

Introduction to Protein Structure

- 1-D world of nucleotide structure and amino acid sequences

➔ now enter to ➔

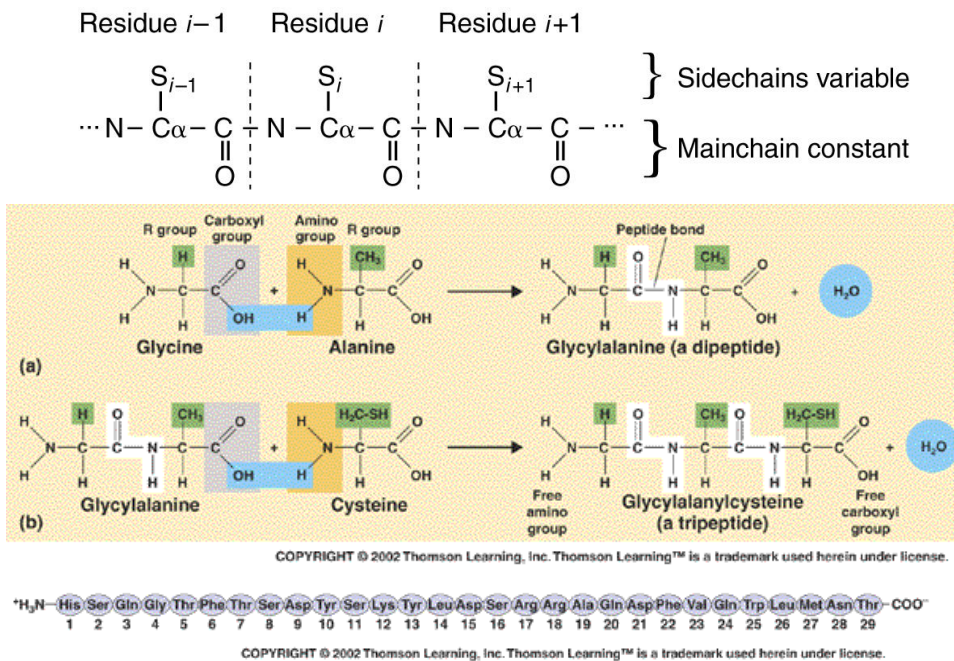
- 3-D world of molecular structures

Proteins play a variety of roles in life process

- Structural proteins
- Enzymes: proteins that catalyze (催化) chemical reactions
- Transport and storage proteins
- Regulatory proteins
- Proteins that control gene transcription
- Proteins that involved in recognition, including cell adhesion (黏著) molecules,
- Antibodies and other protein of the immune system

- Proteins are *large molecules*.
- In many cases only a small part of the structure – an *active site* – is directly functional, the rest existing primarily to create and fix the spatial relationship among the active site residues.
- Proteins evolve by *structural changes*, produced by *mutations* in the amino acid sequence and *genetic rearrangements*, that bring together different combinations of structural subunits.

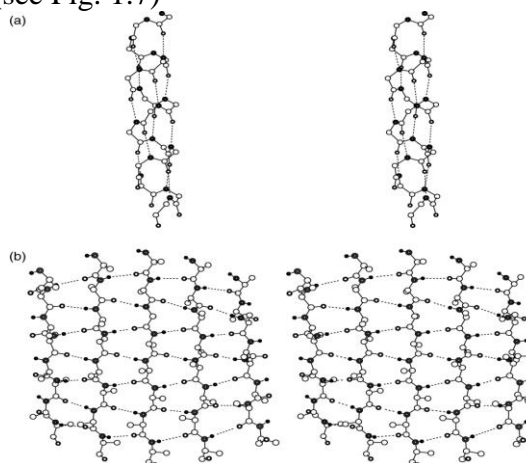
- ~ *85,000* protein structures are now known
- Most were determined by *X-ray crystallography* or *NMR* (nuclear magnetic resonance)
- Few were determined by electron microscopy and others
- Chemically, protein molecules are long polymers typically containing several thousand atoms, composed of a uniform repetitive *backbone* (or *mainchain*) with a particular *sidechain* attached to each residue (see Fig. 1.6)
- Amino acid sequence of a protein records the succession of sidechains.



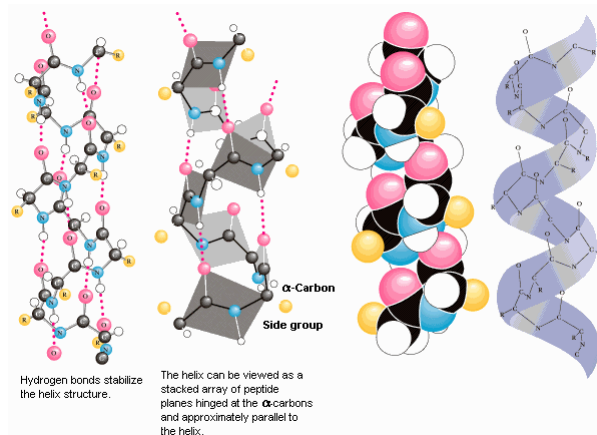
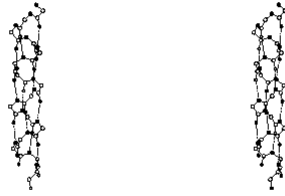
- The polypeptide chain folds into a curve in space
- The course of the chain defining a **folding pattern**
- A great variety of folding patterns: a number of common structural features
- α helices and β sheets (see Fig. 1.7)

- 螺旋 (helix)
- 摺板 (sheet)

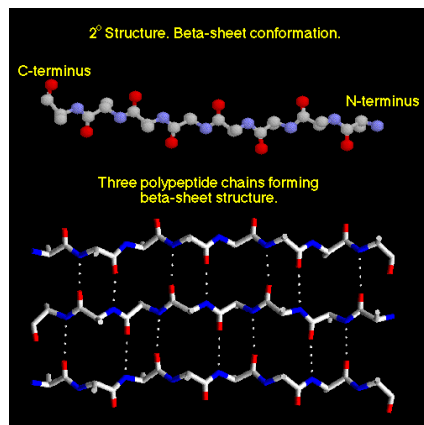
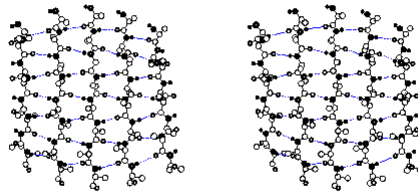
- Folding may be thought of as a kind of intramolecular condensation or crystallization



α helix



β sheet



Hierarchical nature of protein architecture

- **Primary structure**: the amino acid sequence – the set of primary chemical bonds
- **Secondary structure**: the assignment of helices and sheets – the hydrogen-bonding pattern of the mainchain
- **Tertiary structure**: the assembly and interactions of the helices and sheets
- **Quaternary structure**: for proteins composed of more than one subunit, the assembly of the monomers (單體)

Additional levels to the hierarchy

- **Supersecondary structures**: include the alpha-helix hairpin, the beta-hairpin, and the beta-alpha-beta unit. (Fig. 1.8)
- **Domains**: many proteins contain compact units within the folding pattern of a single chain, that look as if they should have independent stability. (Fig. 1.9)
- **Modular proteins**: are multidomain proteins which often contain many copies of closely related domains.
 - Domain recur in many proteins in different structural contexts; that is, different modular proteins can ‘mix and match’ sets of domains.

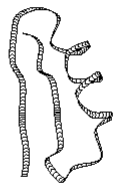
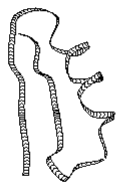
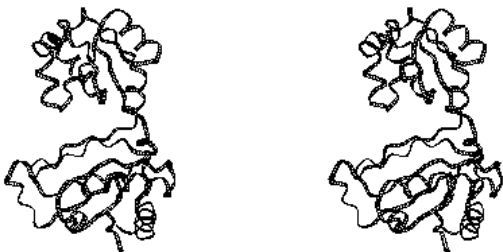
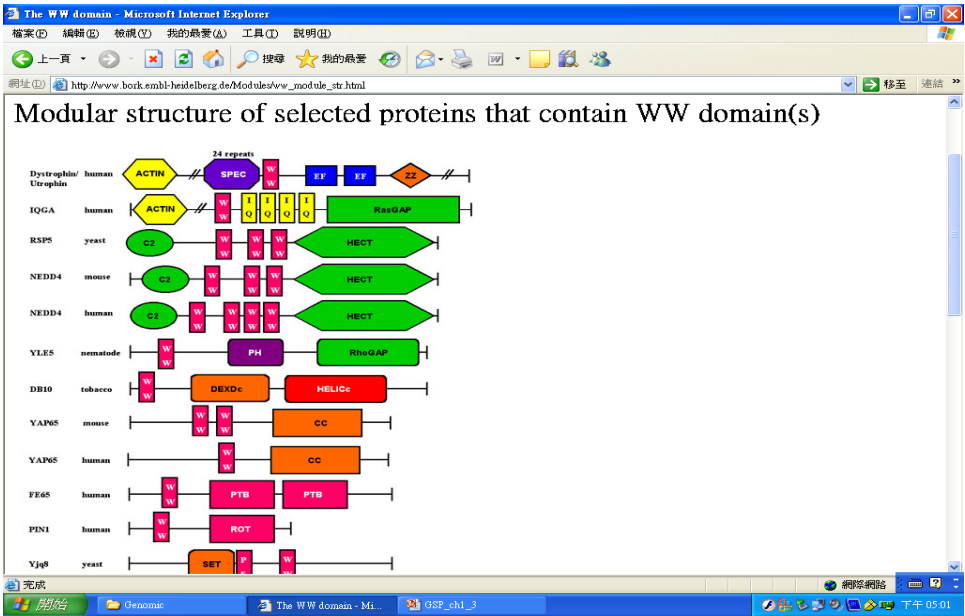


Fig. 1.9 RNA binding protein L1:



Multidomain proteins



Classification of protein structures

- The most general classification of families of protein structures is based on the *secondary and tertiary* structures
- Classification of protein structures occupies a key position in bioinformatics, not least as a bridge between sequence and function.

Class	Characteristic
α -helical	secondary structure exclusively or almost exclusively α -helical
β -sheet	secondary structure exclusively or almost exclusively β -sheet
$\alpha + \beta$	α -helices and β -sheets separated in different parts of the molecule; absence of β - α - β supersecondary structure
α/β	helices and sheets assembled from β - α - β units
α/β -linear	line through centres of strands of sheet roughly linear
α/β -barrels	line through centres of strands of sheet roughly circular
little or no secondary structure	

Fig 1-10a: engrailed homeodomain [1enh]:

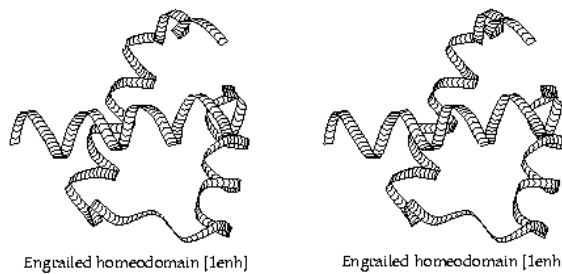


Fig 1-10b: second calponin homology domain from utrophin
[1bhd]:

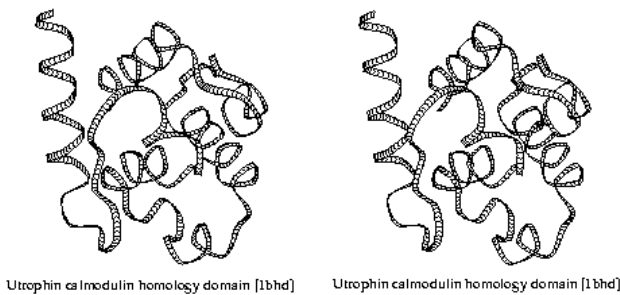
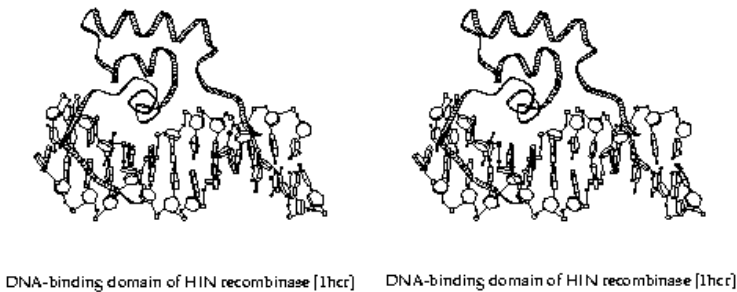
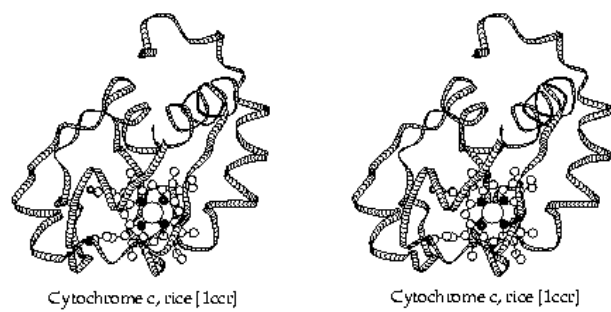


Fig 1-10c: HIN recombinase, DNA-binding domain [1hcr]:



(d) Rice embryo cytochrome c [1ccr]



RCSB PDB: Structure Explorer - Microsoft Internet Explorer

檔案(F) 編輯(E) 檢視(V) 我的最愛(A) 工具(T) 說明(H)

← 上一頁 搜索 我的最愛 移至 連結

網址(U) http://www.rcsb.org/pdb/explore.do?structureId=1CCR

- 1CCR
- Download Files
- FASTA Sequence
- Display Files
- Display Molecule
- Structural Reports
- Structure Analysis
- Help

1CCR

Title STRUCTURE OF RICE FERRICYTOCHROME C AT 2.0 ANGSTROMS RESOLUTION

Authors Ochi, H., Hata, Y., Tanaka, N., Kakudo, M., Sakurai, T., Aihara, S., Morita, Y.

Primary Citation Ochi, H., Hata, Y., Tanaka, N., Kakudo, M., Sakurai, T., Aihara, S., Morita, Y. Structure of rice ferricytochrome c at 2.0 Å resolution. *J.Mol.Biol.* v166 pp.407-418, 1983 [Abstract]

History Deposition 1983-03-14 Release 1983-04-21

Experimental Method Type X-RAY DIFFRACTION Data [EDS]

Resolution(Å)	R-Value	R-Free	Space Group
1.50	0.190 (work)	n/a	P 6 ₁

Unit Cell

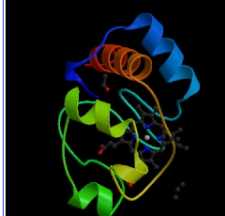
Length (Å)	a	43.78	b	43.78	c	110.05
Angles (°)	alpha	90.00	beta	90.00	gamma	120.00

Molecular Description monomer (protein 112 residues)

Asymmetric Unit Polymer: 1 Molecule: CYTOCHROME C Chains: _

Images and Visualization

Biological Molecule / Asymmetric Unit



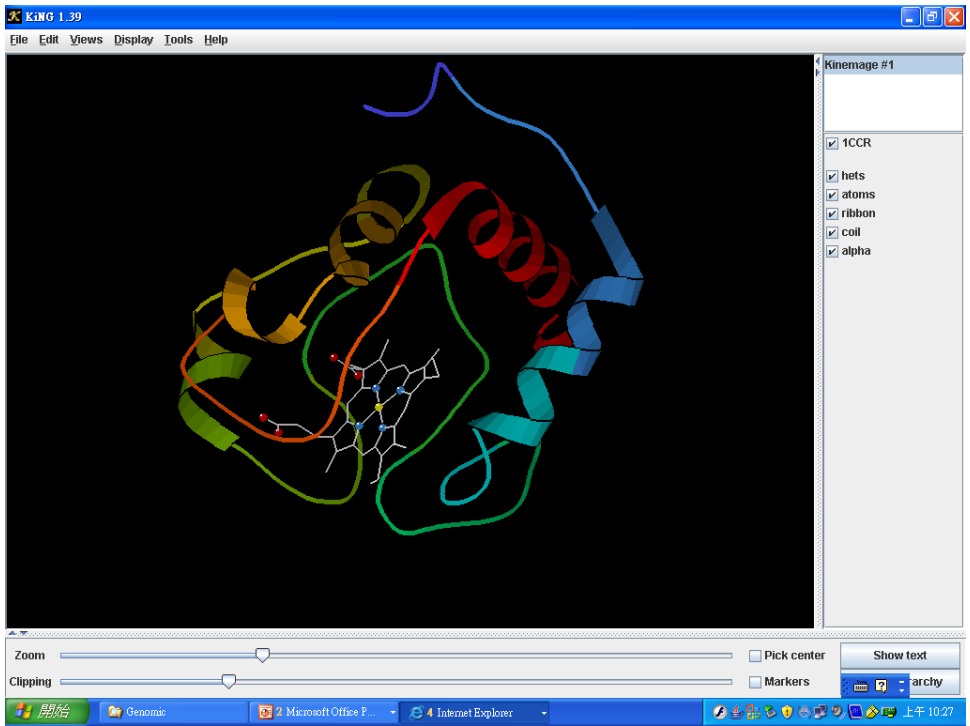
Display Options

- KING
- Jmol
- WebMol
- Protein Workshop
- Quick-PDB
- All Images

完成

Genomic 01_Background GSP_ch1_3 RCSB PDB: Structure

網路網路 上午 09:50



RCSB PDB: Sequence Details Report - Microsoft Internet Explorer

Home Search Structure Queries

Structure Summary Biology & Chemistry Materials & Methods Sequence Details Geometry

Sequence Details 1CCR

Chain _ representative of identical chains Chain _

Description CYTOCHROME C

Type polypeptide(L)

Polymer Id 1

Number of residues 112

Domains [d1ccr_1: Mitochondrial cytochrome c](#)

Sequence and Secondary Structure

Key: = extended strand, = turn, = disulfide bond, = alpha helix, = 310 helix, = pi helix, Greyed out residues have no structural information

XASFSKAPPGNPKAGEKTFKTKCAOCHTVDKGAGHKOGPNLNLPGKSGTTPGYSYSTA

DKNMAVIWEENTLYDYLLNPKKYTPGTRMVFPLKKPQERADLTSLKKEATS

Download Chain _ in Fasta Format

For Sequence Only

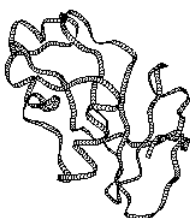
Windows taskbar: 開始, Genomic, RCSB PDB: Sequenc..., GSP_ch1_3, 下午 06:55



TATA-box-binding protein [1cdw]



TATA-box-binding protein [1cdw]



barnase [1bcm]



barnase [1bcm]



OE-domain from Lys-tRNA synthetase [1bbw]



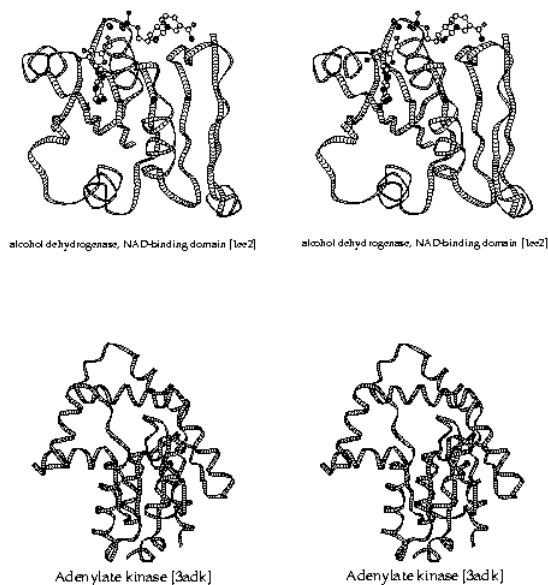
OE-domain from Lys-tRNA synthetase [1bbw]



Scytalone dehydratase [3std]



Scytalone dehydratase [3std]



RCSB PDB : Structure Explorer - Microsoft Internet Explorer

檔案(F) 編輯(E) 檢視(V) 我的最愛(A) 工具(T) 說明(H)

← 上一頁 → 搜索 我的最愛

網址(U) <http://www.rcsb.org/pdb/explore.do?structureId=1EE2> 移至 連結 >>

RCSB PDB

PROTEIN DATA BANK

Contact Us | Help | Print Page

A MEMBER OF THE CPDB

An Information Portal to Biological Macromolecular Structures

As of Tuesday Mar 21, 2006 there are 35701 Structures | PDB Statistics

PDB ID or keyword
Author
SEARCH
Advanced Search

Home Search Structure Results
Structure Summary Biology & Chemistry Materials & Methods Sequence Details Geometry

Queries

1EE2

Download Files

FASTA Sequence

Display Files

Display Molecule

- Image Gallery
- KiNG Viewer
- Jmol Viewer
- WebMol Viewer
- Rasmol Viewer (Plugin required)
- Swiss-PDB Viewer (Plugin required)
- KiNG Help
- Jmol Help
- WebMol Help
- Protein Workshop Help
- QuickPDB
- Asymmetric Unit / Biological Molecule

Title

THE STRUCTURE OF STEROID-ACTIVE ALCOHOL DEHYDROGENASE AT 1.54 Å RESOLUTION

Authors

Adolph, H.W.

Primary Citation

Adolph, H.W., Zwart, P., Meijers, R., Hubatsch, I., Kiefer, M., Lanzin, V., Cedergren-Zeppezauer, E. Structural basis for substrate specificity differences of horse liver alcohol dehydrogenase isozymes. *Biochemistry* v39 pp.12885-12897, 2000

[Abstract]

History

Deposition 2000-01-30 Release 2000-10-27

Experimental Method

Type X-RAY DIFFRACTION Data [EDS]

Parameters

Resolution[Å]	R-Value	R-Free	Space Group
1.54	0.148 (obs.)	0.183	P 2 ₁ (P 1 2 ₁ 1)

Unit Cell

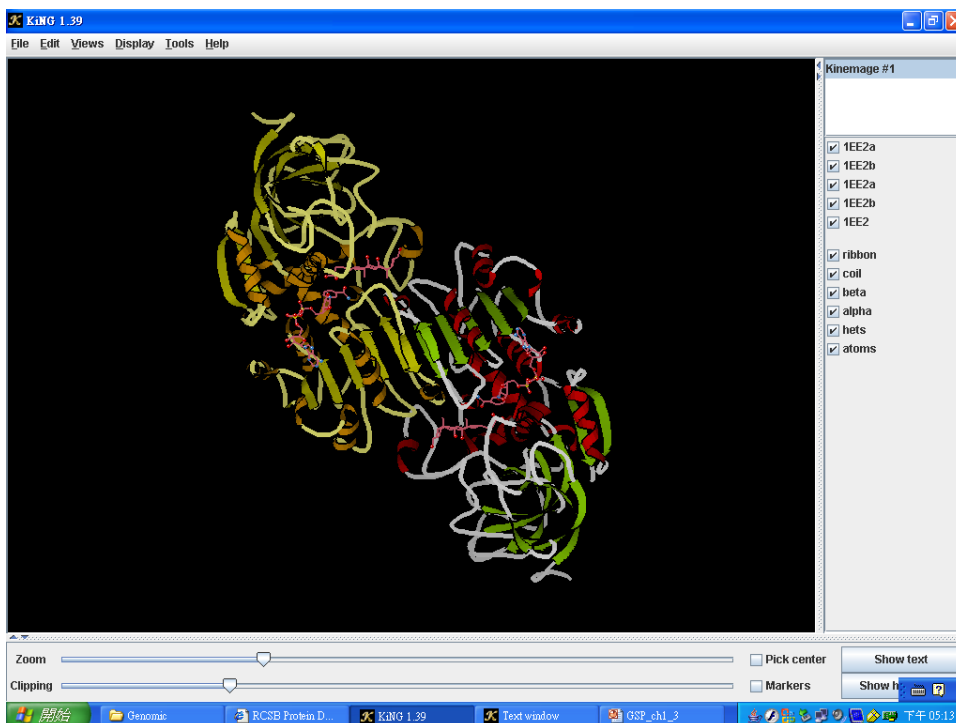
Length [Å]	a	b	c
Å	55.03	73.16	92.49
°	90.00	90.00	90.00

Images and Visualization

Biological Molecule / Asymmetric Unit

Display Options

- KiNG
- Jmol
- WebMol
- Protein Workshop
- QuickPDB
- All Images



RCSB PDB: MarvinView - Microsoft Internet Explorer

檔案(F) 編輯(E) 檢視(V) 我的最愛(A) 工具(T) 說明(H)

← 上一頁 → 搜尋 我的最愛

網址(1) http://www.rcsb.org/pdb/marvin.do?handler=structureExplorer&hetId=NAD&sid=1EE2 移至 連結 >>

RCSB PDB
PROTEIN DATA BANK

A MEMBER OF THE PDB

An Information Portal to Biological Macromolecular Structures

As of Tuesday Mar 21, 2006 there are 35701 Structures | PDB Statistics

Contact Us | Help | Print Page

PDB ID or keyword Author SEARCH Advanced Search

Home Search Structure Results

Queries

- Home
- Tutorial About This Site
- Getting Started
- Download Files
- Deposit and Validate
- Structural Genomics
- Dictionaries & File Formats
- Software Tools
- Educational Resources
- General Information
- Acknowledgements
- Frequently Asked Questions
- Known Problems
- Report Bugs/Comments

Ligand Summary Powered by ChemAxon 1EE2

Right click on the image for animation, color and other options.

Name NICOTINAMIDE-ADENINE-DINUCLEOTIDE

HET ID NAD

Formula $C_{21}H_{27}N_7O_{14}P_2$

SMILES String NC(=O)c1ccc[n+](c1)C2OC(COP([O-])=O)OP([O-])=OCC3OC(C(O)C3O)n4cnc5c(N)ncnc45C(O)C2O

輔一菸鹼胺腺呤雙核酸(Nicotinamide adenine dinucleotide, NAD)

Back to Structure Explorer

© RCSB Protein Data Bank

Applet IDView started

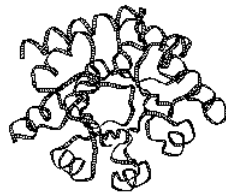
開始 Genomic 01_Background GSP_chl_3 3 Internet Explorer 上午 10:39



Chemotaxis receptor methyltransferase [1af7]



Chemotaxis receptor methyltransferase [1af7]



Thiamine phosphate synthase [2tps]



Thiamine phosphate synthase [2tps]



pancreatic spasmodic polypeptide [2pep]



pancreatic spasmodic polypeptide [2pep]

Web resources

- The Worldwide Protein Data Bank (wwPDB)
<http://www.wwpdb.org/>
- The Research Collaboratory for Structural Bioinformatics (RCSB) (USA) <http://www.rcsb.org/>
- The Macromolecular Structure Database (MSD) (UK)
<http://www.ebi.ac.uk/pdbe/>
- The protein databank Japan <http://www.pdbj.org/>
- BMRB (USA) <http://www.bmrwisc.edu/>
- Structural Classification of Proteins (SCOP)
<http://scop.mrc-lmb.cam.ac.uk/scop/>
- The Molecular Modeling DataBase (MMDB)
<http://www.ncbi.nlm.nih.gov/Structure/MMDB/mmdb.shtml>

Protein structure prediction and engineering

- Amino acid sequence of a protein dictates its 3D structure
- If amino acid sequences contain sufficient information to specify 3D structures of proteins, it should be possible to ***devise an algorithm to predict protein structure from amino acid sequence.***
 - This has proved elusive (難以理解的).

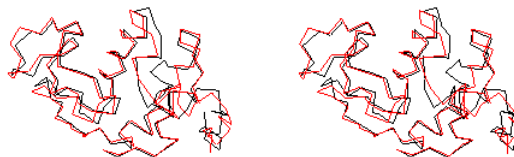
Less-ambitious goals:

- **Secondary structure prediction** — which segments of the sequence form helices and which form strands of sheet?
- **Fold recognition** — Given a library of known protein structures and their amino acids sequences, and the amino acid sequence of a protein of unknown structure, can we find the structure in the library that is most likely to have a folding pattern similar to that of the protein of unknown structure?
- **Homology modelling** — If the sequences of two homologous proteins have 50% or more identical residues in an optimal alignment, the structures are likely to have similar conformations over more than 90% of the model.

Aligned sequences and superposed structures of two related proteins: Alignment of Chicken lysozyme and Baboon alpha-lactalbumin

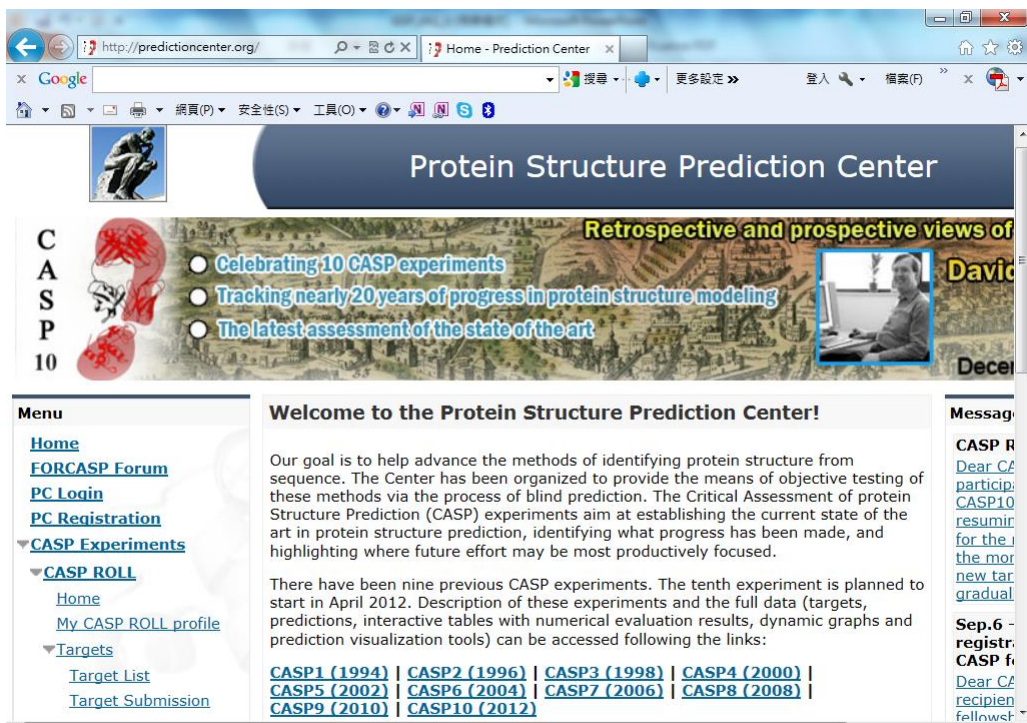
Chicken lysozyme	KVFGRCELAAAMKRHGLDNYRGYSLGNWVCAAKFESNFNTQATNRNTDGS
Baboon alpha-lactalbumin	KQFTKCELSQNLV-DIDGYGRIALPELICTMFHTSGYDTQAIVEND-ES
Chicken lysozyme	TDYGILQINSRWWCNDGRTPGSRNLCNIPCSALLSSDITASVNC AKKIVS
Baboon alpha-lactalbumin	TEYGLFQISNALWCKSSQSPQSRNICDITCDKFLDDDDITDDIMCAKKILD
Chicken lysozyme	DGN-GMNAWVAWRNRCKGTDVQA-WIRGCRL-
Baboon alpha-lactalbumin	I-KGIDYWIAHKALC-TEKL-EOWL-CE-K

Superposition of Chicken lysozyme (black) and Baboon alpha-lactalbumin (red):



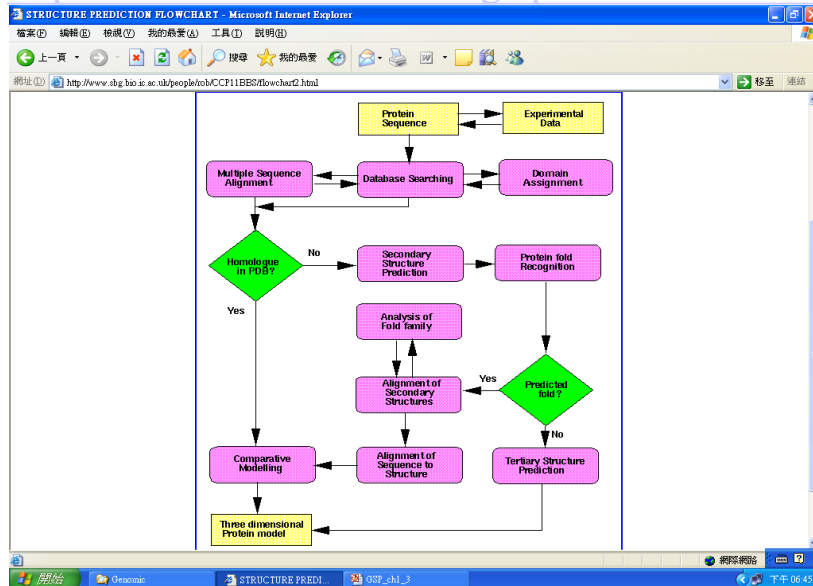
Critical Assessment of Structure Prediction (CASP)

- Judging of techniques for predicting protein structures requires blind test.
- Predictors submit models, which are held until the deadline for release of the experimental structure.
- Then the predictions and experiments are compared – to the delight of a few and the chagrin of most.



STRUCTURE PREDICTION FLOWCHART

<http://www.russell.embl.de/gtsp/flowchart2.html>



A computer game to learn protein folding

- Maintained by University of Washington, Department of Computer Science
- <http://fold.it/>
- Learn to play this game and get a score as high as you can
 - Download the “get started”
 - Register an account



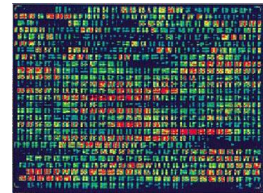
Protein Engineering

- In the laboratory we can manipulate nucleic acids and protein at will.
 - We can probe them by exhaustive mutation to see the effects on function.
 - We can endow (賦予) old proteins with new functions, as in the development of catalytic (催化作用) antibodies
 - We can even create new ones. Engineered proteins must obey the laws of physical chemistry but not the constraints of evolution. With engineered proteins we can explore new territory.

Proteomics

- Combines the census (統計數), distribution, interactions, dynamics, and expression patterns of the proteins within living systems.
- A data-intensive subject, depending on high-throughput measurements
 - Include DNA microarrays, and mass spectrometry.

DNA Microarrays



- Or DNA chips
- Devices for checking a sample simultaneously for the presence of many sequences
- Can be used
 - To determine expression patterns of different proteins by detection of mRNAs
 - For genotyping(遺傳型), by detection of different variant gene sequences, including but not limited to single-nucleotide polymorphisms (SNPs)

Applications of DNA microarrays

- Identifying genetic individuality in tissues or organisms, or genotyping
- Investigating cellular states and processes
- Diagnosis of genetic disease
- Diagnosis of infectious disease
- Specialized diagnosis of disease
- Genetic warning signs
- Drug selection
- Target selection for drug design
- Pathogen (病原體) resistance
- Measuring temporal variations in protein expression

System biology

- Integration – to put all cell part back together
- First aspect:
 - The study of patterns within a cell or an organism: pathways and control cascades, and patterns of protein expression.
 - Patterns have both static and dynamic aspects
 - Identification of pairs of proteins that bind to each other, and assembly of pairwise interactions into a network Static pattern.
 - Dynamic pattern: the flow of metabolites through a network of enzymes, or the flow of information down a control cascade, is a dynamic pattern.

- Second aspect – comparison of occurrence, activities and interactions of genes and proteins *across different species*.
 - The systems we are trying to understand arose through processes of evolution. Different species illuminate one another.
- High-throughput methods of genomics and proteomics provide data about sequences, expression patterns and interactions.
 - Systems biology takes the data as pieces of a jigsaw puzzle that extends in both space and time. To understand the complex and delicate instrument that is the living cell, we must fit the pieces into their frame.

Clinical implications

1. Diagnosis of disease and disease risks
 - DNA sequencing can detect the absence of a particular gene, or a mutation.
2. Genetics of responses to therapy – customized treatment
 - People differ in their ability to metabolize drugs, different patients with the same condition may require different dosages.

3. Identification of drug targets

A target is a protein the function of which can be selectively modified by interaction by a drug, to affect the symptoms or underlying causes of a disease.

4. Gene therapy

If a gene is missing or defective, we'd like to replace it or at least supply its product. If a gene is overactive, we'd like to turn it off.

Direct supply of proteins is possible for many diseases.

Practice

- Huntington disease
- Find out the cause of this disease using the Internet search.
- What is the phenomenon of “anticipation”?
- Answer:
- The same questions for other diseases: 地中海型貧血 Mediterranean anemia , 紅斑性狼瘡 Systemic Lupus Erythematosus