray | 15 May | 26 Jun | Extracytoplasmic domain from histidine kinase, Cellvibrio japonicus |
| T0286 | CelX | 205 | X-ray | 15 May | 26 Jun | Cellulose esterase family, Clostridium thermocellum |
| T0287 | CagS | 199 | X-ray | 11 May | 1 Jun | HP0534, cag pathogenicity island protein, Helicobacter pylori |
| T0288 | PRKCAB | 93 | X-ray | 16 May | 9 Jun | SGC target PRKCAB, Homo sapiens  PDB code 2GZV |
| T0289 | AAH7881 | 312 | X-ray | 16 May | 9 Jun | CESG target, Rattus norvegicus  PDB code 2GU2 |
| T0290 | PPI64 | 173 | X-ray | 17 May | 10 Jun | SGC target PPI64, Homo sapiens  PDB code 2GW2 |
| T0291 | EPHA3 | 310 | X-ray | 17 May | 11 Jun | SGC target EPHA3, Homo sapiens  PDB code 2GSF |
| T0292 | NEK2A | 277 | X-ray | 17 May | 11 Jun | SGC target NEK2A, Homo sapiens  PDB code 2CL1 |
| T0293 | MGC3329 | 250 | X-ray | 18 May | 12 Jun | SGC target MGC3329, Homo sapiens  PDB code 2H00 |

CASP7 Target T0283 - Mozilla Firefox

File   Edit   View   Go   Bookmarks   Tools   Help

http://predictioncenter.gc.ucdavis.edu/casp7/targets/templates/t0283.doc.html          Go

Buffalo   Finance   NCBI   ISI   Dictionary   Excite   Science   Journals   Temp   Vendors   UIUC

Search   NCBI       PubMed       for

You have been registered — The BALL Web...   BALLView — The BALL Website   CASP7 Target T0283

# CASP7 Target T0283

**1. Protein Name**
  BH3980
**2. Organism Name**
  Bacillus halodurans
**3. Number of amino acids (approx)**
  112
**4. Accession number**
  10176605
**5. Sequence Database**
  NCBI NR
**6. Amino acid sequence**
  MSFIEKMIGSLNDKREWKAMEARAKALPKEYHHAYKAIQKYMWTSGGPTDWQDTKRIFGG
  ILDLFEEGAAEGKKVTDLTGEDVAAFCDELMKDTKTWMDKYRTKLNDSIGRD
**7. Additional information**
  DUF1048, more info available at http://www1.jcsg.org/cgi-bin/psat/analyzer.cgi?acc=10176605
**8. X-ray structure**
  yes
**9. Current state of the experimental work**
  refined model
**10. Interpretable map?**
  yes
**11. Estimated date of chain tracing completion**
  complete
**12. Estimated date of public release of structure**
  July

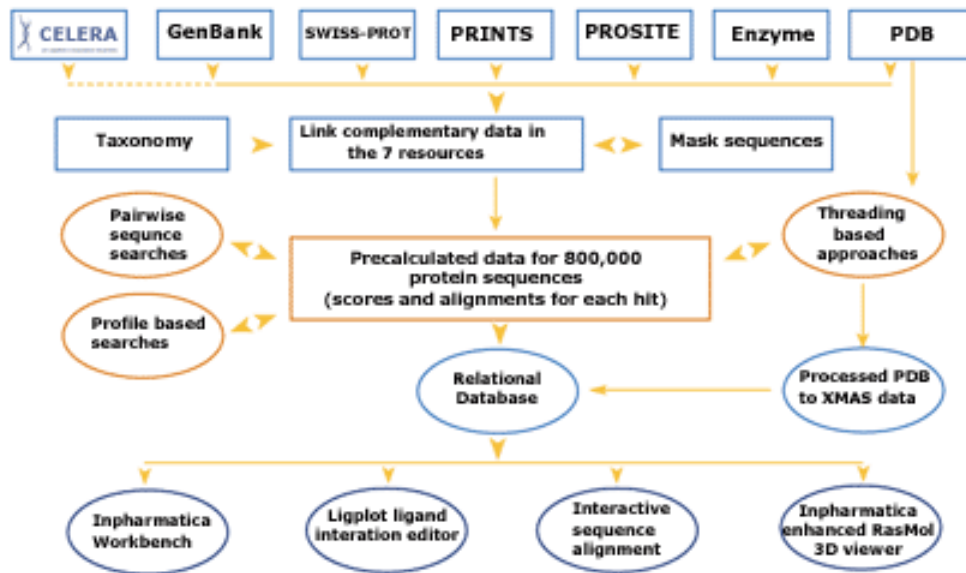## Related Files

Template Sequence file

Template PDB file

# Biopendium™

inpharmatica

http//:www.inpharmatica.com



## Biopendium™



| Family Annotation | Sequence Data | 3D Data | Supporting Data | Derived Information |
|---|---|---|---|---|
| Prints  pfam  SCOP  Prosite | Genbank  Swissprot  PIR | PDB  Cleaned & Processed PDB data | Taxonomy  Enzyme | Signal peptides – von Heijne  Membrane regions - Memsat  Coiled coils - Lupas  Secondary structure – Kabsch-Sander  Accessibility - Lee & Richards  H-bonds/ hydrophobics  Ligand Interactions – Ligplot  Roman Laskowski/ Janet Thornton |

Domain Professor — Profile based domain annotation

Blast Pairwise sequence search

iPSI-Blast Profile based searches

Genome Threader — Profile based Threading — David Jones
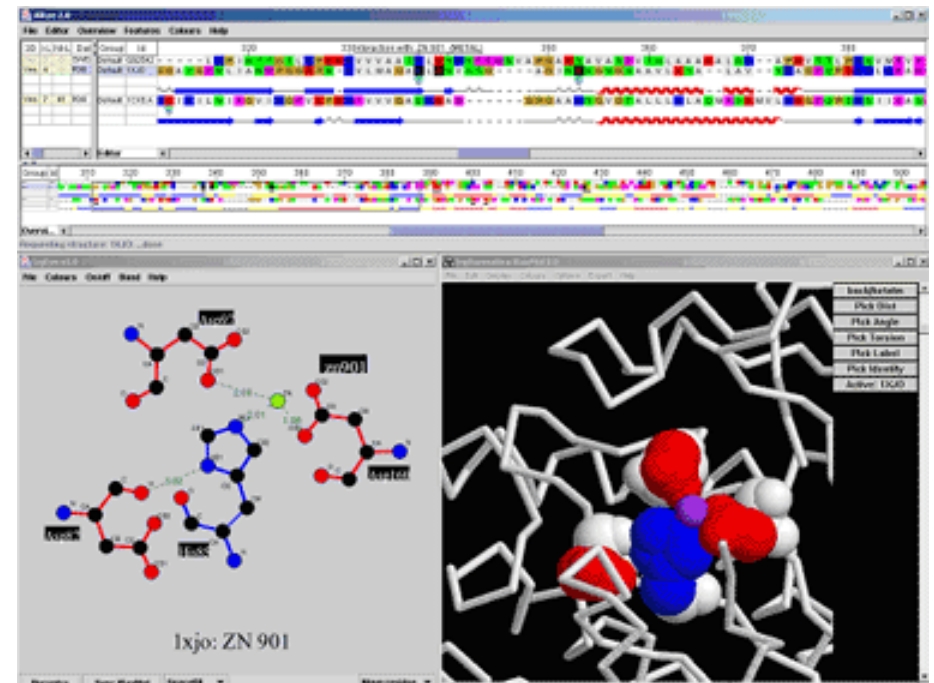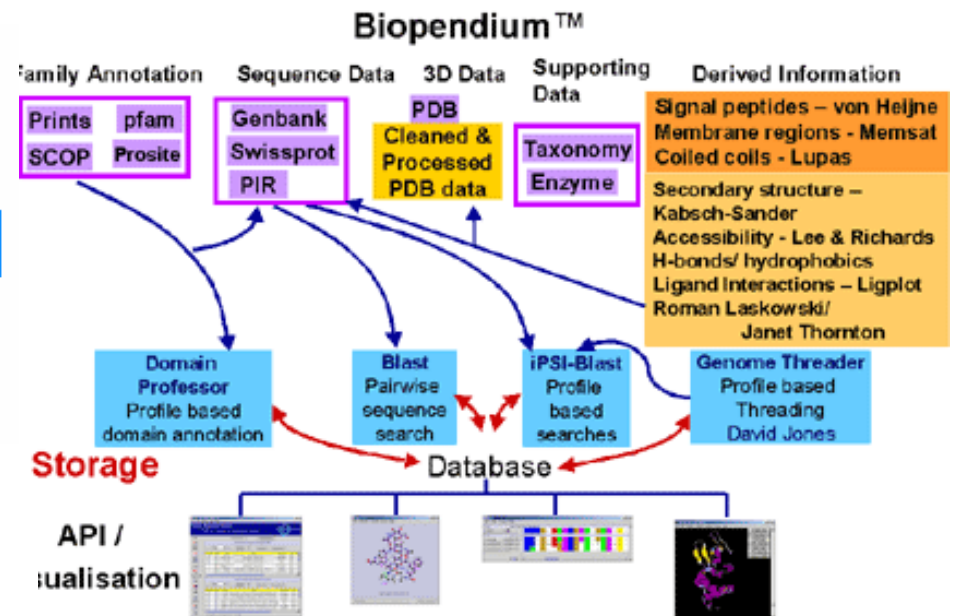
**Storage**

Database

**API / Visualisation**

The Biopendium™*

Celera Edition Biopendium™ with Celera Human and Mouse Data
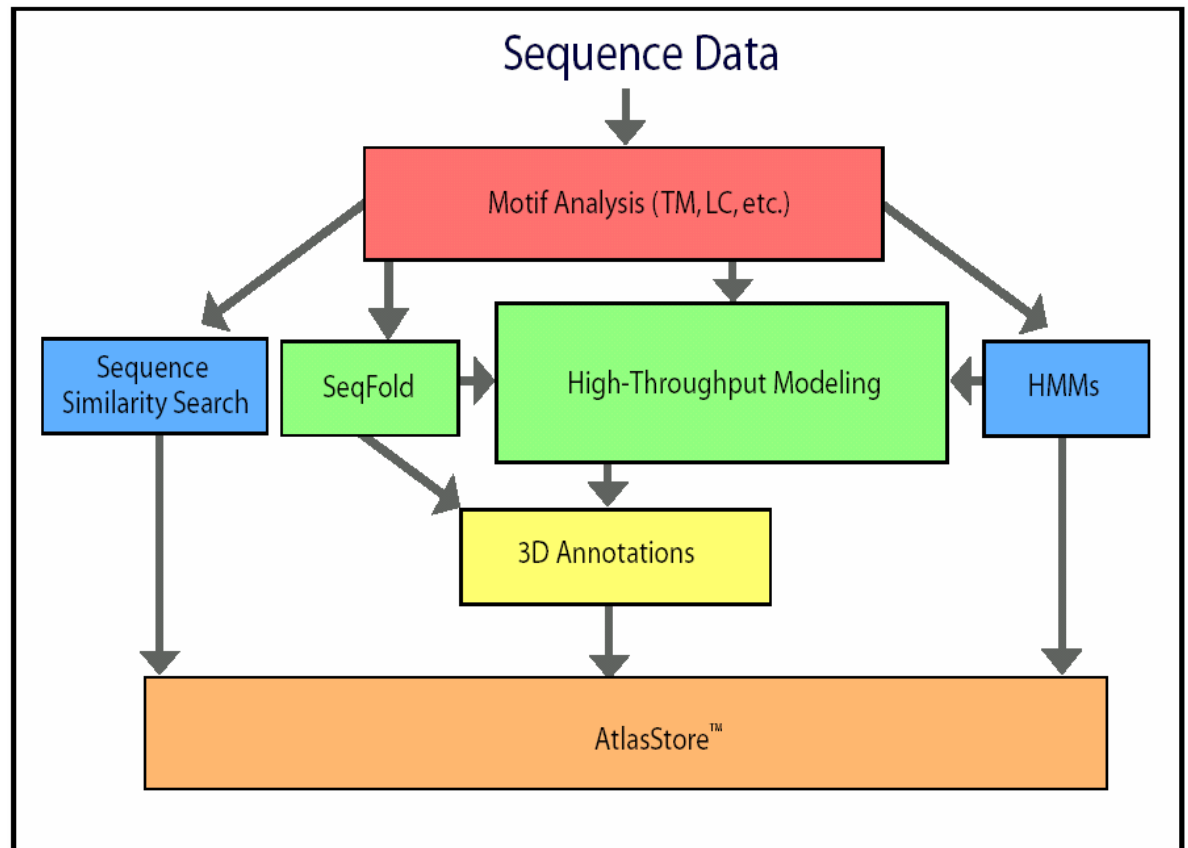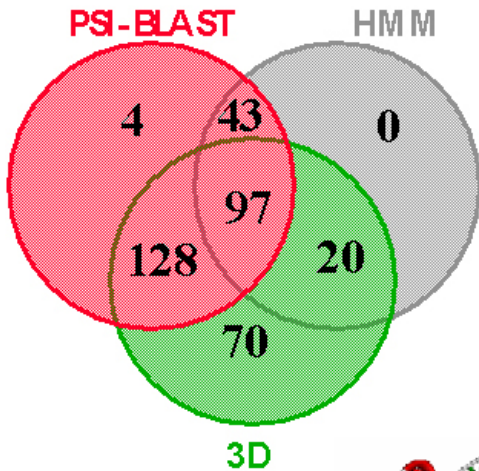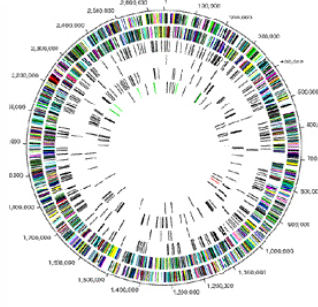


*Available with or without Celera Data



1xjo: ZN 901
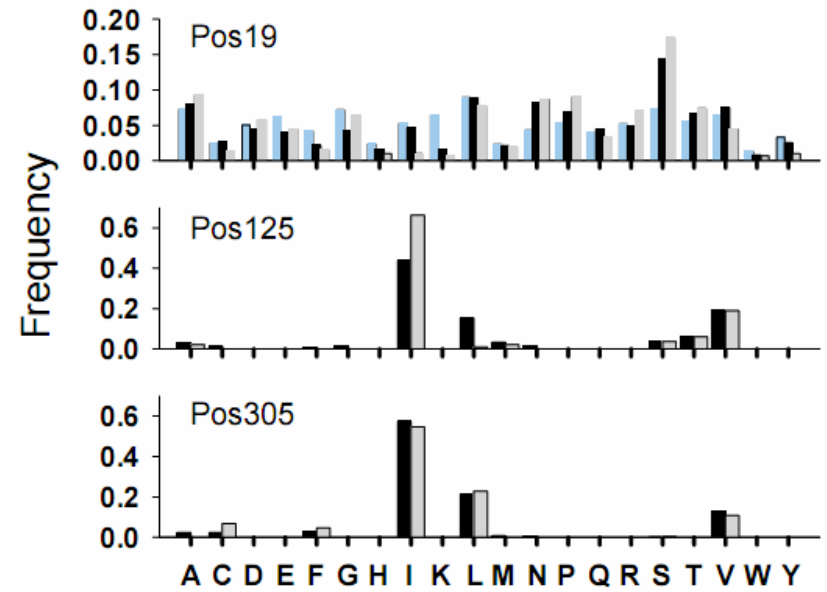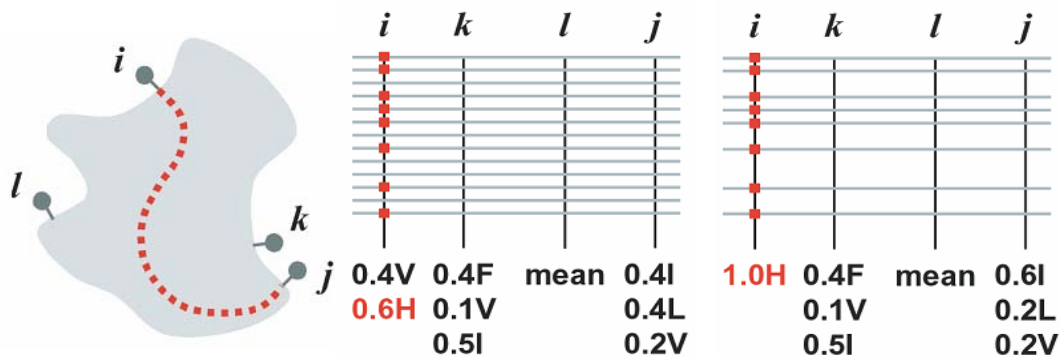
# Accelerating Target Discovery

## Genomes Analyzed with GeneAtlas™

| | |
|---|---|
| *Mycoplasma genitalium* | *Mycobacterium tuberculosis* |
| *Saccharomyces cerevisiae* | *Aquifex aeolicus* |
| *Escherichia coli* | *Haemophilus influenzae Rd* |
| *Caenorhabditis elegans* | *Methanococcus jannaschii* |
| *Borrelia burgdorferi* | *Synechocystis sp.* |
| *Rickettsia prowazekii* | *Bacillus subtilis* |
| *Mycoplasma pneumoniae* | *Helicobacter pylori* |
| *Archaeoglobus fulgidus* | *Pyrococcus horikoshii* |
| *Methanobacterium thermoautotrophicum* | *Treponema pallidum* |
| *Chlamydia trachomatis* | *Deinococcus radiodurans* |
| *Arabidopsis thaliana (partial)* | *Drosophila melanogaster (partial)* |
| *Homo sapiens (partial)* | *Vibrio cholerae* |



accelrys

HYSEQ INC.

GEMA Biotech

GENAISSANCE PHARMACEUTICALS

eXeGeNICS

### PSI-BLAST / HMM / 3D



4    43    0

97

128    20

70

### Sequence Data



Motif Analysis (TM, LC, etc.)

Sequence Similarity Search

SeqFold

High-Throughput Modeling

HMMs

3D Annotations

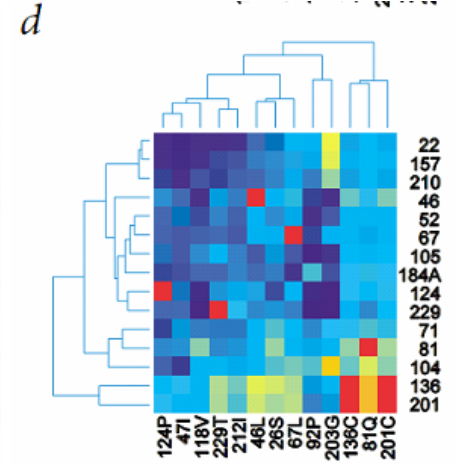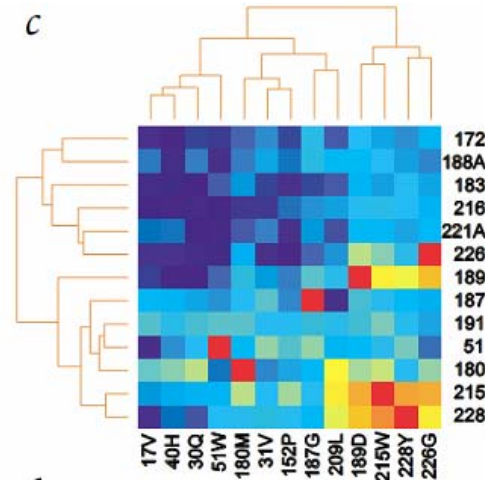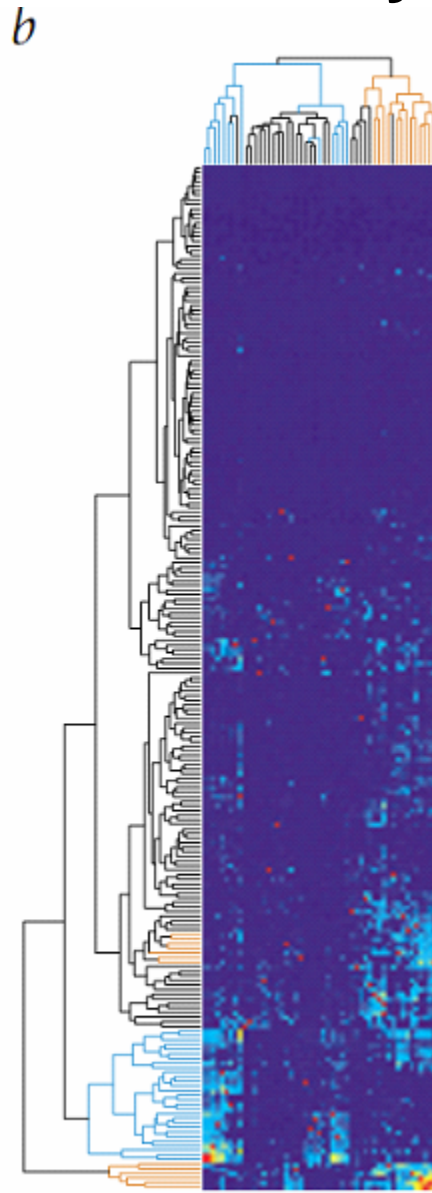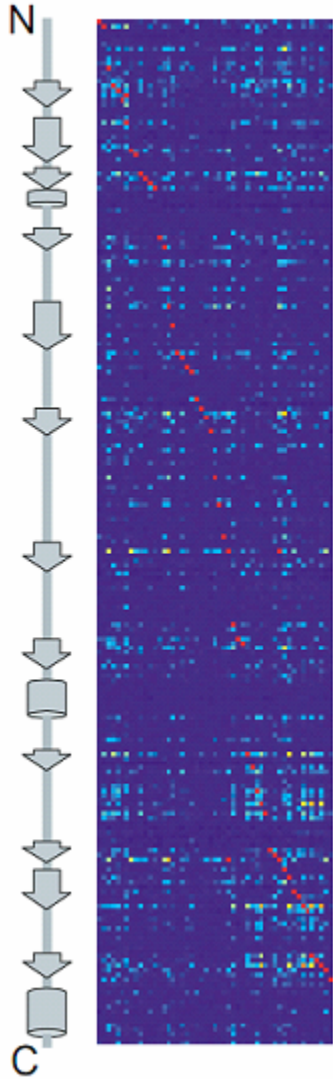AtlasStore™

# Bioinformatics and protein engineering

- Information required for specifying the tertiary structure is contained in the amino acid sequence
- Can we extract the information and use it to specify a protein fold?
- Use statistical information encoded in a multiple sequence alignment

Hypothesis: structural coupled residues would appear more often together than statistically expected



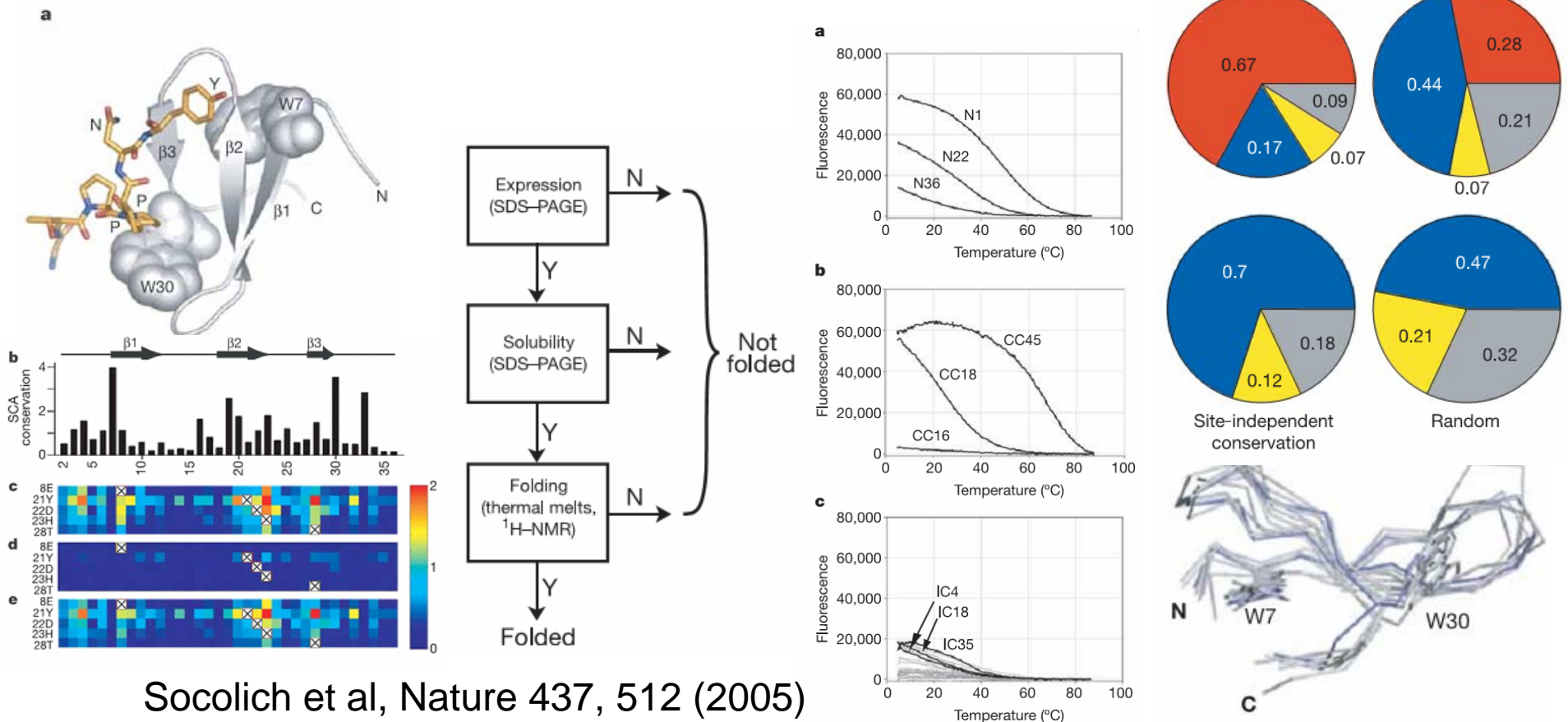Suel et al, NSB (2003)

# Chymotrypsins

# Designing a fold from sequence conservation

Apply statistical analysis to 120 WW domain proteins to identify which residues are structurally coupled

Using simulated annealing Monte Carlo, design sequences that reproduces
  i)  intrinsic amino acid distribution at each position, or
  ii) both the sequence conservation and statistical coupling



Socolich et al, Nature 437, 512 (2005)