

Predikce struktury proteinů

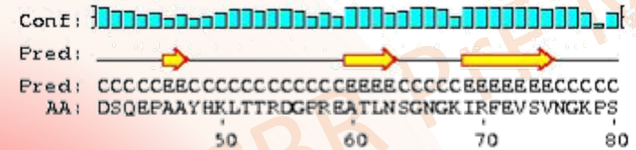
Pokročila bioinformatika NCEB PrF MU

Struktura proteinů

1D

ADSQTSSNRAGEFSIPPNTDFRAIFFANAAE
QQHIKLFIGDSQEPAAAYHKLTTTRDGPREATL
NSGNGKIRFEVSVNGKPSATDARLAPINGK
KSDGSPFTVNFVIVVSEDGHDSYNDGIVV
LQWPIG

**primární
(sekvence)**

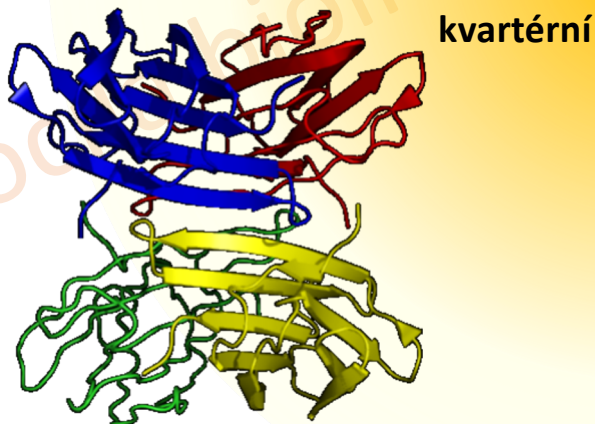


2D

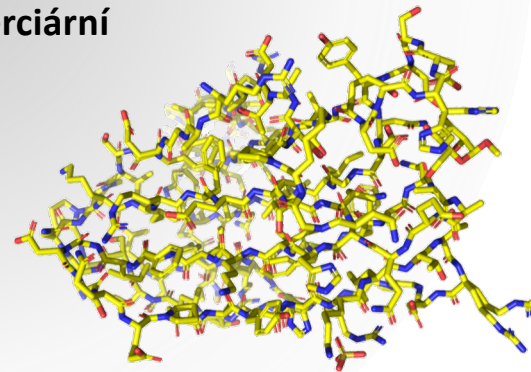


sekundární

4D



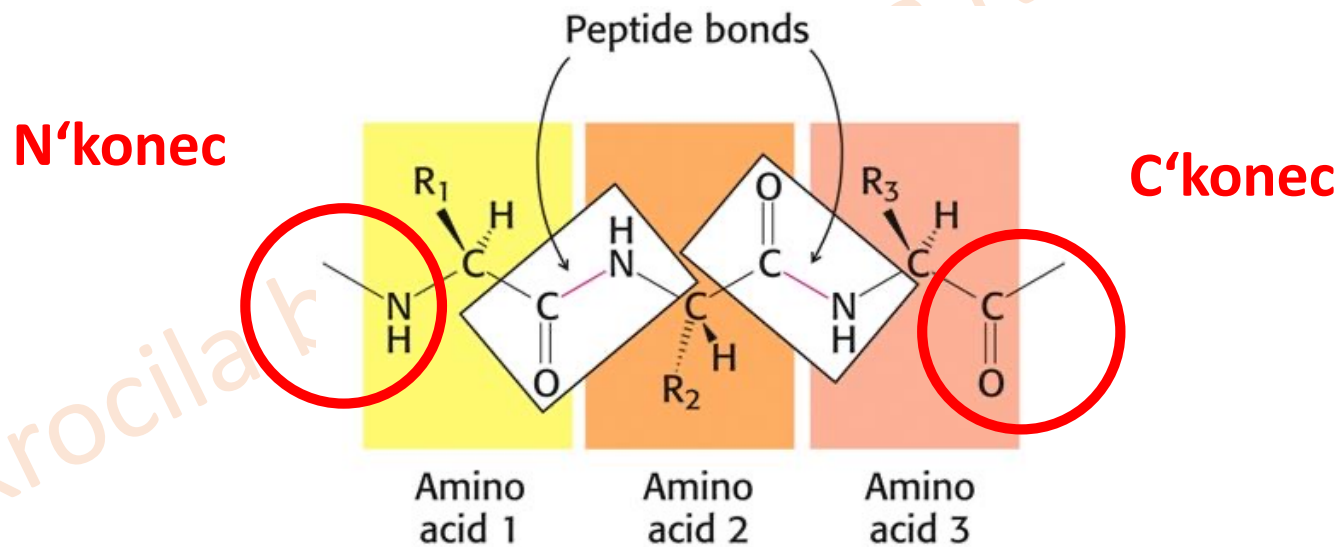
terciární



3D

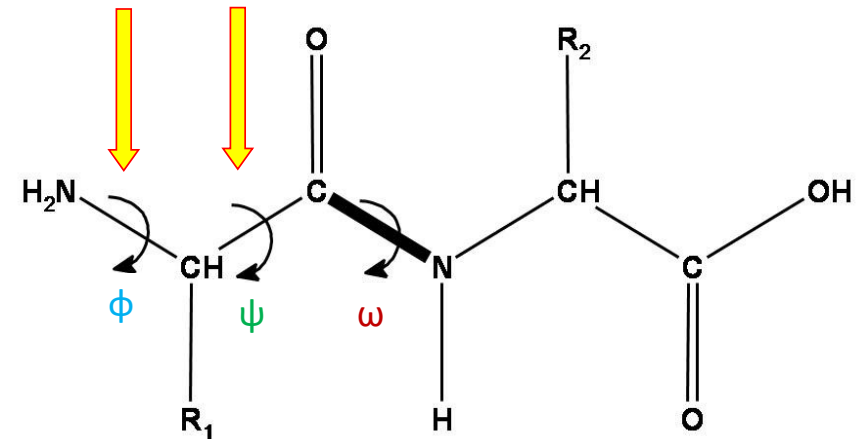
Primární struktura

- Sekvence aminokyselin zapsaná od N' konce k C' konci



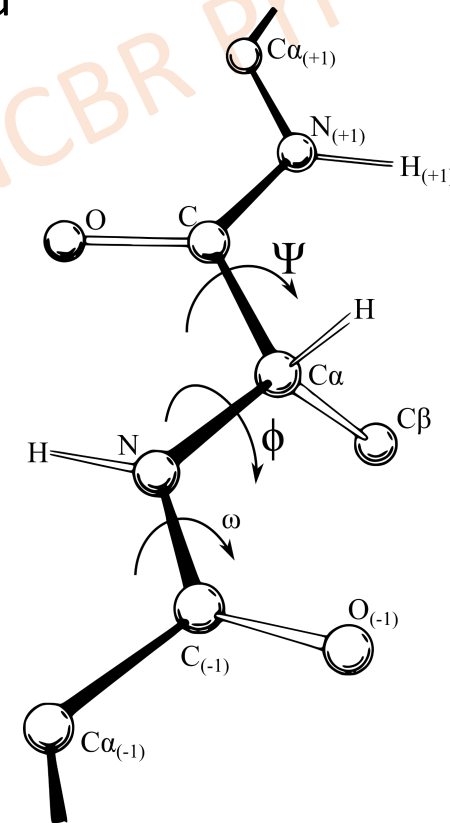
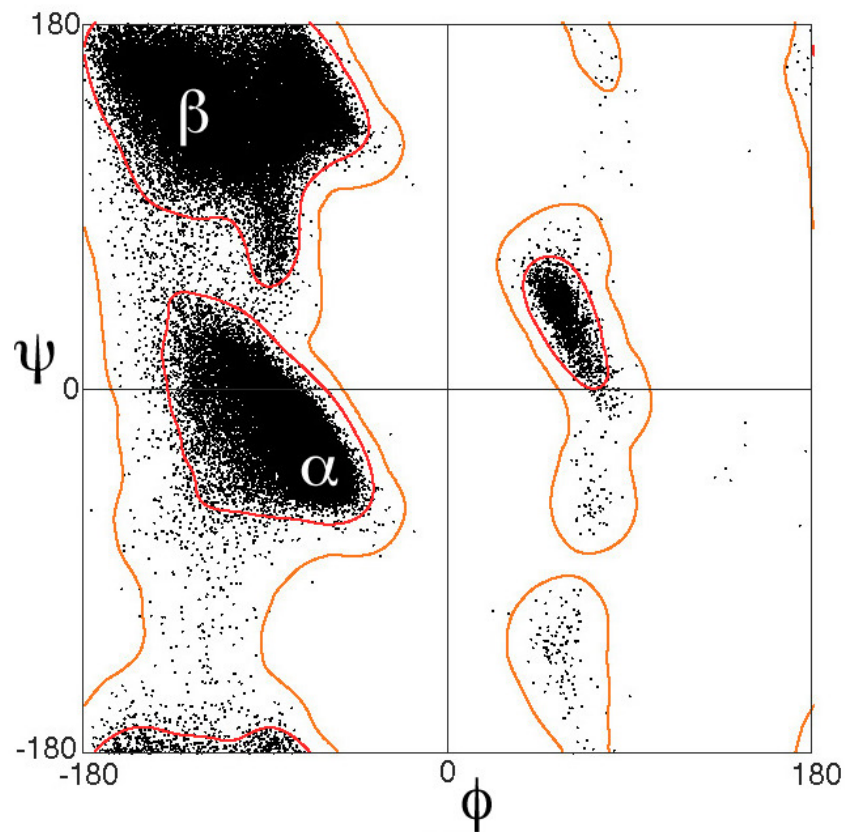
Sekundární struktura

- Definována pomocí **torzních úhlů** peptidové páteře
- Pro každou aminokyselinu lze definovat tři úhly:
 - ϕ – úhel kolem vazby N-C α
 - ψ – úhel kolem vazby C α -C(karb.)
 - ω – úhel kolem peptidové vazby (180°, výjimečně 0°)
- Stabilizována pomocí vodíkových můstků mezi atomy peptidové kostry



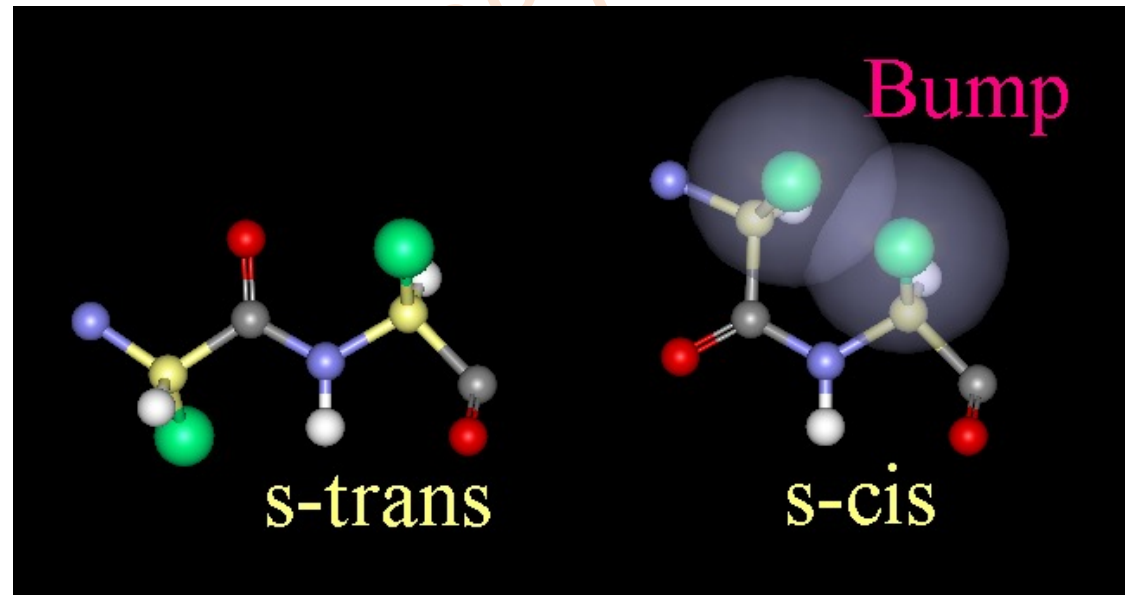
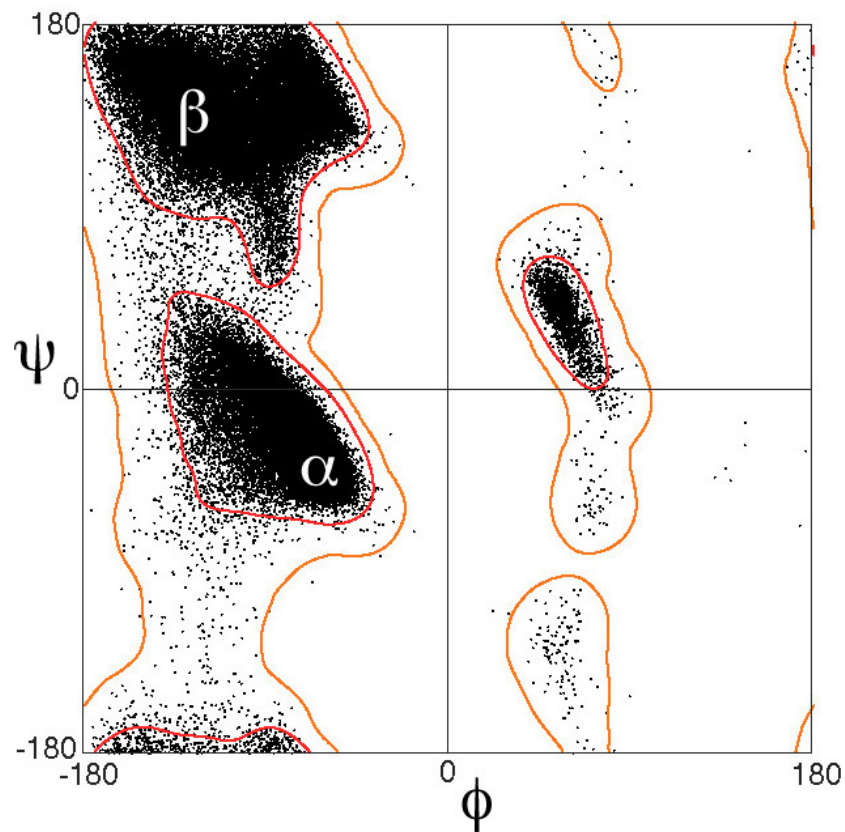
Ramachandranův diagram

➤ Každé aminokyselině odpovídá jeden bod v diagramu



Ramachandranův diagram

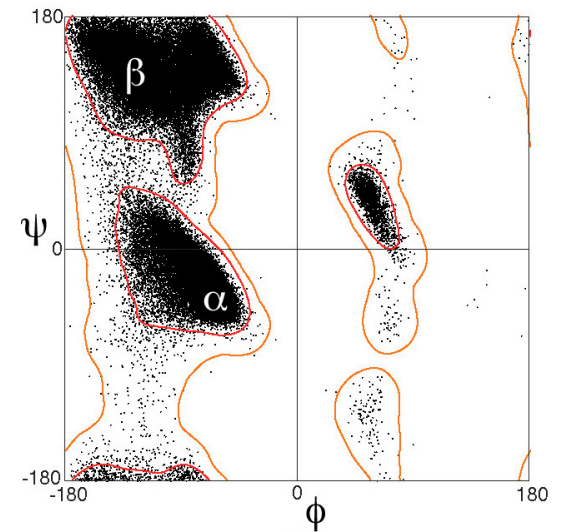
➤ Každé aminokyselině odpovídá jeden bod v diagramu



Sekundární struktura

2D

- Stabilní konformace polypeptidového řetězce
- Důležité pro udržení 3D struktury
- α -šroubovice (helix), β -skládaný list (sheet), otáčky, smyčky
- Cca 50 % aminokyselin je součástí α a β struktur

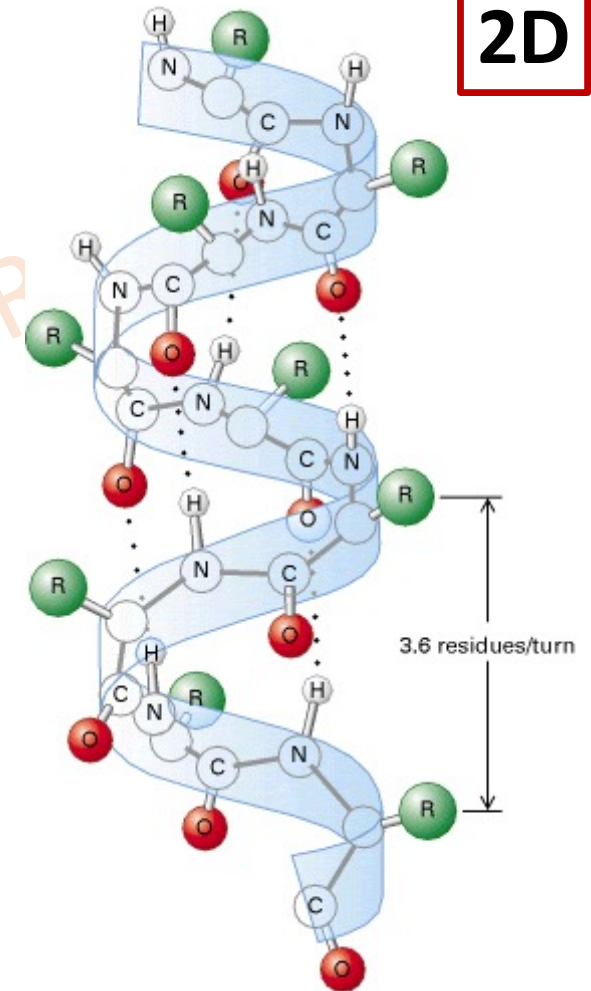


Šroubovice (helix)

2D

- α -helix – nejčastější
- 3_{10} -helix – obvykle na začátku nebo na konci α -helixu
- π -helix – málo stabilní, málo častý

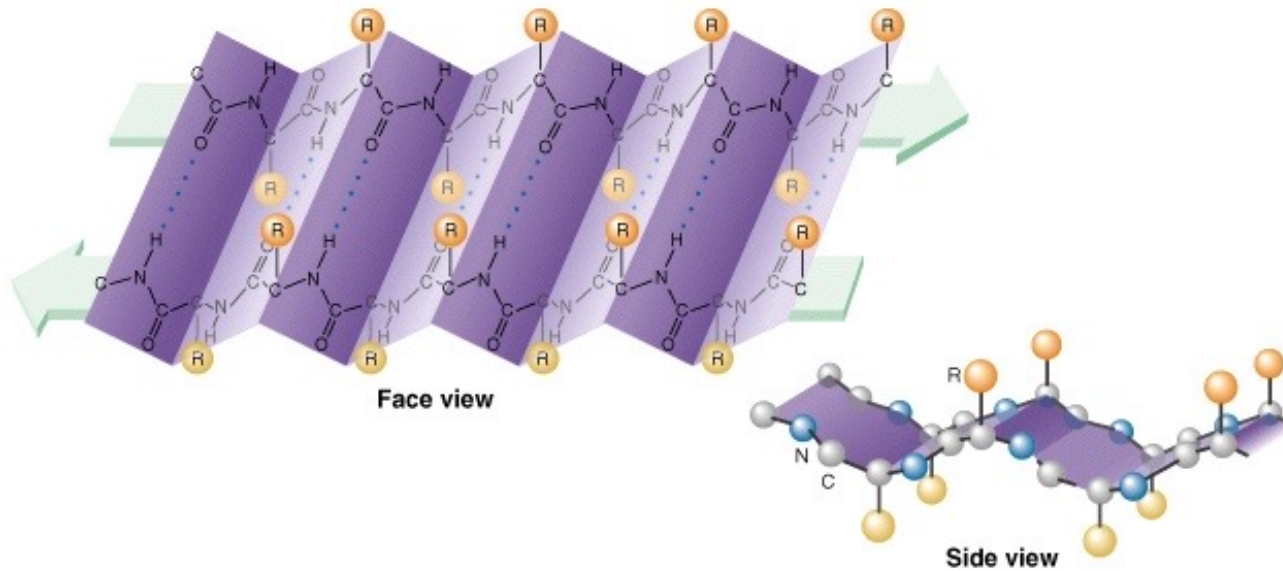
	α -helix	3_{10} -helix	π -helix
Vodíkové můstky	n ... n+4	n ... n+3	n ... n+5
Residua na otáčku	3,6	3	4,4
Vinutí (Å na 1 AK)	1,5	2	1,15



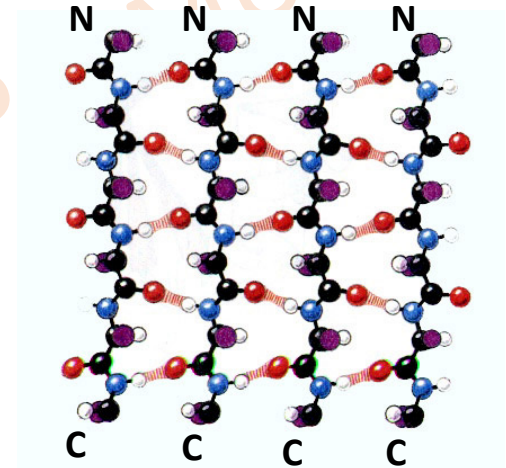
Skládaný list (extended β -sheet)

2D

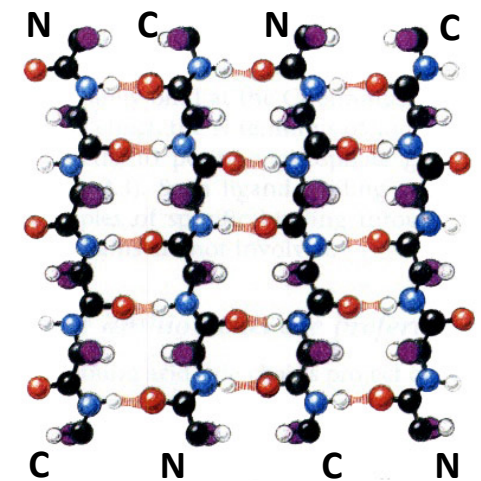
➤ Paralelní, antiparalelní, mix



Paralelní



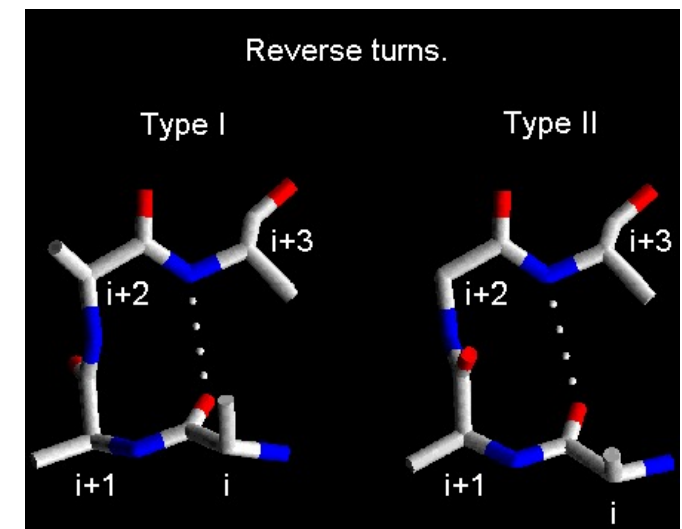
Antiparalelní



Ostatní

2D

- Úseky které nespádají do kategorií helix nebo list
- **Kombinace** povolených torzních úhlů
- **Nestabilní** konformace
- **Nestandardní** konformace (glycin, prolin)
- **Otáčky** (turns), „náhodné klubko“ (random coil)

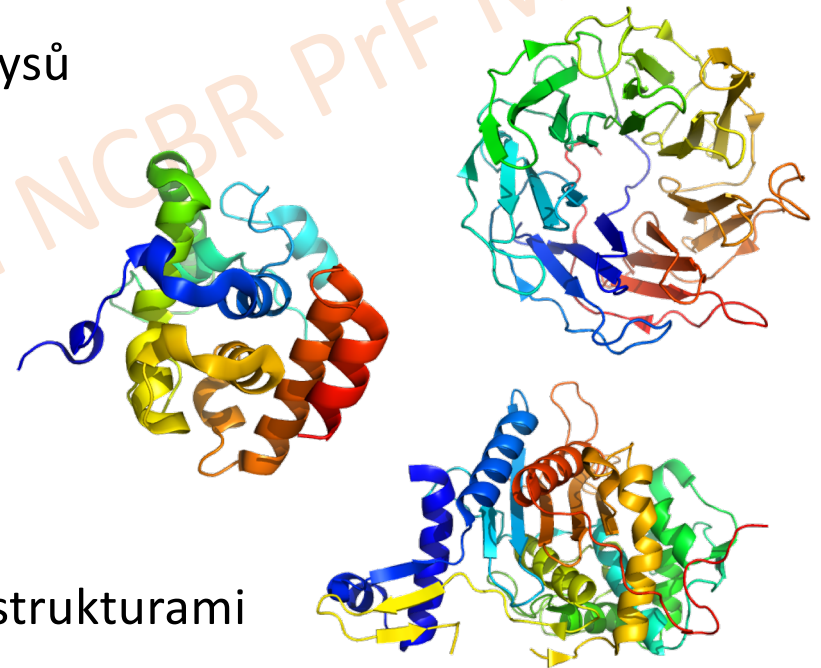


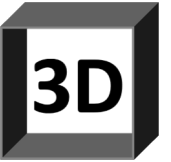
Dělení proteinů dle 2D struktury

Zejména pro účely klasifikace, hledání společných rysů

Každý protein obsahuje mj. smyčky a ohyby

- Jen α struktury
- Jen β struktury
- α/β
 - Motivy kombinující α i β struktury
- $\alpha + \beta$
 - Oddělené domény tvořené jen α nebo jen β strukturami
- **Malé proteiny**
 - Speciální případy, např. obsahující ionty kovů, stabilizované disulfidickými můstky



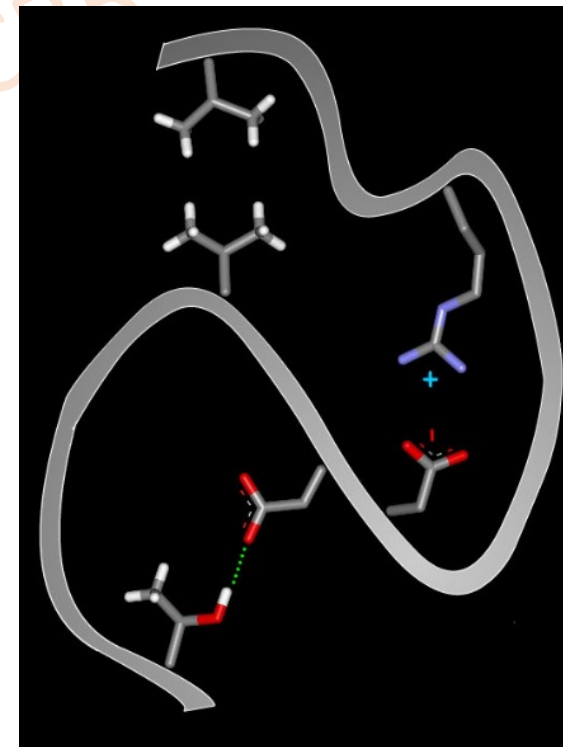


Terciární struktura

➤ Konkrétní umístění jednotlivých atomů polypeptidového řetězce v prostoru

➤ Stabilizována pomocí:

- **Vodíková vazba** (H-můstek)
mezi polárními AK, mezi hlavním řetězcem
- **Iontová** interakce – nabité AK
- **Hydrofobní** interakce – nepolární AK
- „**Stacking**“ (π - π , CH- π interakce) – aromatické AK
- Kovalentní vazba **síra-síra** – cystein / cystin
- Vazba **iontů kovů**



Od 2D ke 3D



➤ **Motivy**

- 2-3 prvky sekundární struktury

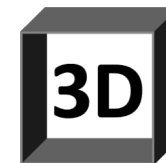
➤ **Foldy**

- Kombinace jednoduchých motivů

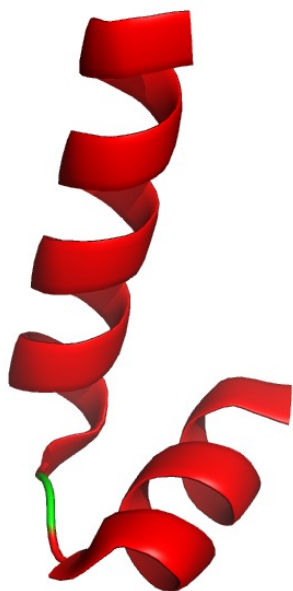
➤ **Domény**

- Tvořeny motivy/foldy
- Část struktury s vlastní funkcí (nejmenší funkční jednotka)
- Nezávislá jednotka (alespoň částečně nezávislá)

Jednoduché motivy



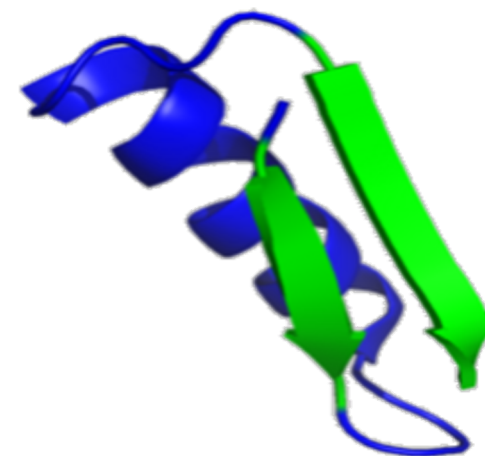
Helix-turn-helix



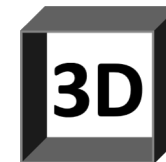
β -vlásenka



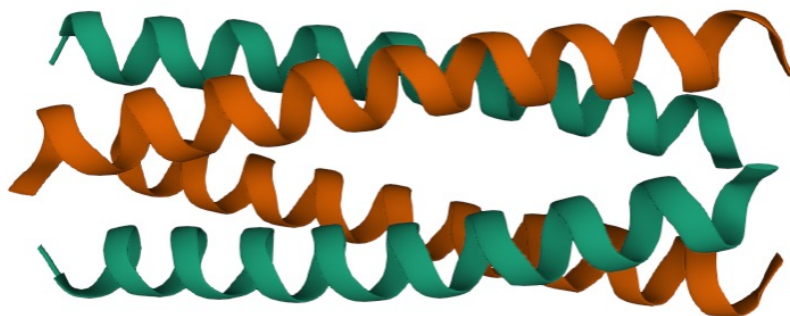
β - α - β



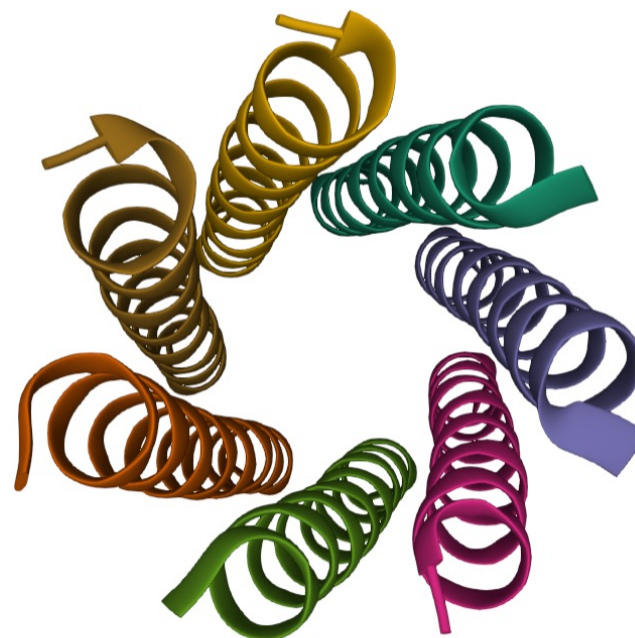
Složené α -motivy/foldy



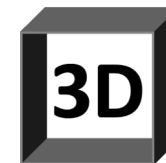
4-helix bundle



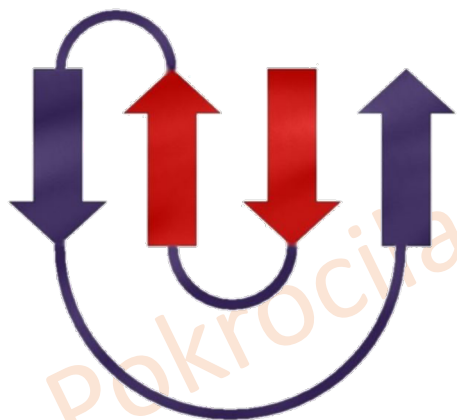
7-helix barrel



Složené β -motivy/foldy



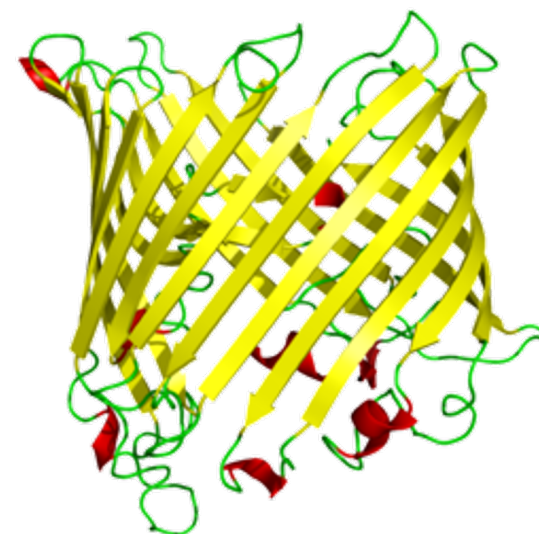
Řecký klíč



β -meandr



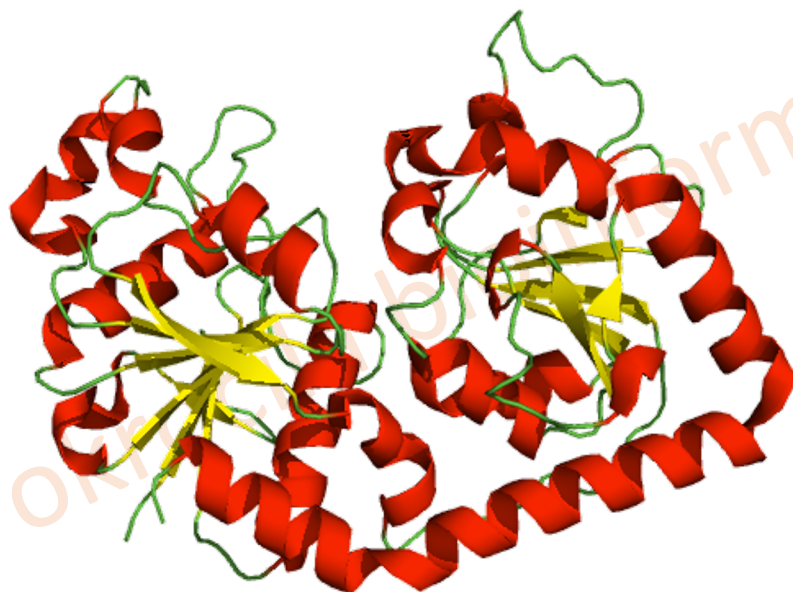
β -barel



Složené α/β -motivy/foldy



Rossmannův fold



TIM-barel



Structural classification of proteins (SCOP)

2D

<https://scop.mrc-lmb.cam.ac.uk/>

The screenshot shows the SCOP 2 website interface. At the top, there is a navigation bar with the SCOP logo, links for 'About', 'Contact', and 'Download', and a search input field containing 'dnmt1'. Below this is a banner for legacy SCOP websites. The main content area features the 'SCOP 2' title, a 'Learn More' button, and a descriptive paragraph about the database. A search bar with 'Keyword and ID search' and 'Sequence search' options is present, with a 'Go' button. At the bottom, there are two columns of links for browsing by structural class and protein type.

SCOP 2

SCOP: Structural Classification of Proteins

Nearly all proteins have structural similarities with other proteins and, in some of these cases, share a common evolutionary origin. The SCOP database, created by manual inspection and abetted by a battery of automated methods, aims to provide a detailed and comprehensive description of the structural and evolutionary relationships between all proteins whose structure is known. As such, it provides a broad survey of all known protein folds, detailed information about the close relatives of any particular protein, and a framework for future research and classification.

Latest update on **2020-03-31** includes **44,218** non-redundant domains representing **532,428** protein structures. Folds, superfamilies and families statistics [here](#).

Keyword and ID search Sequence search

Enter free text, SCOP ID, PDB ID or UniProt ID

Go

Browse by structural class

- All alpha proteins
- All beta proteins
- Alpha and beta proteins(a/b)
- Alpha and beta proteins(a+b)
- Small proteins

Browse by protein type

- Globular proteins
- Membrane proteins
- Fibrous proteins
- Non-globular/Intrinsically unstructured proteins

CATH – Protein structure classification database

➤ Domény jsou klasifikovány podle CATH hierarchie

➤ Třída (Class)

- Podle sekundární struktury
- Jen α , jen β , α i β , minimum sekundární struktury

➤ Architektura

- 3D uspořádání sekundární struktury

➤ Topologie/fold

- Jak jsou prvky sekundární struktury uspořádané za sebou

➤ Homologní nadrodina

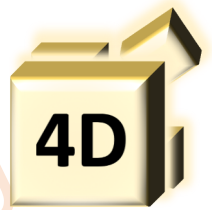
- V případě, že jsou domény evolučně příbuzné (homologní proteiny)

<https://www.cathdb.info/>

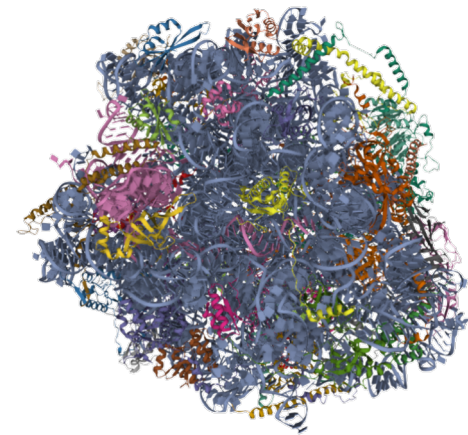
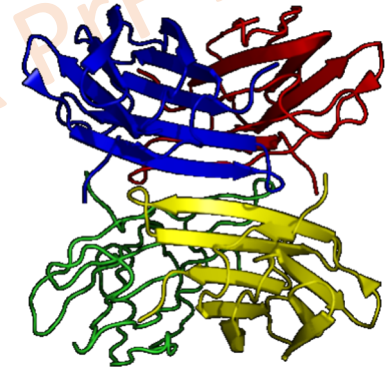
The screenshot shows the CATH database interface with three sections:

- Matching CATH Superfamilies:** Displays a protein structure visualization, the CATH ID **2.120.10.70**, and the name **Fucose-specific lectin**. It includes a list of PDB codes and a "View all entries" button.
- Matching CATH Domains:** Displays a protein structure visualization, the CATH ID **4agiA00**, and the name **PDB code 4agi, chain A, domain 00**. It includes a list of PDB codes and a "View all entries" button.
- Matching PDB Structures:** Displays a protein structure visualization, the CATH ID **4agi**, and the name **PDB code 4agi**. It includes a list of PDB codes and a "View all entries" button.

Kvartérní struktura



- Vzájemná kombinace více řetězců (monomerů)
- Podle typu podjednotek:
 - **Homooligomery** (identické jednotky)
 - **Heterooligomery** (alespoň dva různé typy jednotek)
- **Komplexy** proteinů s dalšími makromolekulami
 - Ribosom, proteasom, replikační komplex,...
- **Nadmolekulární komplexy**
 - Virové částice, buněčná membrána, organely,...



Způsob uložení 3D (4D) strukturních dat



- Veřejně dostupné **databáze**
 - **Protein Data Bank** (PDB), Biological Magnetic Resonance Data Bank, EMDataBank
- **Koordináty** atomů, doplňkové informace (**meta data**)
- Definovaný **formát**
 - PDB
 - mmCIF

Predikce struktury

- Predikce struktury znamená přiřazení strukturních atributů jednotlivým aminokyselinám (2D struktura, koordináty – tvorba 3D modelu)
- Struktura 2D a 3D je konzervovaná více než samotná sekvence
- **Vstupní informace:**
 - Sekvence
 - Fyzikálně-chemické parametry
 - Informace v databázích
- **Výstup:**
 - Model struktury (2D, 3D, 4D)

Proč predikovat strukturu?

- **Klasifikace** proteinů
- Vytvoření modelu struktury pro další studium
- **Předpověď funkce** proteinu
 - Homologní struktury
 - Vazebná místa
- **Analýza povrchu**
 - Přístupnost solventu, tunely, kavity

Predikce sekundární struktury

- Predikce 3 základních typů: H (helix), E (β -list), C/- (smyčka/vše ostatní)
- 1. GENERACE
 - *ab-initio*
 - Vychází z fyzikálně-chemických vlastností a ze statistik pro jednotlivé aminokyseliny

Pokročila bioinformatika NCBR PrF MU

1. Generace – *ab inicio*

Relative Amino acid Propensity Values for Secondary Structure Elements Used in the Chou-Fasman Methods

Amino Acid	(α -Helix)	P (β -Strand)	P (Turn)
Alanine	1.42	0.83	0.66
Arginine	0.98	0.93	0.95
Asparagine	0.67	0.89	1.56
Aspartic acid	1.01	0.54	1.46
Cysteine	0.70	1.19	1.19
Glutamic acid	1.51	0.37	0.74
Glutamine	1.11	1.11	0.98
Glycine	0.57	0.75	1.56
Histidine	1.00	0.87	0.95
Isoleucine	1.08	1.60	0.47
Leucine	1.21	1.30	0.59
Lysine	1.14	0.74	1.01
Methionine	1.45	1.05	0.60
Phenylalanine	1.13	1.38	0.60
Proline	0.57	0.55	1.52
Serine	0.77	0.75	1.43
Threonine	0.83	1.19	0.96
Tryptophan	0.83	1.19	0.96
Tyrosine	0.69	1.47	1.14
Valine	1.06	1.70	0.50

$$R_i(SS)$$

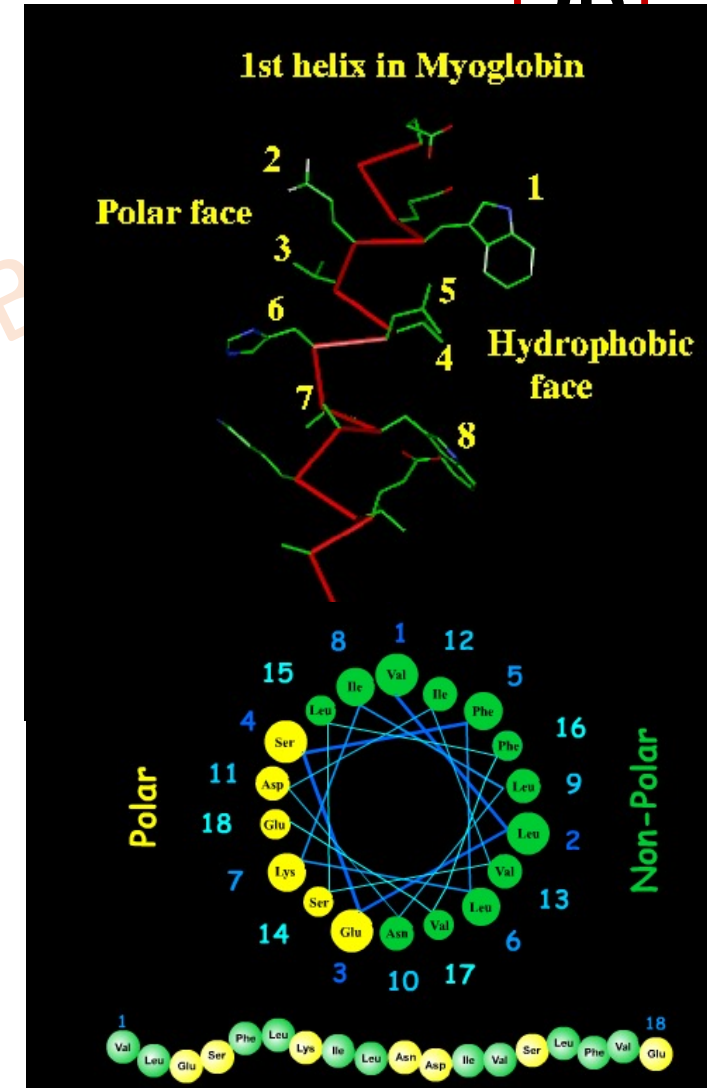
$$R_t(SS)$$

$$\Sigma R_i$$

$$\Sigma R_t$$

Typické znaky α -šroubovice

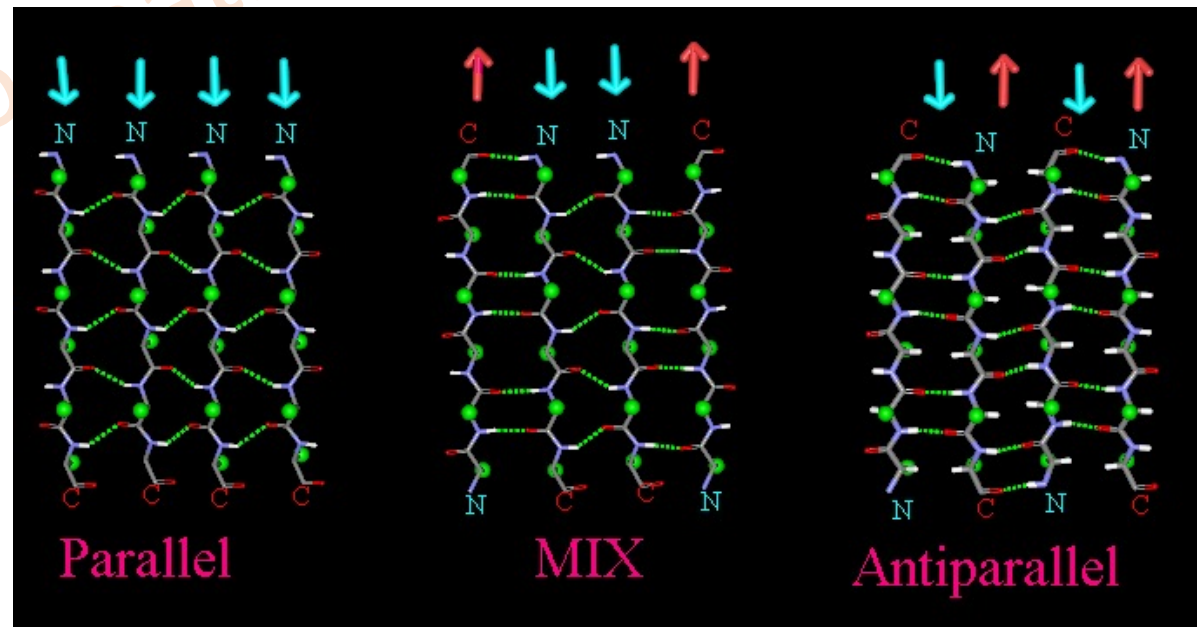
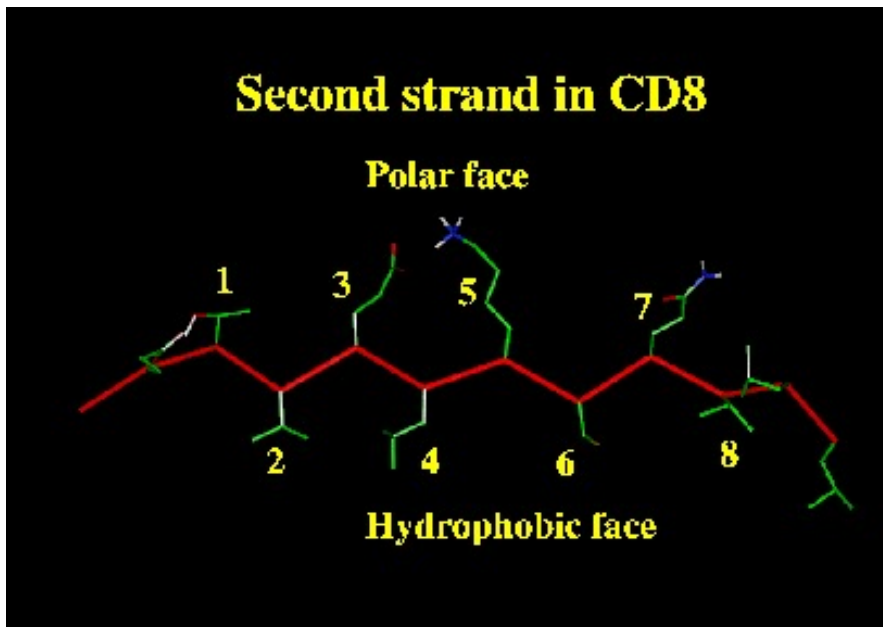
- Často je **částečně exponovaná**
 - Jedna strana je otočená dovnitř proteinu (hydrofobní) a druhá ven (hydrofilní)
 - Residuum (aminokyselina) n , $n+3$, $n+4$, $n+7$ míří na stejnou stranu
- **Transmembránový helix**
 - Všechny aminokyseliny hydrofobní



Typické znaky β -list (musí být stabilizován jinou částí polypeptidového řetězce!)

U β -listu se střídají boční řetězce po 180°

pro částečně zanořený β -list platí, že každé liché reziduum je polární, každé sudé nepolární, u plně zanořeného jsou všechna nepolární... tj. residua směřující na stejnou stranu by měla mít stejný charakter



α -šroubovice nebo β -list?

ELKAHIRVDLTQ

α

ELKAHIRVDLTQ

ELKAHIRVDLTQ

β

Polární

Nepolární

α -šroubovice nebo β -list?

ELKAHIRVDLTQ

ELKAHIRVDLTQ

α



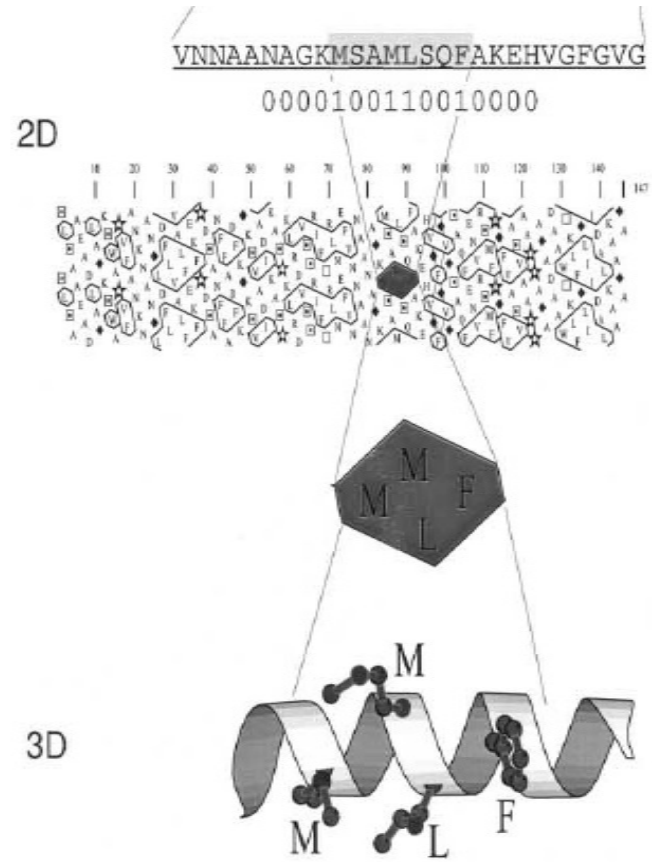
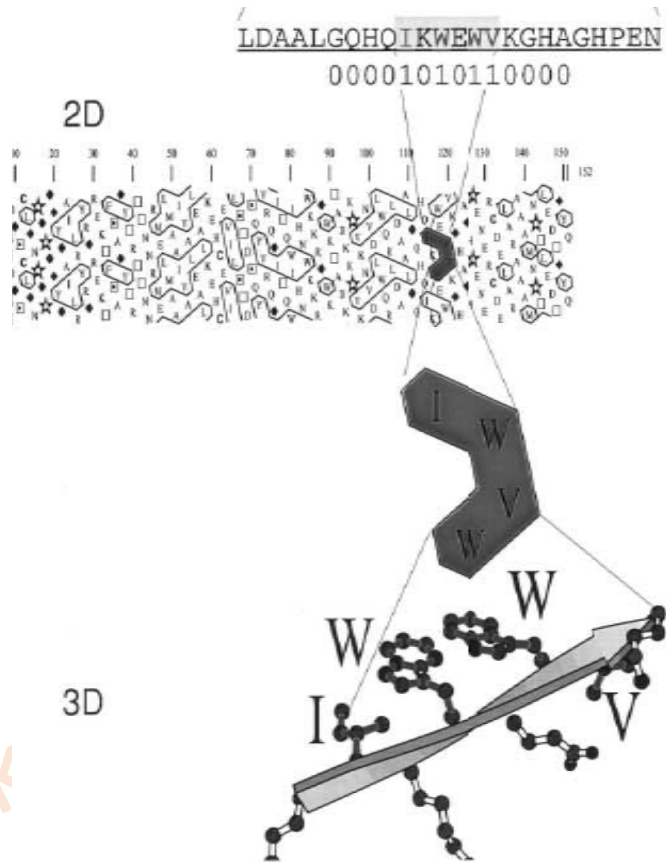
ELKAHIRVDLTQ

β



Polární

Nepolární



Pok

MU

2D

RPBS Web Portal – HCA

<https://mobyale.rpbs.univ-paris-diderot.fr/cgi-bin/portal.py?form=HCA#forms::HCA>

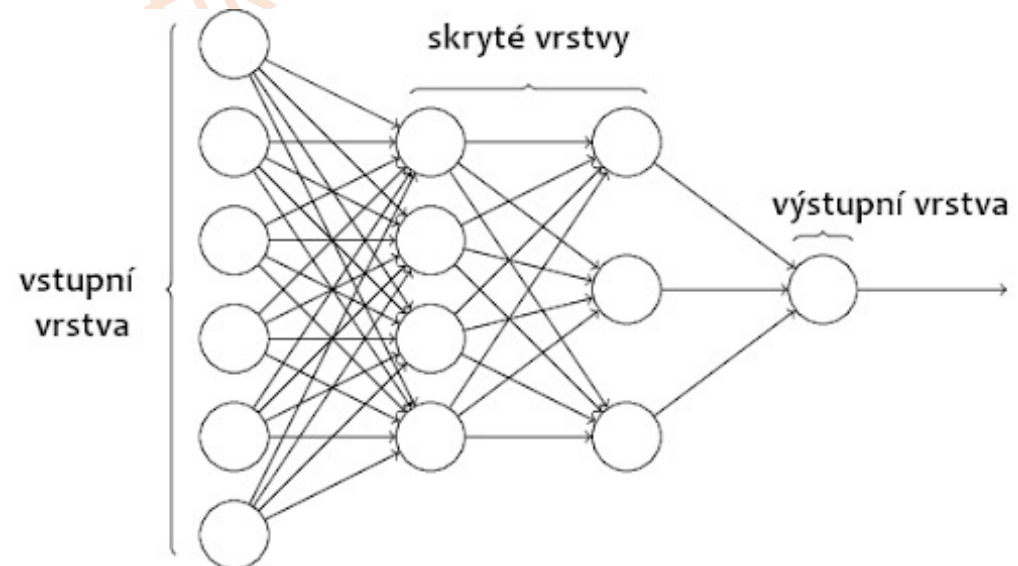
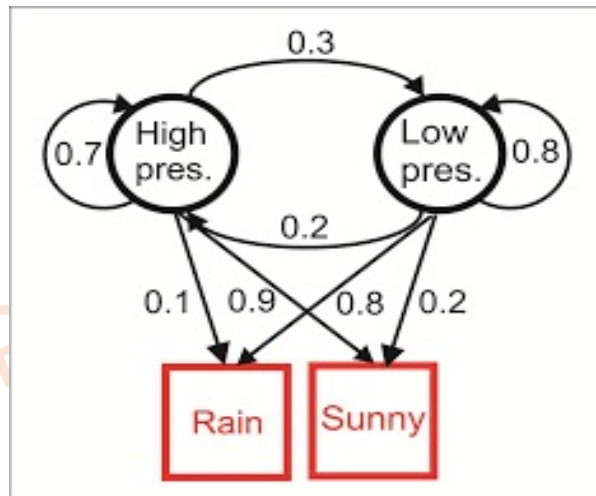
The screenshot displays the RPBS Web Portal interface. At the top, the header includes the RPBS logo, the text "RPBS Web Portal", and user options: "(guest)", "set email", "sign-in", "sign-out", and "refresh workspace". Below the header is a navigation bar with tabs for "Welcome", "Forms", "Data Bookmarks", "Jobs", and "Tutorials". The "Forms" tab is active, showing a sub-tab for "HCA" with a red 'x' icon. The main content area features the title "HCA 1.0.2" and the description "Hydrophobic Cluster Analysis." Below this are buttons for "Run", "Reset", and "Help pages". A section labeled "Input Data" is visible, containing a visualization of a protein sequence. The visualization is titled "query.data.seq" and "Drawn by Luc Canard". It shows a sequence of amino acids (A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y) represented by colored letters (green, blue, red, black) and symbols (squares, circles, stars) arranged in a complex, folded pattern. The sequence is indexed from 10 to 310. The bottom of the page shows a "Black and white" option and a "41" page number.

Predikce sekundární struktury

- Predikce 3 základních typů: H (helix), E (β -list), C/- (smyčka/vše ostatní)
- 1. GENERACE
 - *ab-initio*
 - Vycházela z fyzikálně-chemických vlastností a ze statistik pro jednotlivé aminokyseliny
- 2. GENERACE
 - Zahrnuje i vliv okolních aminokyselin
- 3. GENERACE
 - *Homology-based models*
 - Metody strojového učení
 - Využívá multiple sequence alignmentu a toho, že 2D struktura je více konzervovaná než sekvence

Metody strojového učení (*Machine learning*)

- Model, který je natrénovaný na známé sadě dat
- Neuronové sítě
- Skryté Markovovy modely



PSIPRED

2D

- Predikce sekundární struktury pomocí 2 neuronových sítí
- Časově náročnější
- Ve srovnání s většinou programů na predikci sekundární struktury má lepší výsledky

<http://bioinf.cs.ucl.ac.uk/psipred/>

Choose prediction methods

Popular Analyses

- PSIPRED 4.0 (Predict Secondary Structure)
- MEMSAT-SVM (Membrane Helix Prediction)
- DISOPRED3 (Disopred Prediction)
- pGenTHREADER (Profile Based Fold Recognition)

Contact Analysis

- DeepMetaPSICOV 1.0 (Structural Contact Prediction)
- MEMPACK (TM Topology and Helix Packing)

Fold Recognition

- GenTHREADER (Rapid Fold Recognition)
- pDomTHREADER (Protein Domain Fold Recognition)

Structure Modelling

- Bioserf 2.0 (Automated Homology Modelling)
- DMPfold 1.0 Fast Mode (Protein Structure Prediction)
- Domserf 2.1 (Automated Domain Homology Modelling)

Domain Prediction

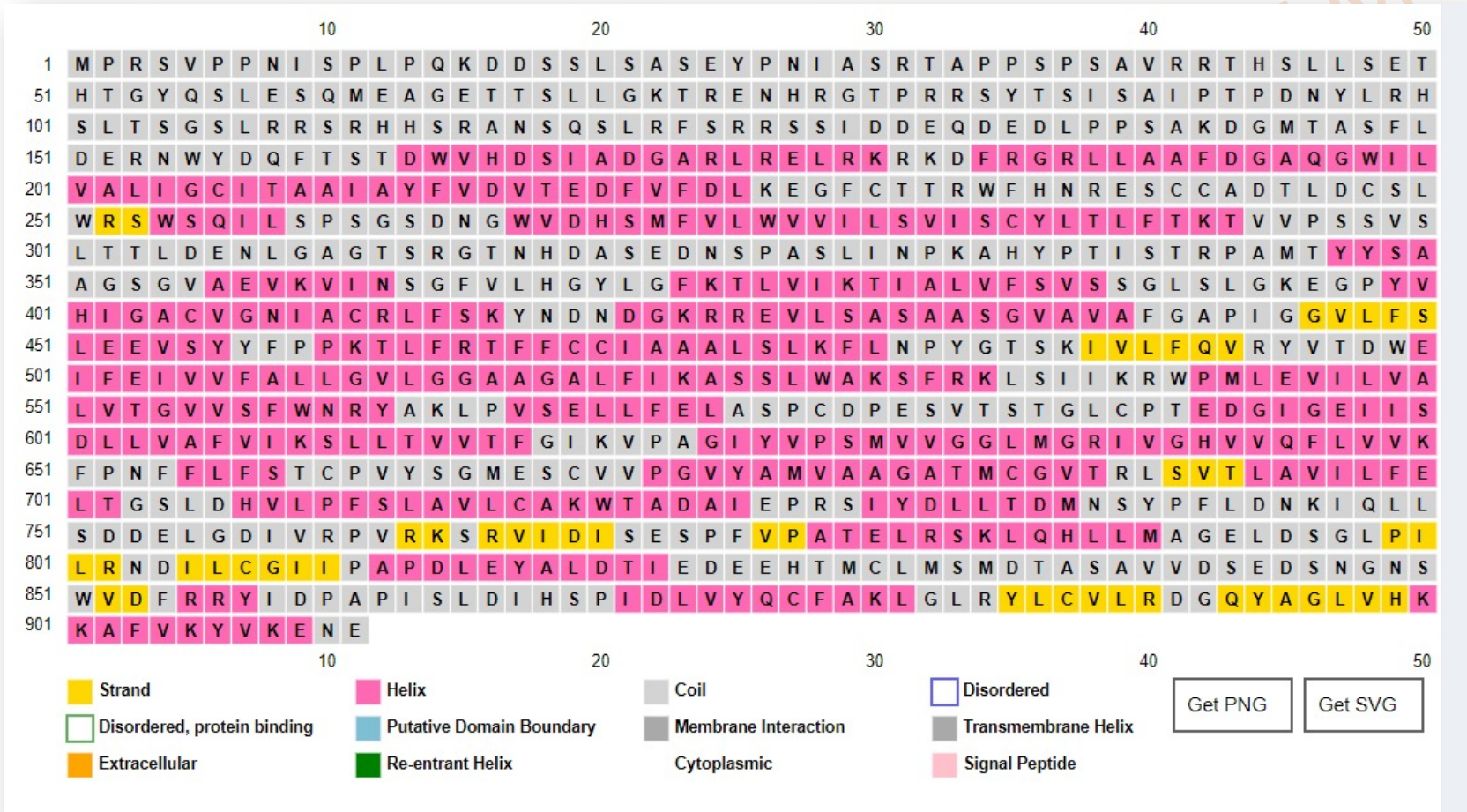
- DomPred (Protein Domain Prediction)

Function Prediction

- FFPred 3 (Eukaryotic Function Prediction)

[Help](#)

PSIPRED



Rozšíření predikce 2D struktury

- Predikce **více typů** 2D struktury (dle DSSP – Database of Secondary Structure Assignments)
 - α -helix (H)
 - 3_{10} -helix (G)
 - π -helix (I)
 - β -řetězec, extended strand (E)
 - β -bridge (B)
 - turn (T)
 - bend (S)
 - ostatní, coil (C)
- Predikce **přístupnosti solventu**
- Predikce **transmembránových helixů**

Predikce terciární struktury

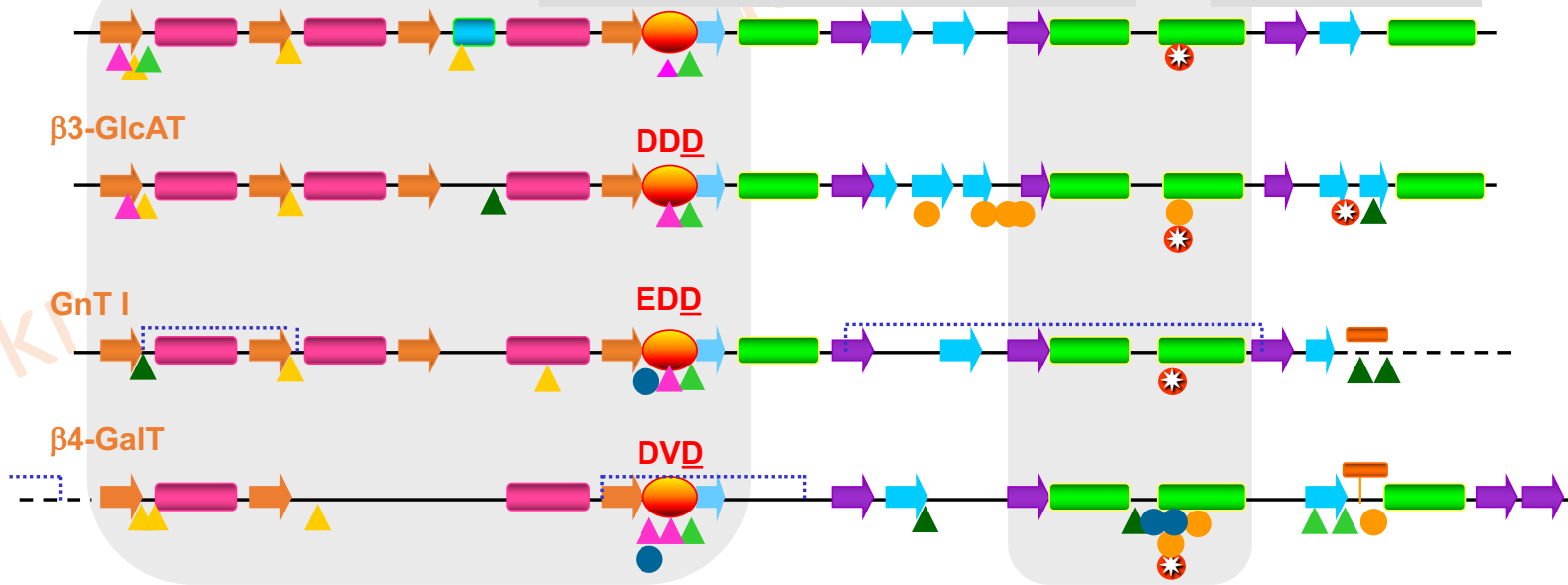
- Klasifikace proteinů
 - Předpověď funkce
 - Vytvoření modelu pro další studium
-
- ***Ab initio***
 - **Homologní modelování**
 - **Threading („navlékání“)**



Metody pro predikci funkce

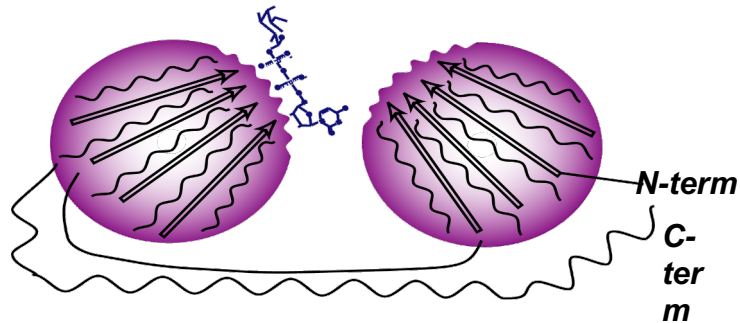
„klasické“ metody: vícenásobné aminokyselinové přiložení
 pozitivní alignment pouze mezi sekvencemi stejné rodiny

Enzym	Gen	Organismus	Sekvence
Gal α 1,4-Gal β -R	LgtD	<i>N. men</i>	CDKVLVYLDIDVLRVRSITPLWDTDLGDNWLGACID YFNAGVLLINLKKWR
Glc α 1,3-Glc α -R	RfaI	<i>E. coli</i>	APKVLVYLDADIICQGTIEPLINFSFPDDKVA MVVT YFNSGFLLINTAQWA
Gal α 1,3-Glc α -R	RfaI	<i>S. typh</i>	QIKVLVYLDADIACKGSIQELIDLNFAENEIAAVVA YFNAGFILIXIPLWT
Glc α 1,2-Glc α -R	RfaJ	<i>E. coli</i>	LDRLLYLDADVCKGDISQLLHLGLN-GAVAAVVK YFNSGVVYLDLKKWA
Gal α 1,6-Man α -R	LpcA	<i>R. leg</i>	IERLLYLDADVLA VSPVDELFTRNFAQKALAAVDD YFNAGVLLFDWSACR
Glc α 1,3-Man α -R	DUGT	<i>D. mel</i>	VRKIIFVDADAIVRTDIKELYDMDLGGAPYAYTPF YHISALYVVDLKRFR



Dvě pozorované topologie 3D struktur glykosyltransferas

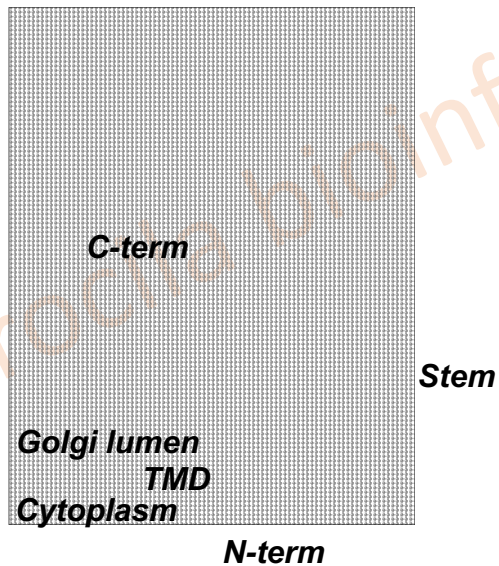
BGT-fold



(Prokaryotes/Phage)

β -GlcT (BGT, phage T4)	n.c.	inv
β 4-GlcNAcT (MurG, <i>E.coli</i>)	GT28	inv
β -GlcT (GtfB, <i>M. orientalis</i>)	GT1	inv

SpsA-fold



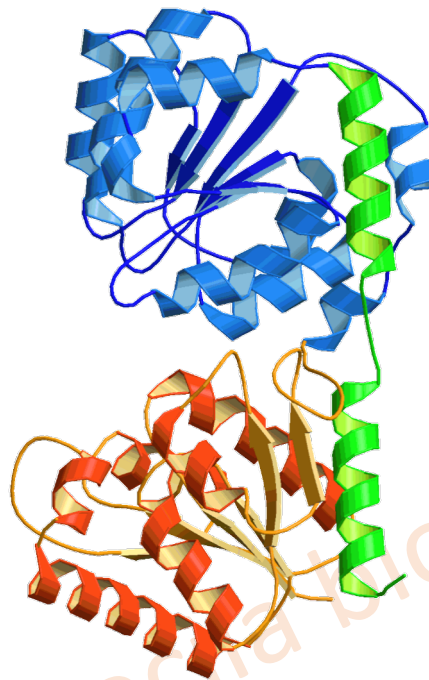
(Prokaryotes)

SpsA (<i>B. subtilis</i>)	GT2	inv
α 4-GalT (LgtC, <i>N.meningitis</i>)	GT8	ret

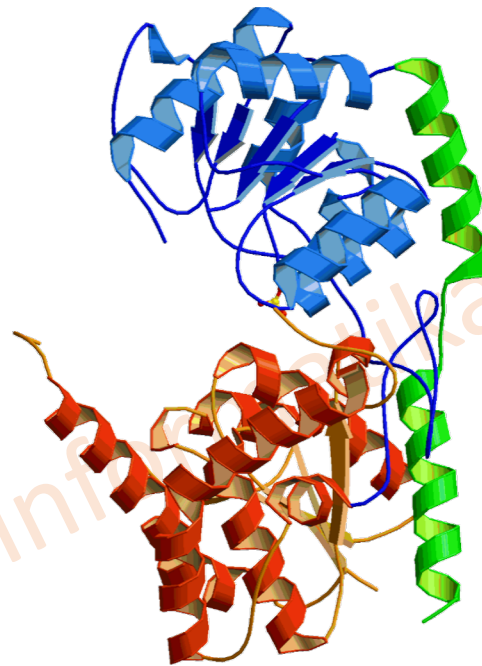
(Eucaryotes)

β 4-GalT1 (bovine)	GT7	inv
β 2-GlcNAcT (GnT I, rabbit)	GT13	inv
β 3-GlcAT I (human)	GT43	inv
α 3-GalT (bovine)	GT6	ret
Glycogenin (rabbit)	GT8	ret
α 3-GalNAcT (GTA, human)	GT6	ret
α 3-GalT (GTB, human)	GT6	ret

Nadrodina s BGT foldem



MurG (β -GlcNAcT)
GT28
E. coli
Ha *et al.*, 2000



GtfB (β -GlcT)
GT1
A. orientalis
Mulichak *et al.*, 2001

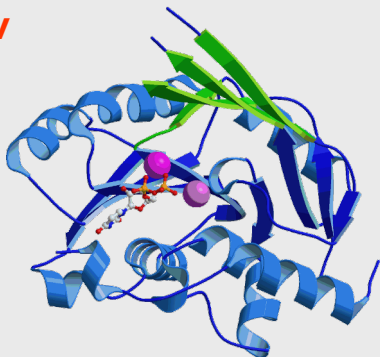


BGT (β -GlcT)
n.c.
Phage T4
Vrieling *et al.*, 1994

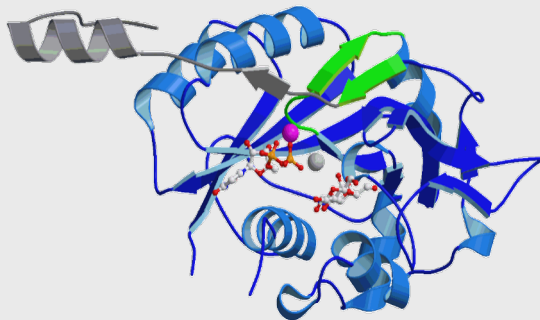
Nadrodina s SpsA foldem

Společná NBD

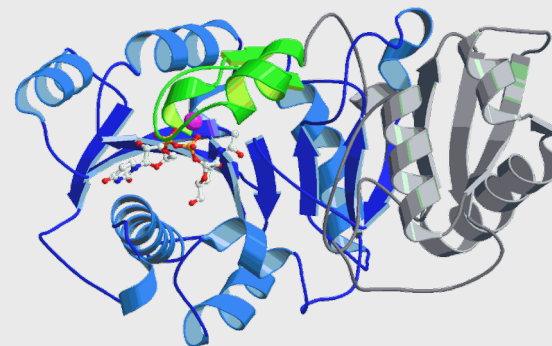
Inv



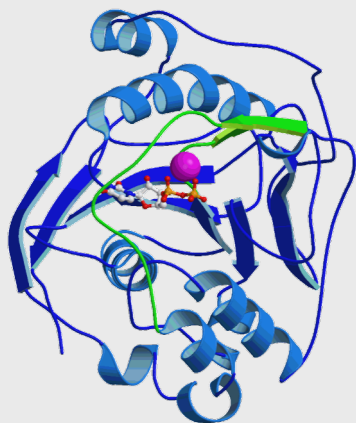
SpsA [GT2]
Charnok *et al*, 1999, 2001



Hum β3-GlcAT [GT43]
Pedersen *et al*, 2000

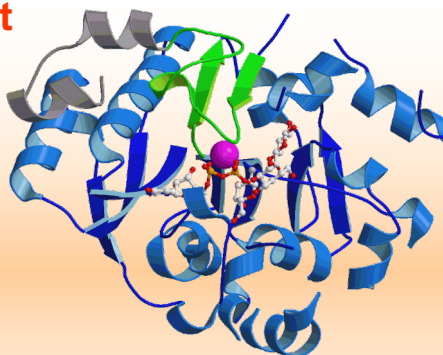


Rabbit GnT I [GT13]
Ünlügil *et al*, 2000

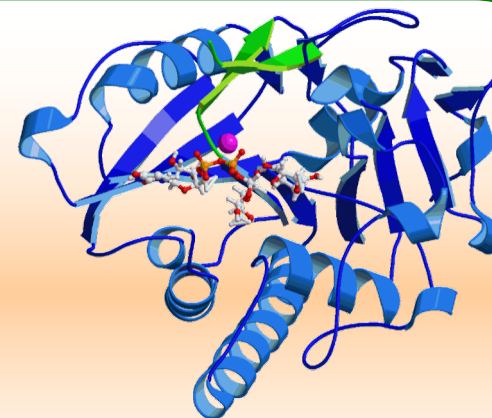


Bovine β4-GalT [GT7]
Gastinel *et al*, 1999
Ramakrishnan *et al*, 2001, 2002

Ret



LgtC (α4-GalT) [GT8]
Neisseria meningitidis
Persson *et al*, 2001



Bovine α3-GalT [GT6]
Gastinel *et al*, 2001
Boix *et al*, 2001, 2002

Predikce terciární struktury

- Klasifikace proteinů
 - Předpověď funkce
 - Vytvoření modelu pro další studium
-
- ***Ab initio***
 - **Homologní modelování**
 - **Threading („navlékání“)**



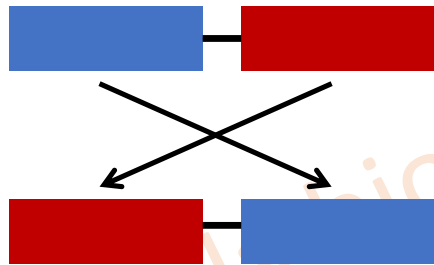
Threading



- Porovnává možnost přiložení sekvence na proteiny známých **foldů**
- „navlékání“ = rozpoznání a přiřazení proteinového foldu aminokyselinové sekvenci
- S využitím strukturních databází (PDB, SCOP, CATH) je vytvořena databáze existujících foldů - sekvence je porovnávána s touto databází (3D profilů) a na jejich základě jsou konstruovány 3D-modely
- 3D profil - každému reziduu v 3D struktuře je přiřazena environmentální proměnná (obsah polárních atomů v postranním řetězci, skrytá plocha, sekundární elementy, apod.) vycházející z předpokladu, že okolí rezidua je více konzervováno než aminokyselina samotná.
- Reziduum může být také popsáno pomocí svých interakcí
- Výsledná kvalita modelu shoda je popsána pomocí Z-skóre nebo energie
- **U multidoménných struktur je potřeba aminokyselinovou sekvenci rozdělit na jednotlivé domény a analyzovat je separátně**

PLLSASIVSAPVVTSETYVDIPGLYLDVAKAGIRDGKLQVILNVPTPYATGNNFPGIYFAIATNQG VVADGCFTYSSKV
 PESTGRMPFTLVATIDVGSVTFVKGQWKSVRGSAMHIDSYASLSAIWGTAA PSSQGSNGQAETGGTGAGNIG
 GGERDGT FNLPPHIKFGVTALHAANDQTIDIYIDDDPKPAATFKGAGA QDQNLGTKVLDSGNGRVRVIMANGR
 PSRLGSRQVDIFKKS YFGIIGSEDGADDDYNDGIVFLNWPLG

ERDGT FNLPPHIKFGVTALHAANDQTIDIYIDDDPKPAATFKGAGA QDQNLGTKVLDSGNGRVRVIMANGRPSR
 LGSRQVDIFKKS YFGIIGSEDGADDDYNDGIVFLNWPLG PLLSASIVSAPVVT SQTYVDIPGLYLDVAKAGIRDGKLQ
 VILNVPTPYATGNNFPGIYFAIATNQG VVADGCFTYSSKVPESTGRMPFTLVATIDVGSVTFVKGQWKSVRGSAM
 HIDSYASLSAIWGTAA PSSQGSNGQAETGGTGAGNIGGGGKLA AALEIKRASQPELAPEDPEDVEHHHHHH



#	=====	#
EMBOSS_001	1 -----	0
EMBOSS_001	1 ERDGT FNLPPHIKFGVTALHAANDQTIDIYIDDDPKPAATFKGAGA QDQ	50
EMBOSS_001	1 -----	0
EMBOSS_001	51 NLGTVLDSGNGRVRVIMANGRPSRLGSRQVDIFKKS YFGIIGSEDGAD	100
EMBOSS_001	1 -----PLLSASIVSAPVVTSETYVDIPGLYLDVAKAGIRD	35
EMBOSS_001	101 DDYNDGIVFLNWPLGPLLSASIVSAPVVT SQTYVDIPGLYLDVAKAGIRD	150
EMBOSS_001	36 GKLVILNVPTPYATGNNFPGIYFAIATNQG VVADGCFTYSSKVPESTGR	85
EMBOSS_001	151 GKLVILNVPTPYATGNNFPGIYFAIATNQG VVADGCFTYSSKVPESTGR	200
EMBOSS_001	86 MPFTLVATIDVGSVTFVKGQWKSVRGSAMHIDSYASLSAIWGTAA PSSQ	135
EMBOSS_001	201 MPFTLVATIDVGSVTFVKGQWKSVRGSAMHIDSYASLSAIWGTAA PSSQ	250
EMBOSS_001	136 GSGNQAETGGTGAGNIGGGGERDGT FNLPPHIKFGVTALHAANDQTID	185
EMBOSS_001	251 GSGNQAETGGTGAGNIGGGG-----	271
EMBOSS_001	186 IYIDDDPKPAATFKGAGA QDQNLGTVLDSGNGRVRVIMANGRPSRLGS	235
EMBOSS_001	272 -----KLAAA-----LEIK-----RAS-----	283
EMBOSS_001	236 RQVDIFKKS YFGIIGSEDGADDDYNDGIVFLNWPLG	271
EMBOSS_001	284 -QPE-----LAPEDPEDVEHHH-----HHH	302

Threading

PHYRE2 (3D-PSSM)

<http://www.sbg.bio.ic.ac.uk/phyre2>

Threading at 2D level and scoring at 3D level :

matching of secondary structure elements, and propensities of the residues in the query sequence to occupy varying levels of solvent accessibility

The PSIPRED Protein Sequence Analysis Workbench

<http://bioinf.cs.ucl.ac.uk/psipred/>

GenTHREADER Rapid fold recognition, matching your sequence against a library of whole PDB chains.

pGenTHREADER Highly sensitive fold recognition using profile-profile comparison (whole chain library).

pDomTHREADER Highly sensitive homologous domain recognition using profile-profile comparison (domain library).

I-TASSER

<https://zhanglab.ccmb.med.umich.edu/I-TASSER/>

a hierarchical approach to protein structure and function prediction. It first identifies structural templates from the PDB by multiple threading approach LOMETS, with full-length atomic models constructed by iterative template fragment assembly simulations. Function insights of the target are then derived by threading the 3D models through protein function database BioLiP.

Phyre2



- Server pro 3D predikci struktur pomocí **threadingu**
- Vysoce výkonný – poměrně spolehlivá detekce foldu i při **nízké homologii** (i pod 15%)

<http://www.sbg.bio.ic.ac.uk/phyre2/>

#	Template	Alignment Coverage	3D Model	Confidence	% i.d.	Template Information
1	c1r1zB	 Alignment		100.0	31	PDB header: sugar binding protein Chain: B; PDB Molecule: ergic-53 protein; PDBTitle: the crystal structure of the carbohydrate recognition2 domain of the glycoprotein sorting receptor p58/ergic-533 reveals a novel metal binding site and conformational4 changes associated with calcium ion binding
2	c2a6yA	 Alignment		100.0	21	PDB header: sugar binding protein Chain: A; PDB Molecule: emp47p (form1); PDBTitle: crystal structure of emp47p carbohydrate recognition domain2 (crd), tetragonal crystal form
3	d2a6za1	 Alignment		100.0	20	Fold: Concanavalin A-like lectins/glucanases Superfamily: Concanavalin A-like lectins/glucanases Family: Lectin leg-like
4	c2dupB	 Alignment		100.0	42	PDB header: protein transport Chain: B; PDB Molecule: vesicular integral-membrane protein vip36; PDBTitle: crystal structure of vip36 exoplasmic/luminal domain, metal-free form

Phyre2

Phyre2

ARDLVIPMIYCGHG

User sequence



Homologous
sequences

Search the 10 million known
sequences for homologues
using PSI-Blast.

Pokrocila@informatics NCBR PrF MU

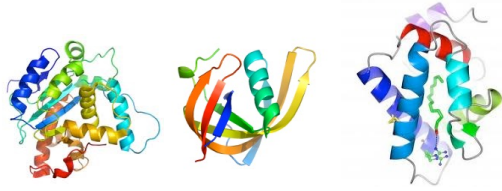
Phyre2



Capture the mutational propensities at each position in the protein

An evolutionary fingerprint

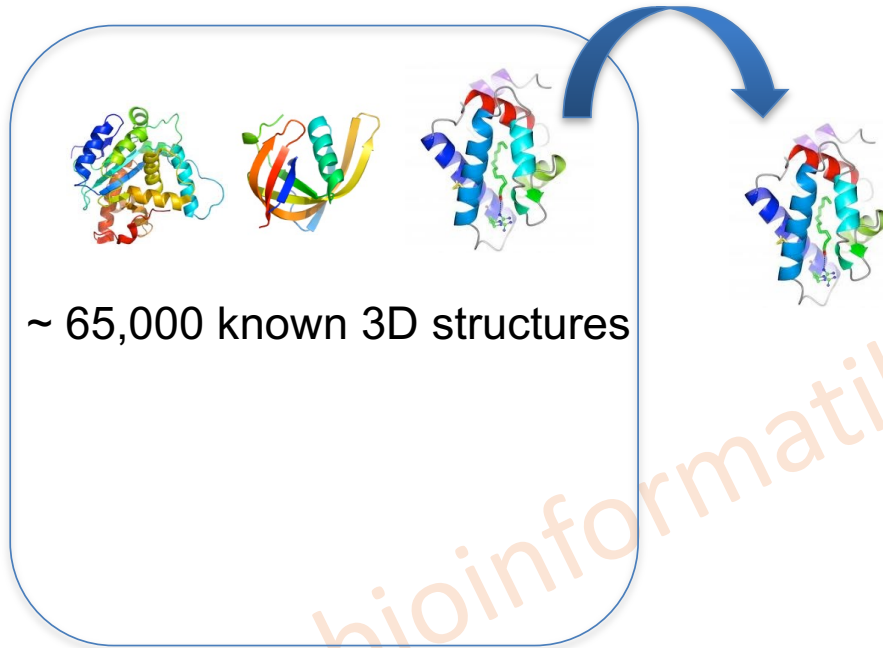
Phyre2



~ 65,000 known 3D structures

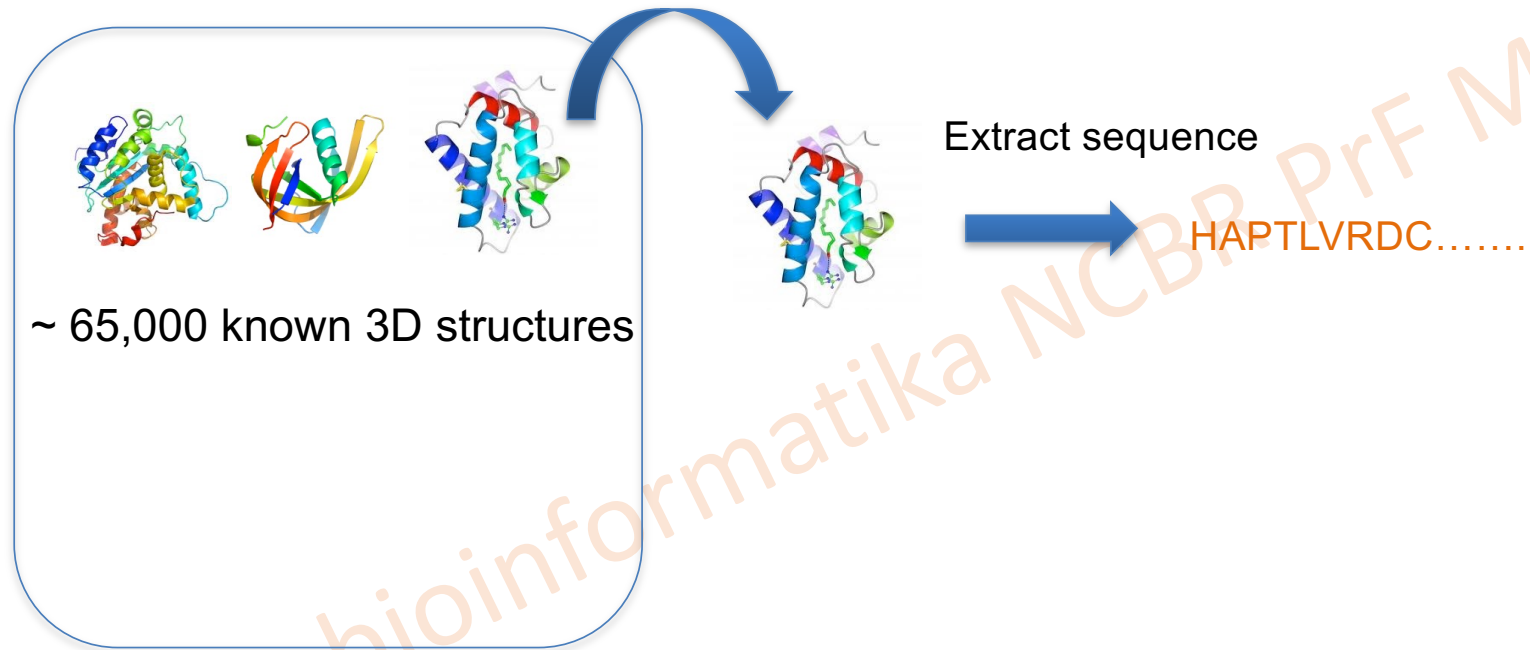
Pokročila bioinformatika NCBR PrF MU

Phyre2



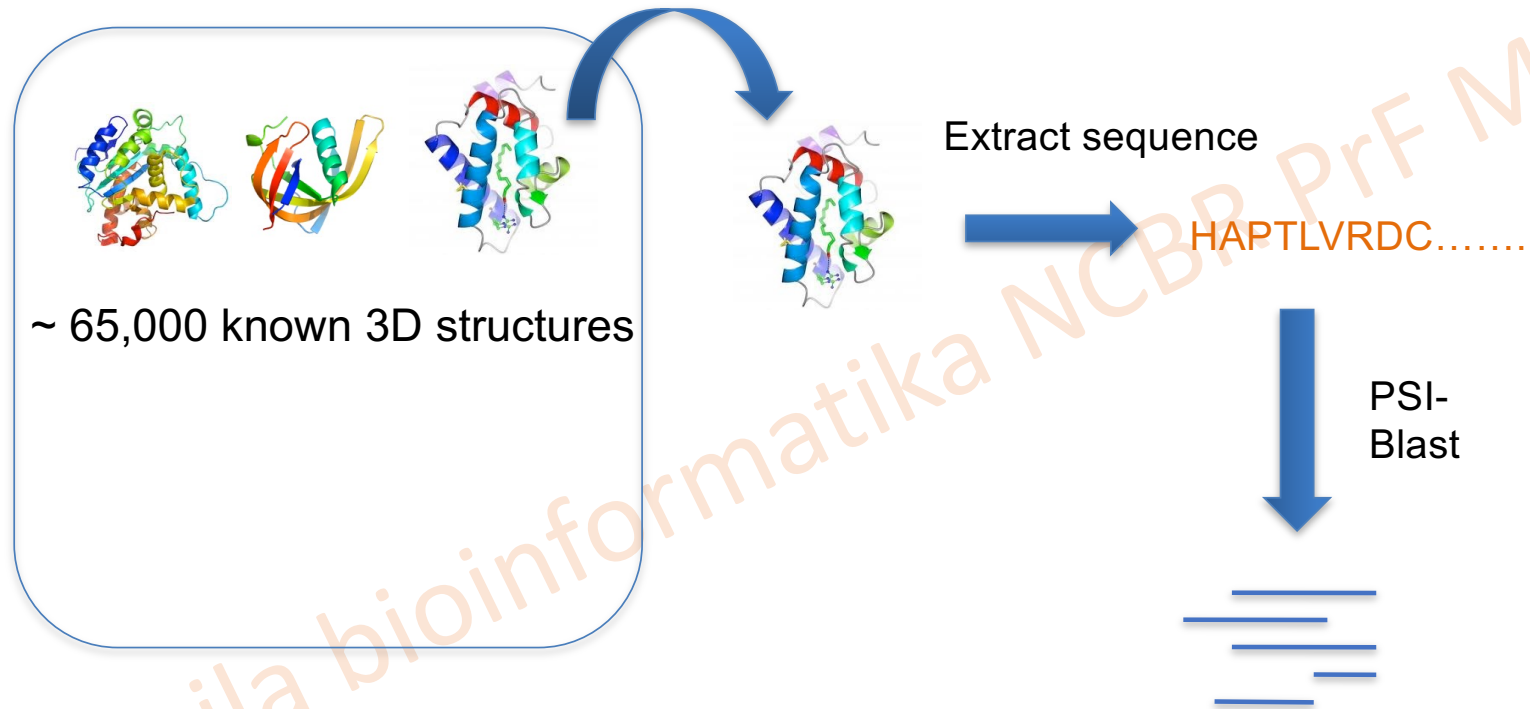
Pokročila bioinformatika NCBR PrF MU

Phyre2



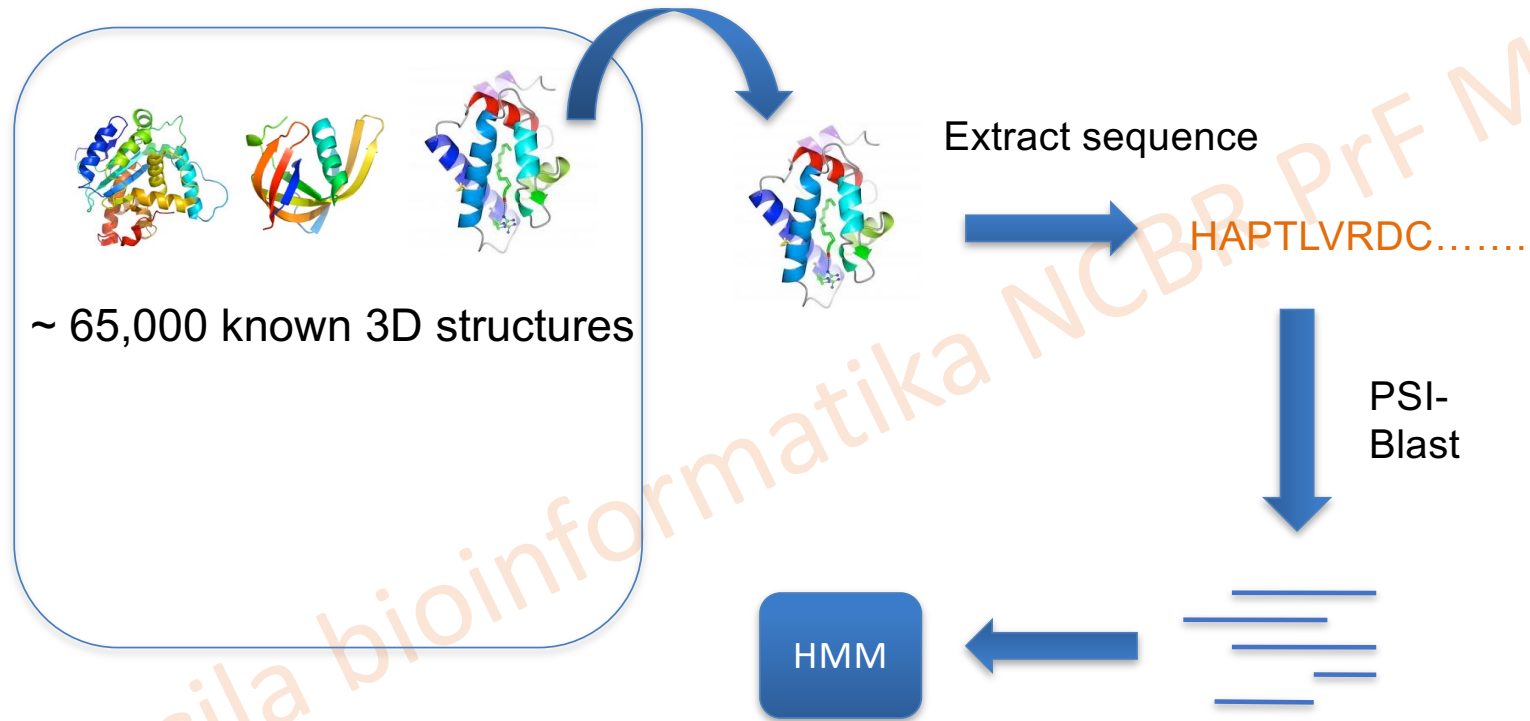
Pokrocila bioinformatika NCBB PrF MU

Phyre2



Pokrocila bioinformatika NCBP PrF MU

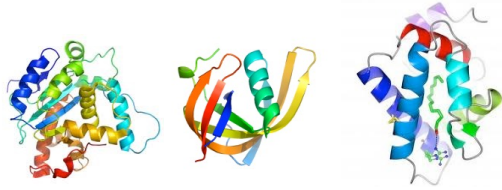
Phyre2



Hidden Markov model
for sequence of KNOWN structure

Pokrocila bioinformatika NCBP PrF MU

Phyre2



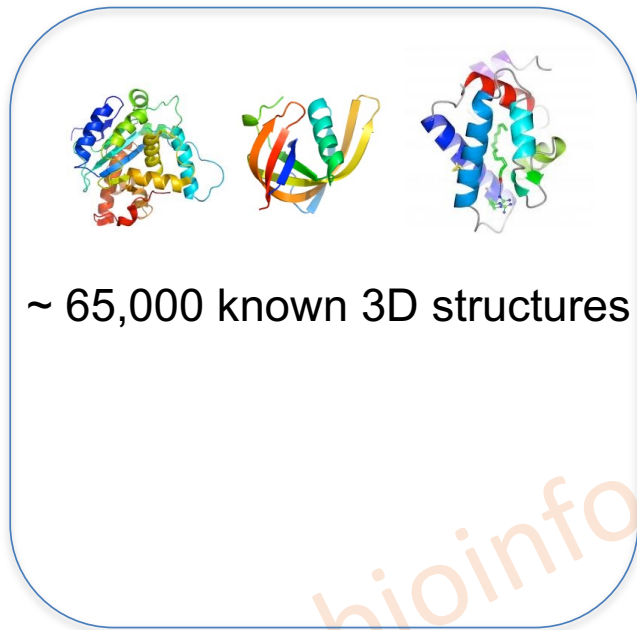
~ 65,000 known 3D structures



~ 65,000 hidden Markov models

Pokročila bioinformatika NCBH FMU

Phyre2

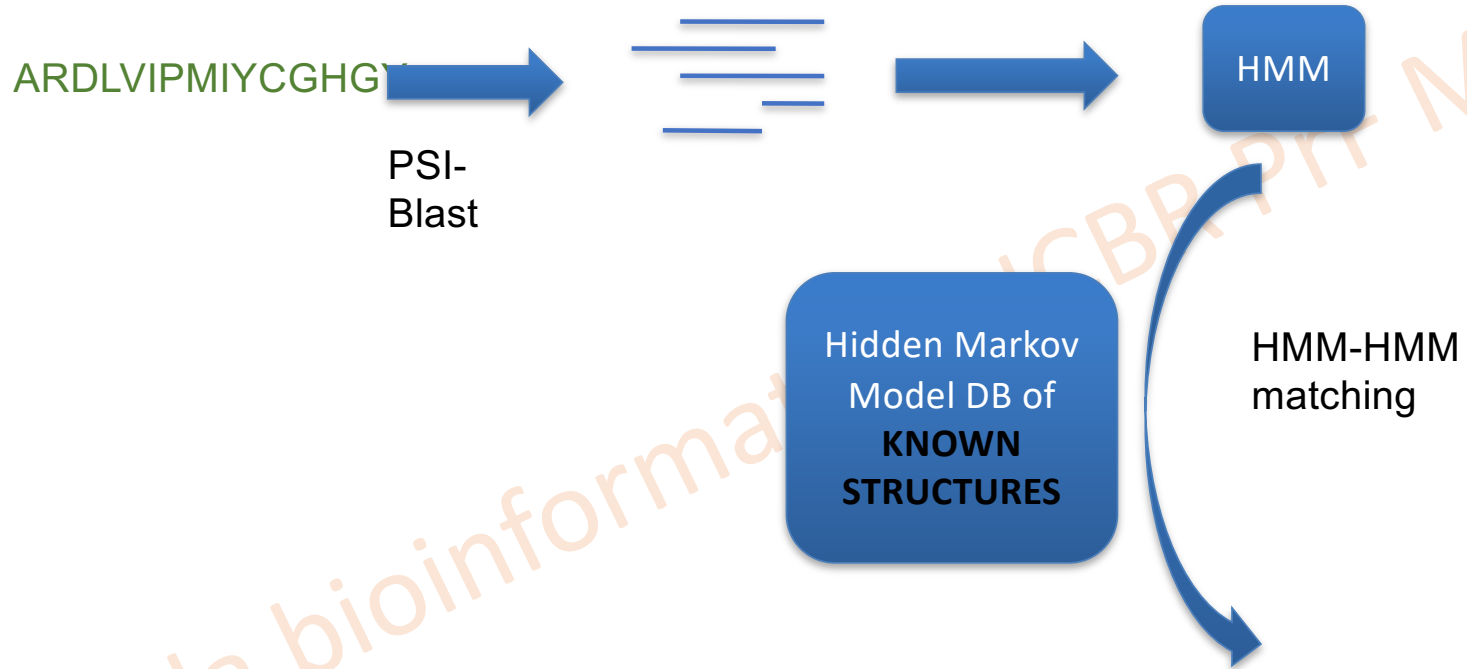


Hidden Markov Model
Database of
**KNOWN
STRUCTURES**

Pokročila bioinformatika

MU

Phyre2



Alignments of user sequence to known structures ranked by confidence.

ARDL--VIPMIYCGHGY
AFDLCDLIPV--CGMAY

Sequence of known structure

Phyre2

ARDLVIPMIYCGHGY

PSI-
Blast



HMM

Hidden Markov
Model DB of
**KNOWN
STRUCTURES**

HMM-HMM
matching

3D-Model



ARDL--VIPMIYCGHGY
AFDLCDLIPV--CGMAY

Sequence of known structure

bioinformatica

Phyre2

ARDLVIPMIYCGHGY

PSI-
Blast



HMM

**Very powerful –
able to reliably detect extremely
remote homology**

**Routinely creates accurate models even
when sequence identity is <15%**

Hidden Markov
Model DB of
**KNOWN
STRUCTURES**

HMM-HMM
matching

3D-Model



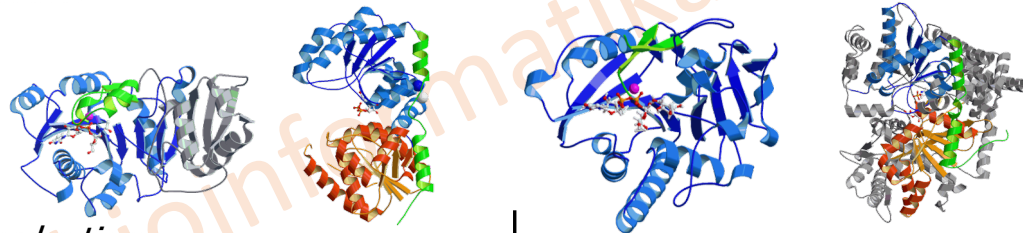
ARDL--VIPMIYCGHGY
AFDLCDLIPV--CGMAY

Sequence of known structure

SDVDIEAGQTLVQVVNISNGETWVAIQLP AQYRSFDLVFENVSPSTSGSVLVAQMAPQSGGVYGSNYS
GSGWGN DLGGGGFYGYSEAKWMCLWPANRSGPNSKTGIYGTCKLMNLNQSN AVPSVTSNLFAPTAY
KNEPGYANVGGCCQKIRGLASSIQFAFALHGGNVPQNTDTFSGGTIKVYGVWN

*3D-fold calculation based
on known structures*

Model quality evaluation



pair
residue-residue
interactions

surface
residue-solvent
interactions

pair/surface
residue-residue and
residue-solvent interaction


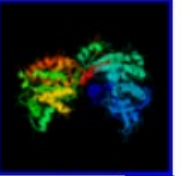


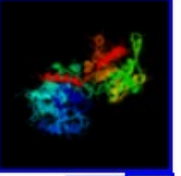
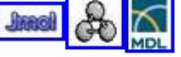

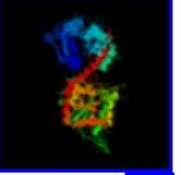

“Quality” scores

Glykogensynthasa – rodina GT3 (v rodině v době analýzy nebyla vyřešena 3D-struktura)

[http://www.sbg.bio.ic.ac.uk/phyre/qphyre_output/95cbaa7600a9bfff/su
mmary.html](http://www.sbg.bio.ic.ac.uk/phyre/qphyre_output/95cbaa7600a9bfff/su
mmary.html)

To predict functional residues and GO classification, try [ConFunc](#)

Fold Recognition

View Alignments	SCOP Code	View Model	E-value	Estimated Precision	BioText	Fold/PDB descriptor	Superfamily
	d2bisa1 (length: 437) 18% i.d.	 	3.9e-36	100 %	n/a	UDP-Glycosyltransferase/glycogen phosphorylase	UDP-Glycosyltransferase/glycogen phosphorylase
	d1rzua (length: 477) 14% i.d.	 	6.1e-36	100 %	n/a	UDP-Glycosyltransferase/glycogen phosphorylase	UDP-Glycosyltransferase/glycogen phosphorylase
	c3c48A (length: 438) 11% i.d.	 	6.1e-31	100 %	n/a	PDB header: transferase	Chain: A: PDB Molecule: predicted glycosyltransferases;



A co protein, který nemá v sekvenčních databázích žádný homolog

Pokročila bioinformatika I/CBR PrF MU

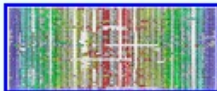
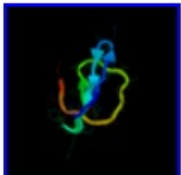


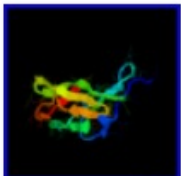

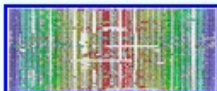


RS-20L

No sequence homology
in databases !



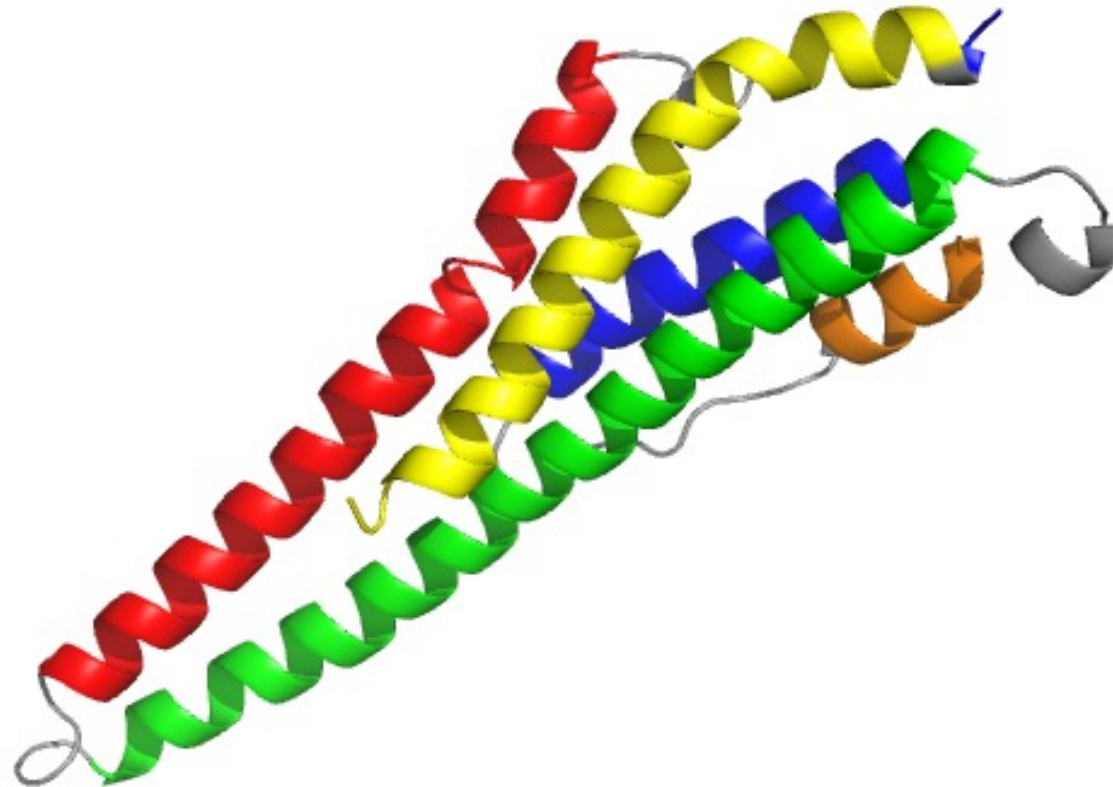
P

Fold Recognition

View Alignments	SCOP Code	View Model	E-value	Estimated Precision	BioText	Fold/PDB descriptor	Superfamily	Family	(beta-test)
	d1eh9a2 (length:67) 24% i.d.	 	50	0 %	n/a	Glycosyl hydrolase domain	Glycosyl hydrolase domain	alpha-Amylases, C-terminal beta-sheet domain	n/a
	c2fsdA (length:142) 19% i.d.	 	50	0 %	n/a	PDB header: virus/viral protein	Chain: A: PDB Molecule: putative baseplate protein;	PDB Title: a common fold for the receptor binding domains of 2 lactococcal phages? the crystal structure of the head3 domain of phage bil170	n/a
	c2ct4A (length:70) 11% i.d.	 	56	0 %	n/a	PDB header: signaling protein	Chain: A: PDB Molecule: cdc42-interacting protein 4;	PDB Title: solution structure of the sh3 domain of the cdc42-2 interacting protein 4	n/a

AB2L structure overview

Structure: 4 helical bundle



Pokroc

F MU

Top model

Model (left) based on template
[d2ja9a1](#)

Top template information

Fold:OB-fold

Superfamily:Nucleic acid-binding
proteins

Family:Cold shock DNA-binding
domain-like

Confidence and coverage

Confidence: **24.1%** Coverage: **20%**

38 residues (20% of your sequence)
have been modelled with 24.1%
confidence by the single highest
scoring template.



You may wish to submit your sequence
to [Phyrealarm](#). This will automatically
scan your sequence every week for
new potential templates as they
appear in the Phyre2 library.

Please note: You must be registered
and logged in to use Phyrealarm.

3D viewing

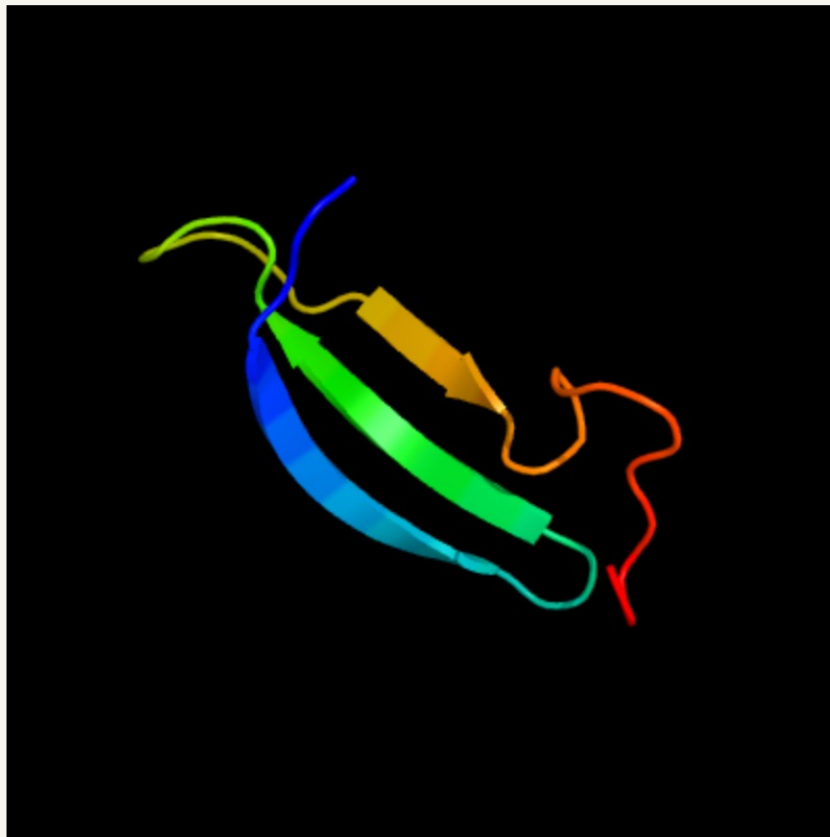


Image coloured by rainbow N → C terminus

Model dimensions (Å): **X:24.236 Y:23.853 Z:38.403**

Prozkoumání možností a principů fungování I-TASSERu
bude domácím úkolem

Pokročila bioinformatika NCBR PrF MU

Homologní modelování



- Je založeno na existenci blízkého **strukturního homologu** (typicky 50 % sekvenční podobnosti a více, minimálně 30%)
- Využívá skutečnosti, že dva proteiny ze stejné rodiny a s podobnou sekvencí mají i podobnou 3D strukturu
- Kromě sekvence našeho proteinu potřebujeme znát strukturu homologního proteinu = **templát**
- Pro vysoce homologní sekvence je spolehlivost velmi vysoká

MODELLER

Mostly used program in academic environment for serious homology modeling

SWISS-MODEL

An automated knowledge-based protein modelling server

Homologní modelování



1. **Alignment** zadané sekvence a sekvence templátu
2. Extrakce proteinové **páteře** ze struktury templátu a umístění **postranních řetězců**
3. **Modelování** otoček a smyček
4. **Minimalizace** energie
5. **Validace** namodelované struktury

Swiss-Model



- Výběr modelu (manuální, automatický)
- Podle vybraného modelu pak predikuje strukturu zadané sekvence
- Součástí výstupu je sada parametrů hodnotících **kvalitu** modelu. Při využití více templátů je tak možno porovnat jednotlivé modely

<http://swissmodel.expasy.org/>

The screenshot displays the Swiss-Model web interface. On the left, the 'Start a New Modelling Project' form is visible, including fields for 'Target Sequence(s)', 'Project Title', and 'Email', along with an 'Upload Target Sequence File...' button and a 'Search For Templates' button. The 'Template Results' table shows a list of templates with columns for Name, Title, Coverage, GMQE, OSQE, Identity, Method, and Oligo State. The 'Model Results' panel on the right provides a detailed view of the selected model, including a 3D structure, quality estimates (GMQE 0.99, OMEAN -1.17), and a comparison plot. The 'Model-Template Alignment' section shows the sequence alignment between the model and the template.

Sort	Name	Title	Coverage	GMQE	OSQE	Identity	Method	Oligo State
2ezy	1.A	BARRIER-TO-AUTOINTEGRATION FACTOR	0.99	0.98	100.00	NMR	homo-dimer	
2ezx	1.A	BARRIER-TO-AUTOINTEGRATION FACTOR	0.99	0.98	100.00	NMR	homo-dimer	
2ezy	1.A	BARRIER-TO-AUTOINTEGRATION FACTOR	0.99	0.70	100.00	NMR	homo-dimer	
2odp	1.A	Barrier-to-autointegration factor	0.99	0.53	100.00	NMR	hetero-trimer	
2bzf	1.A	BARRIER-TO-AUTOINTEGRATION FACTOR	0.99	0.66	100.00	X-ray, 2.9Å	homo-dimer	
6phd	1.B	Barrier-to-autointegration factor	0.98	0.66	95.45	X-ray, 2.1Å	hetero-tetramer	
6pxl	1.B	barrier to autointegration factor (BAF)	0.98	0.67	95.40	X-ray, 2.3Å	hetero-tetramer	
2zud	1.B	DNA repair and recombination protein rcsA	0.40	-	22.03	X-ray, 3.2Å	homo-dimer	
2bka	1.A	DNA REPAIR AND RECOMBINATION PROTEIN RADA	0.40	-	22.03	X-ray, 3.2Å	monomer	



SWISS-MODEL

An automated knowledge-based protein modelling server

- Start SMR-Pipeline in automated mode on BC2-cluster at Thu May 2 08:51:47 2013
- Start BLAST for highly similar template structure identification
- No suitable templates found!
- Run HHSearch to detect remotely related template structures
- Unfortunately, we could not identify useful template structures
- For troubleshooting, please see our article in Nature Protocols:
 - Bordoli, L., Kiefer, F., Arnold, K., Benkert, P., Battey, J. and Schwede, T. (2009). Protein structure homology modelling using SWISS-MODEL Workspace. Nature Protocols, 4, 1.

Computation of this workunit has stopped.

Please see the following log report for details:

Started: Wed May 13 06:59:31 2009 (sms_automode) Reading user input sequence **No Templates found.**

=====

Simple automated template selection could not identify suitable templates. Please use advanced Template Selection under **[Tools]** to select a template and prepare a workunit using the project mode.

Ab initio



- Nejuniverzálnější – vychází pouze ze sekvence
- Výpočetně **nejnáročnější**
- Zahrnuje řadu kroků:
 - Predikce 2D struktury
 - Modelování jednotlivých fragmentů
 - Kombinace fragmentů navzájem
 - Doplnění smyček a flexibilních úseků
- **Nízká spolehlivost** zejm. pro větší proteiny

De novo modelling with Rossetta

(David Baker lab, Univ. of Washington)

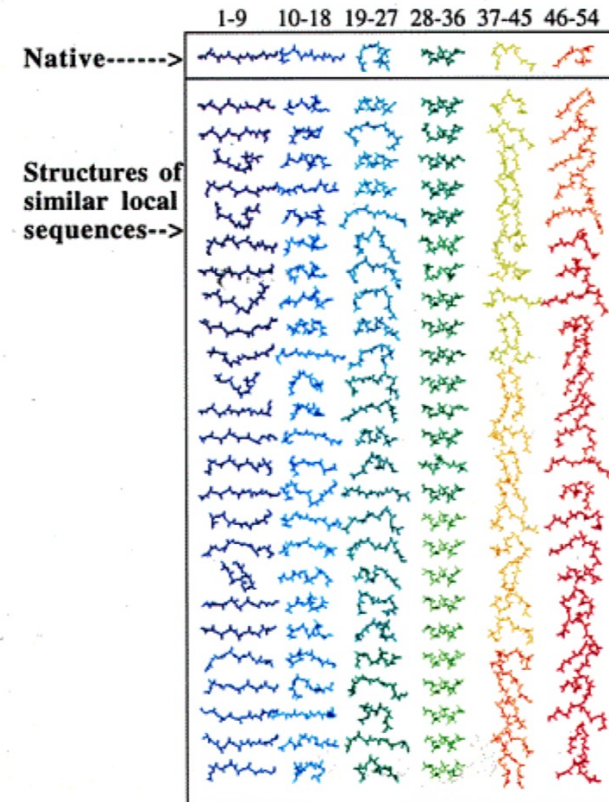
- In contrast to threading, Rosetta does *de novo* prediction – doesn't use templates/homologous structures
- instead performs Monte Carlo search through space of conformations to find minimal energy conformation

Pokročila bioinformatika NCBR PrF MU

De novo modelling with Rossetta

MAU

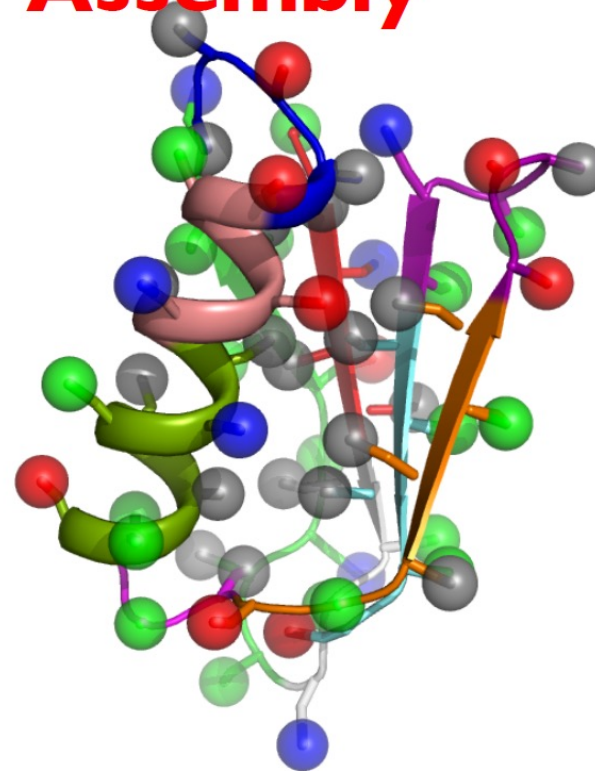
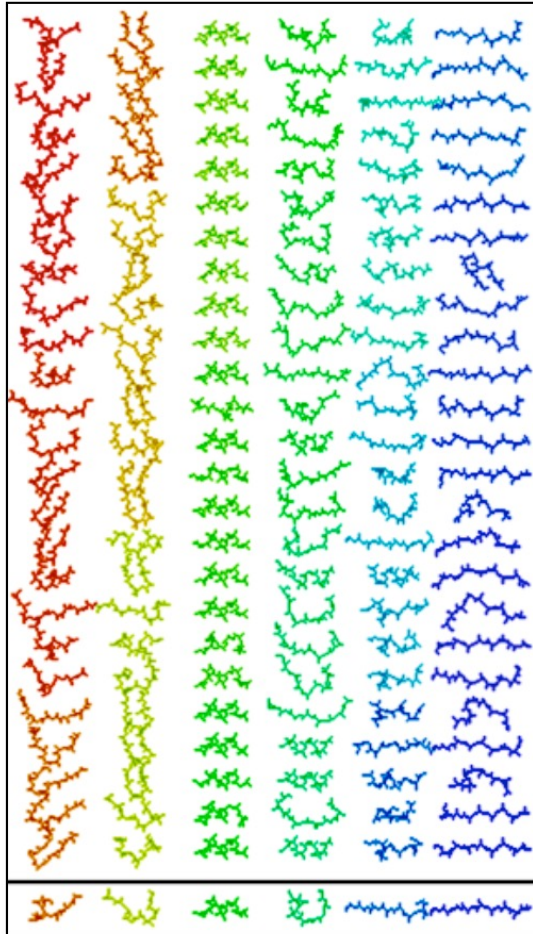
- fragments are selected from known structures
- the window-fragment matches are calculated using
 - PSI-BLAST to build a profile model of the sequence
 - the predicted secondary structure of the sequence



P

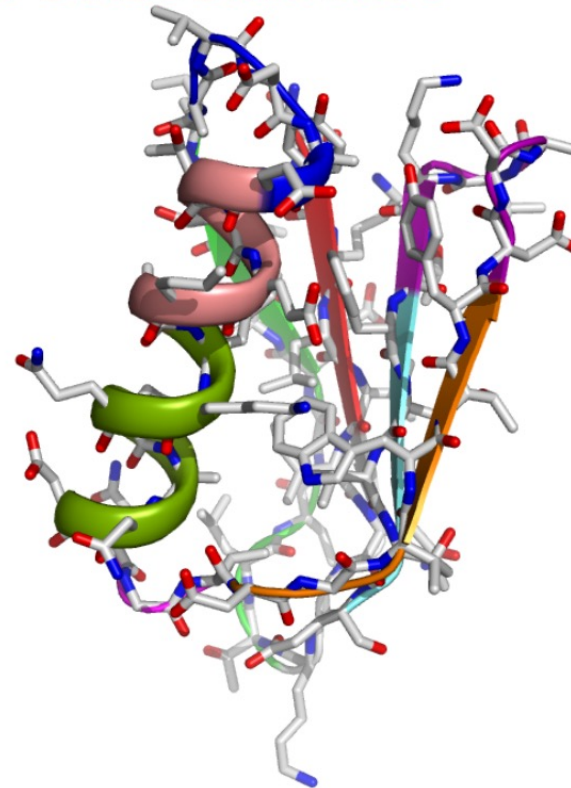
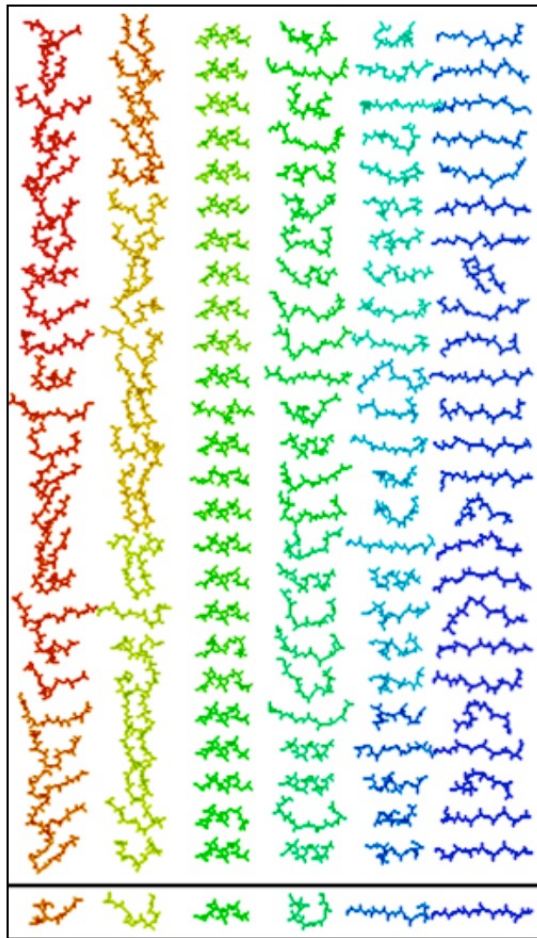
De novo Modeling with Rosetta

Stage I. Fragment Assembly



De novo Modeling with Rosetta

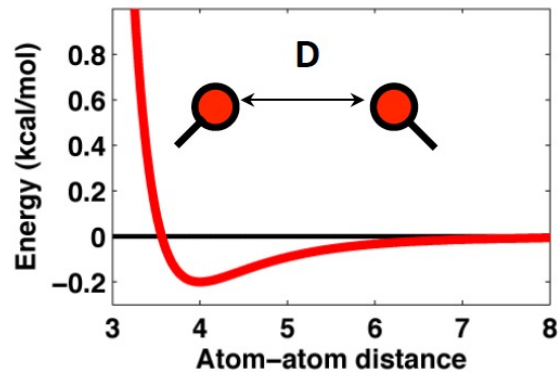
Stage II. All-atom refinement



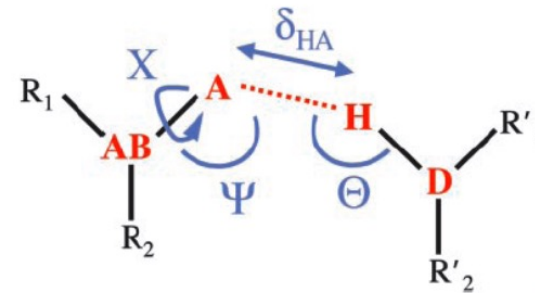
Pc

Ingredients of a high resolution potential

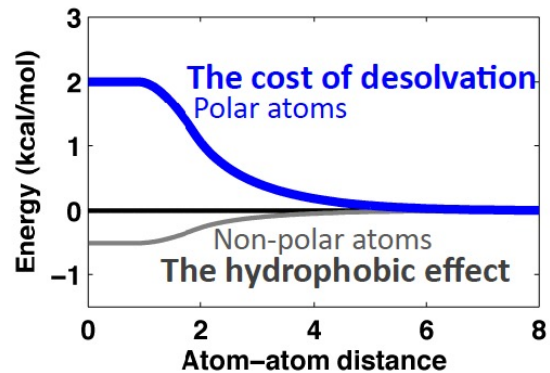
1. Van der waals packing



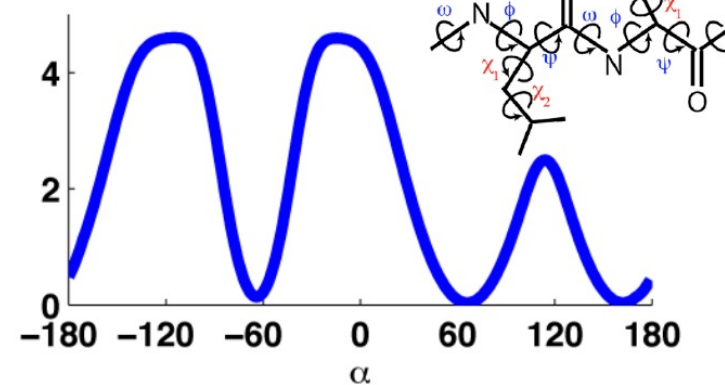
2. Hydrogen bonds



3. Manifestations of water



4. Torsional potential



Scoring Function Takes Into Account

MU

- residue environment (solvation)
- residue pair interactions (electrostatics, disulfides)
- strand pairing (hydrogen bonding)
- strand arrangement into sheets
- helix-strand packing
- steric repulsion
- etc.
- scoring function search progressively adds terms during search
 - initially on the steric overlap term is used
 - then all but “compactness” terms are used
 - etc.
- search is initiated from different random seeds

WEB server - Robetta

<http://robetta.bakerlab.org>

Response Times

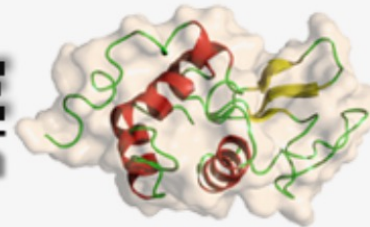
To prevent unnecessary usage we require two manual steps for full structure predictions. The first step is to submit your sequence for domain and template detection. The second step is to continue for 3-D models. You may only select one domain at a time for structure predictions. The second step is computationally expensive so please continue with this step only if necessary. You may help increase computing resources for this service by joining our distributed computing project [Rosetta@HOME](#) and spreading the word out to friends and colleagues.

- ~10 minutes - hours for domain and template detection.
- ~1 day - weeks for high accuracy homology models (templates detected with high confidence > 0.8 and sequence identity > 40%).
- ~1 week - months for difficult targets.

Zhang Lab - QUARK



QUARK ONLINE
Ab Initio Protein Structure Prediction



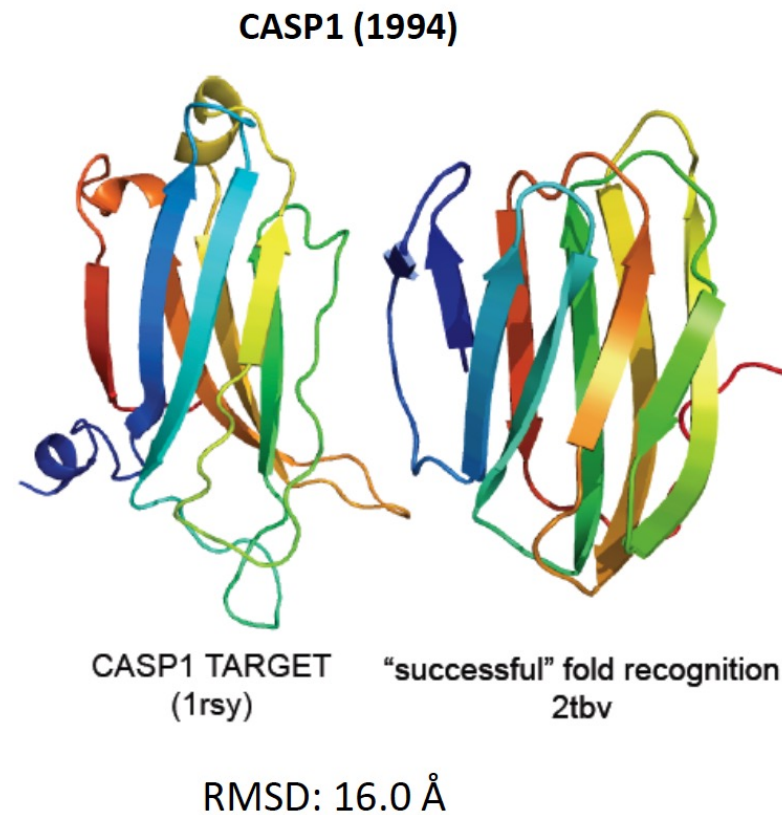
QUARK is a computer algorithm for ab initio protein structure prediction and protein peptide folding, which aims to construct the correct protein 3D model from amino acid sequence only. QUARK models are built from small fragments (1-20 residues long) by replica-exchange Monte Carlo simulation under the guide of an atomic-level knowledge-based force field. QUARK was ranked as the No 1 server in Free-modeling (FM) in [CASP9](#) and [CASP10](#) experiments. Since no global template information is used in QUARK simulation, the server is suitable for proteins that do not have homologous templates in the PDB library. Go to [example](#) to view an example of QUARK output. The server is only for non-commercial use. Questions about the QUARK server can be posted at the [Service System Discussion Board](#).

Cut and paste your sequence (in [FASTA format](#), less than 200 AA. [Example input](#))

Driving innovation in
protein structure
prediction:
“CASP”

Critical Assessment of
Structure Prediction

**Five *blind*
predictions per
target**



CASP 11 (2014)

CASP11 in numbers



Number of groups registered	208
including: expert groups	123
prediction servers	85
Number of regular targets released	100
including all-group (human) targets	55
Targets canceled for all/manual prediction	7 / 10
Number of refinement targets released	37
Number of assisted prediction targets released	71
Number of targets received from	
Joint Center for Structural Genomics (JCSG):	32
Structural Genomics Consortium (SGC):	4
Midwest Center for Structural Genomics (MCSG):	8
Northeast Structural Genomics Consortium (NESG):	5
New York Structural Genomics Research Center (NYSGRC):	6
Non-SGI research Centers and others (Others):	40
Seattle Structural Genomics Center for Infectious Disease (SSGCID):	4
NatPro PSI:Biology (NatPro):	1

PC

<http://predictioncenter.org/casp11/results.cgi>

CASP12 in numbers

Number of groups registered	192
including: <i>expert groups</i>	<i>112</i>
<i>prediction servers</i>	<i>80</i>
Number of regular targets released	82
including <i>all-group (human) targets</i>	<i>56</i>
Targets canceled and not re-released for all/manual prediction	11 / 11
Number of refinement targets released	42
Number of assisted prediction targets released	14

Prediction category	Number of groups/servers contributing	Number of models designated as 1	Total number of models
Tertiary structure predictions	128 / 43	8362	37672
Data assisted predictions	16 / 1	109	528
Residue-residue contacts	38 / 30	3077	3077
Accuracy estimation	47 / 32	3700	7400
Interface accuracy	3 / 0	65	66
Refinement	39 / 5	1457	6227
All (unique):	188 / 80	16770	54970

<http://predictioncenter.org/casp12/results.cgi>

13th Community Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction



CASP13 in numbers

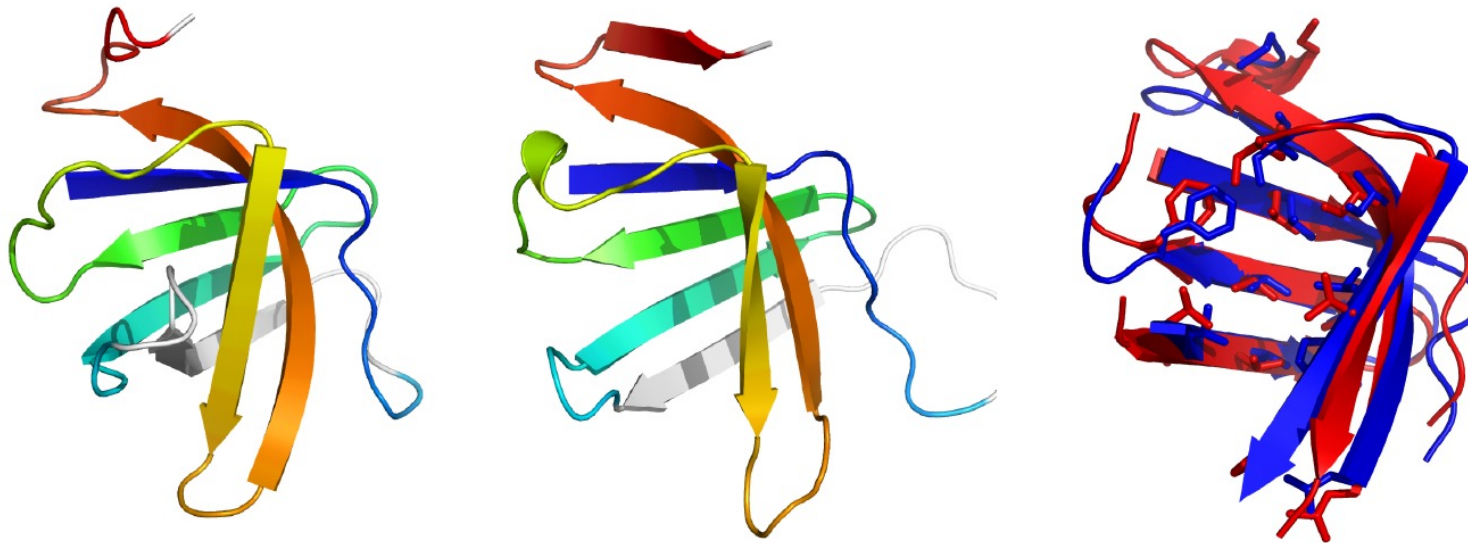
Number of groups registered	210
including: <i>expert groups</i>	123
<i>prediction servers</i>	87
Number of tertiary structure prediction targets released	90
(including <i>all-group targets</i>)	(82)
Number of hetero-multimer targets released	13
Number of refinement targets released	31
Number of assisted prediction targets released	60
Targets canceled (all / human)	(10 / 12)
Targets available/expired for manual non-QA prediction	0 / 72
Targets available/expired for server non-QA prediction	0 / 80
Targets available/expired for QA prediction	0 / 80
Targets available/expired for assisted prediction	0 / 59
Targets available/expired for multimer prediction	0 / 12

Prediction category	Number of groups/servers contributing	Number of models designated as 1	Total number of models
Tertiary structure predictions	107 / 39	7542	35982
Oligomeric predictions	40 / 9	662	2861
Data assisted predictions	24 / 5	456	2017
Residue-residue contacts	46 / 25	3914	3914
Accuracy estimation	52 / 41	4332	8687
Refinement	33 / 6	847	3788
All (unique):	185 / 87	17753	57249

<http://predictioncenter.org/casp13/results.cgi>

De novo successes: all- β

CASP7 target T0316 (domain 3)



Native

Model

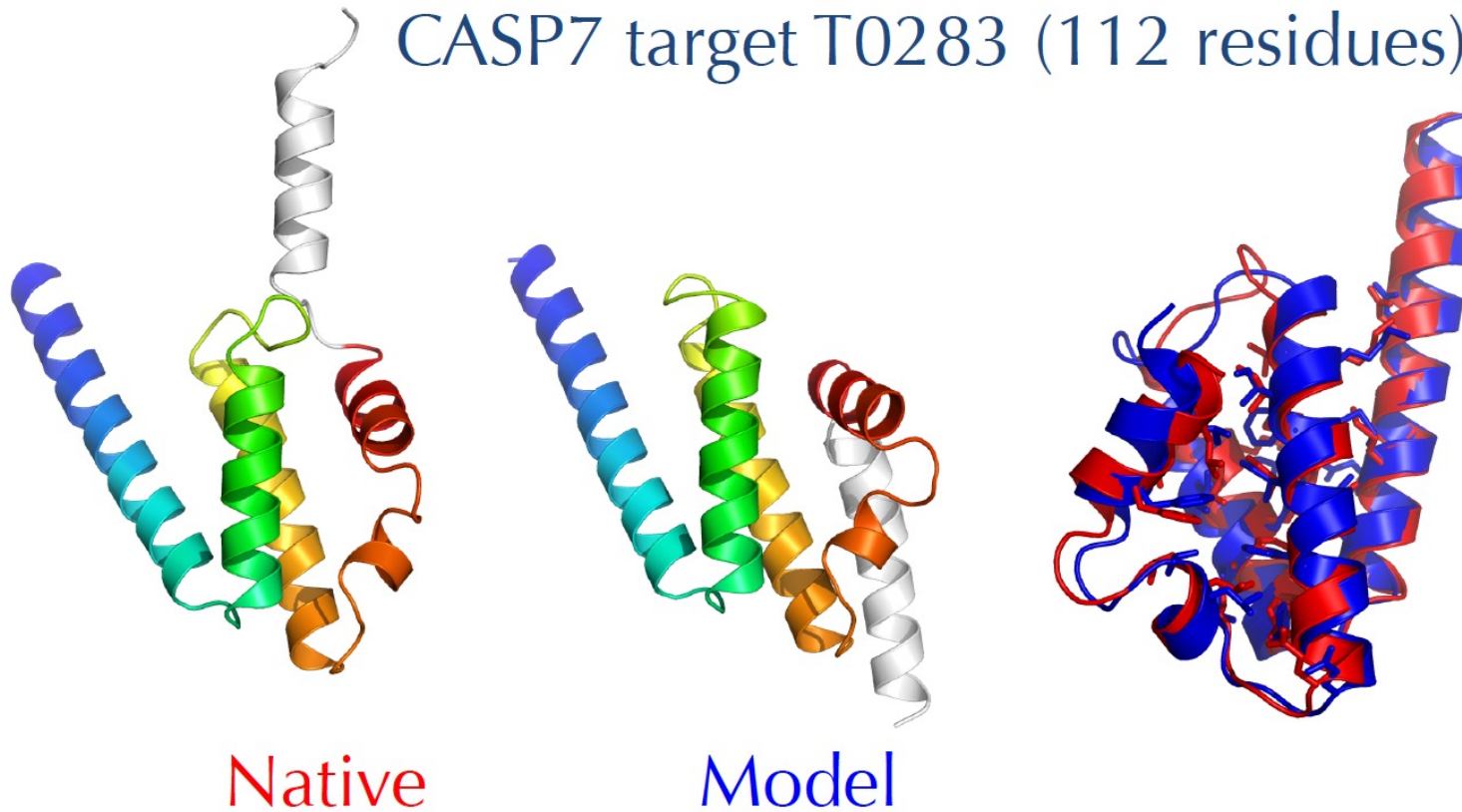
2.0 Å over 61 residues

PC

]

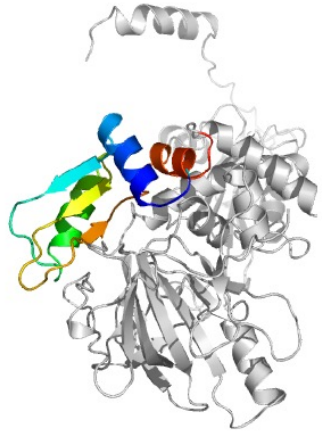
De novo successes: all- α

CASP7 target T0283 (112 residues)

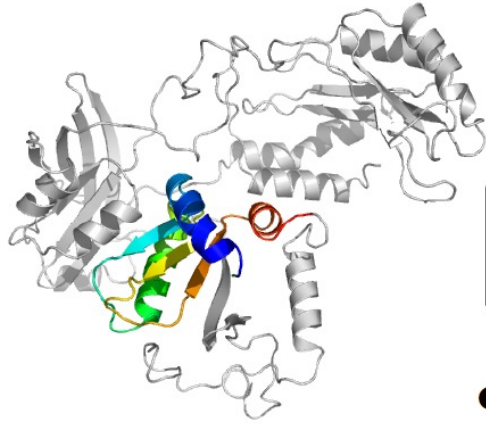


1.4 Å over 90 residues

Is protein folding *solved*?



Native

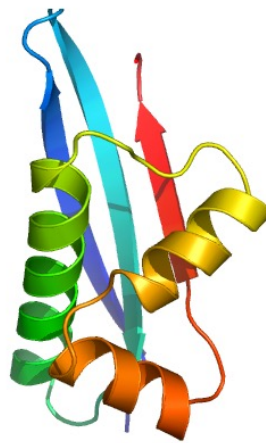
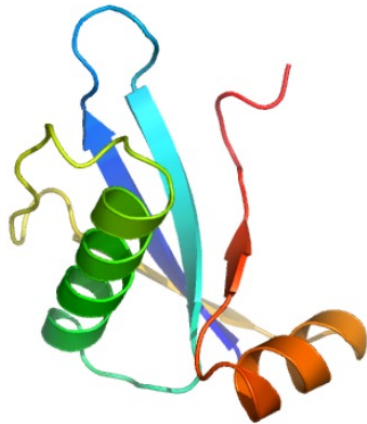


Model

NO! til now?

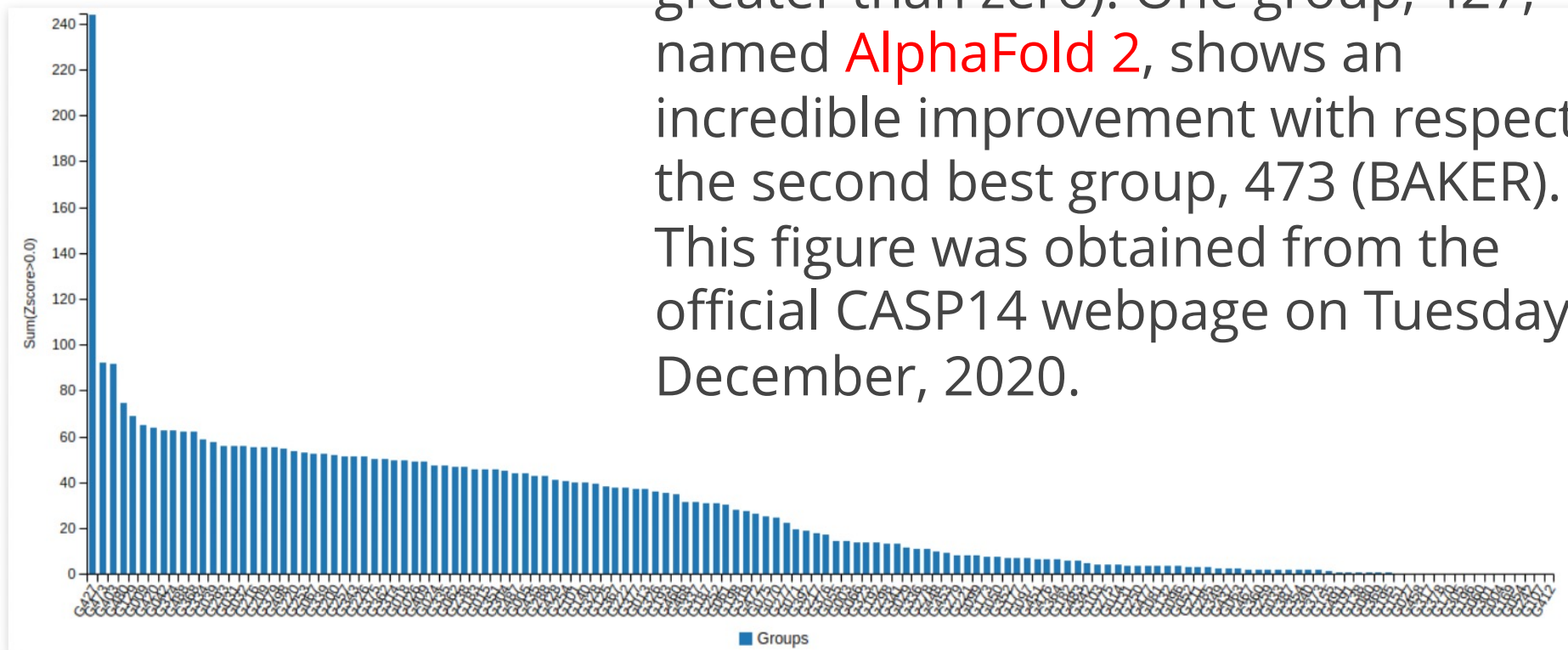
- Success in $<1/3$ of cases.
- Conformational sampling still a huge issue

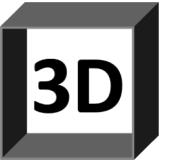
Pok'



CASP 14 (2020)

Ranking of participants in CASP14, as per the sum of the Z-scores of their predictions (provided that these are greater than zero). One group, 427, named **AlphaFold 2**, shows an incredible improvement with respect to the second best group, 473 (BAKER). This figure was obtained from the official CASP14 webpage on Tuesday 1st December, 2020.

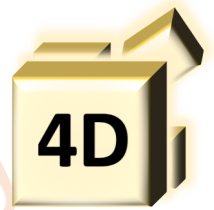




Jakou metodu zvolit?

1. Mám homologní protein se známou strukturou → homologní modelování
2. Využiji experimentální data
 - Threading
 - Kombinace více templátů pro jednotlivé části struktury
 - Různé predikční nástroje
3. *Ab initio* modelování smyček a částí sekvence bez vhodného templátu
4. Mám unikátní sekvenci – *ab initio*

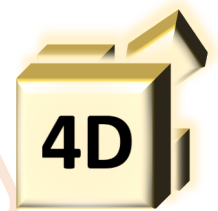
Predikce kvartérní struktury



Zahrnuje různé úrovně, např.:

- Predikce vazebných míst
 - Predikce aminokyselin podílejících se na interakci
 - Odhad oligomerního stavu
 - Protein-protein docking (protein-nukleová kyselina docking)
- SW dosud často nedokonalý, **nízká spolehlivost** predikce
- Složitější postupy většinou nejsou automatizované

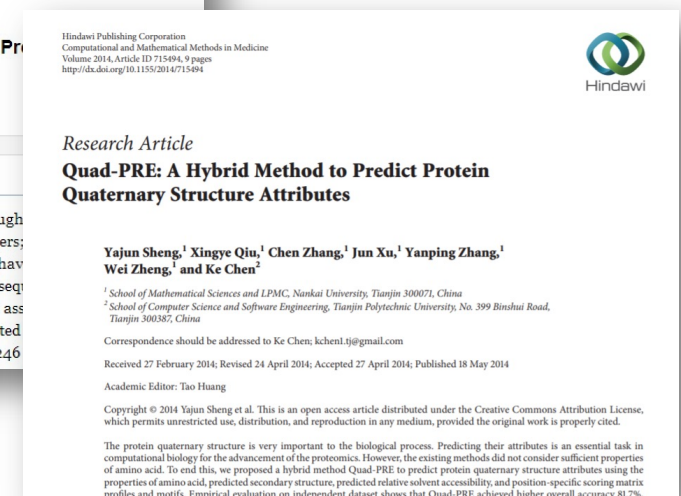
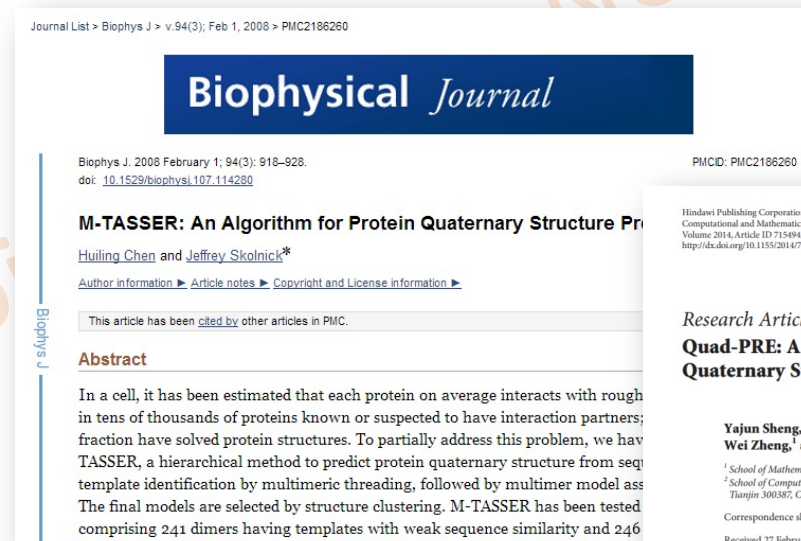
Predikce kvartérní struktury



Programy většinou vycházejí z podobnosti sekvence a/nebo 3D struktury se známými proteiny

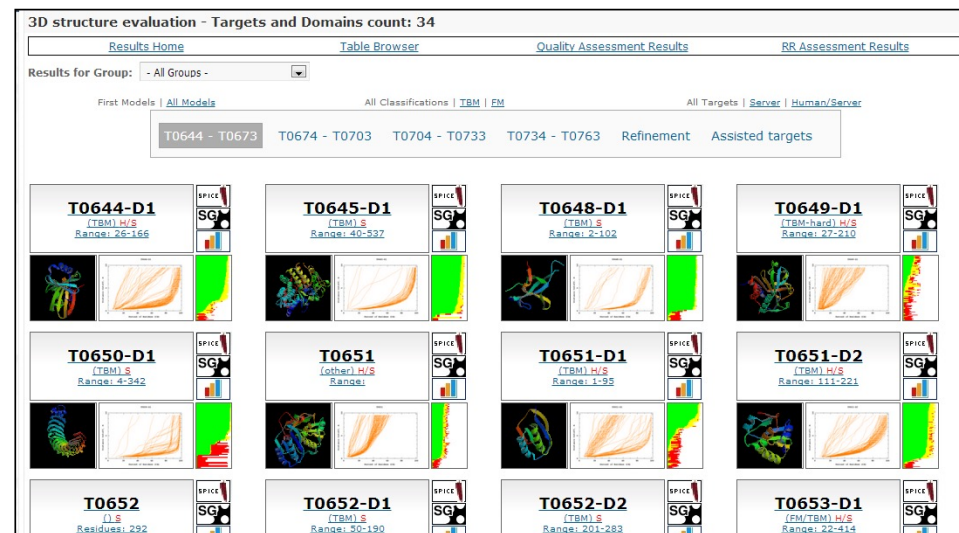
Příklady SW:

- QuatIdent
- QuaBingo
- M-TASSER
- Quad-PRE



Hodnocení kvality predikčních nástrojů - CASP

- *Critical Assessment of Techniques for Protein Structure Prediction*
- 2020 – CASP14
- Predikce vyřešených, ale zatím nepublikovaných struktur
- **Rozsáhlá analýza predikčních programů**
 - Predikce terciárních struktur
 - Identifikace neuspořádaných oblastí
 - Funkční predikce (predikce vazebných míst)
 - Interakce mezi doménami, podjednotkami a proteiny
 - Hodnocení spolehlivosti



Ale!

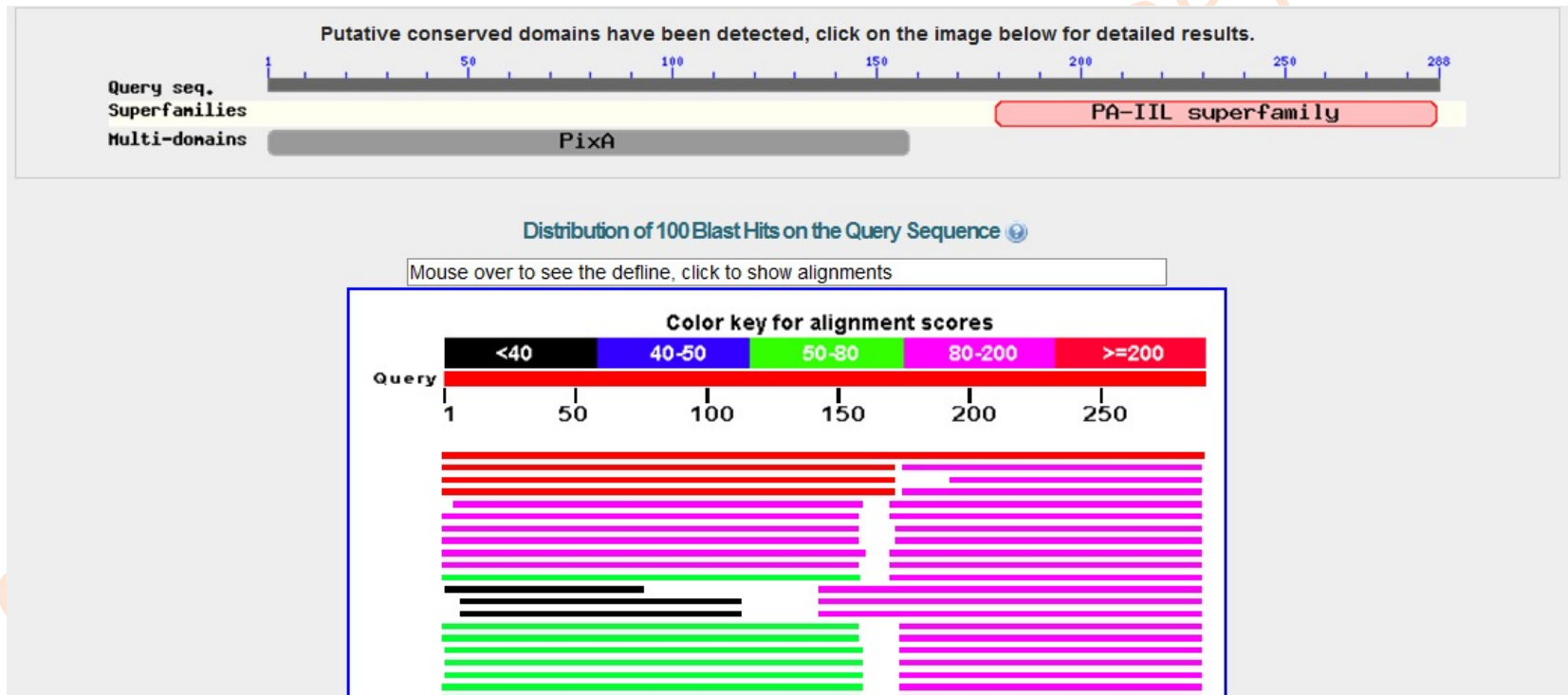
! pozor na domény !

Pokročila bioinformatika NCBR PrF MU

NCBI – Blast (Basic Local Alignment Search Tool) (National Centre for Biotechnology Information)

Prohledávání databází známých aminokyselinových sekvencí

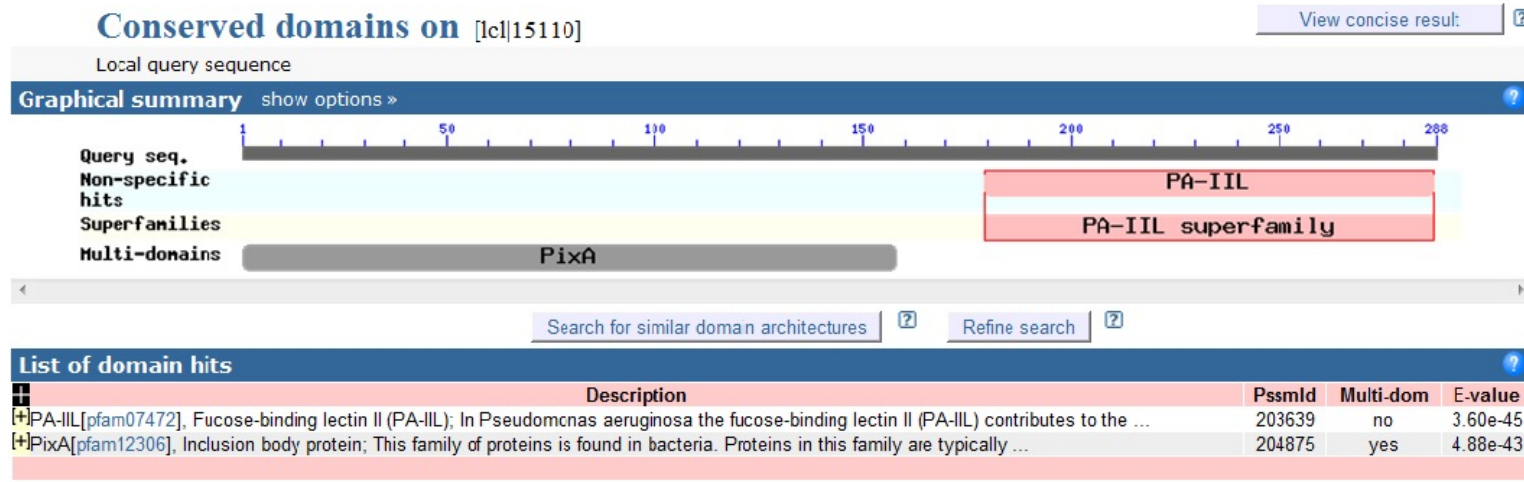
➤ celý protein



NCBI – Blast

Prohledávání databází známých aminokyselinových sekvencí

➤ celý protein



NCBI – Blast

Prohledávání databází známých aminokyselinových sekvencí

➤ celý protein

Conserved Domains

pfam07472: PA-IIL

Fucose-binding lectin II (PA-IIL)

In *Pseudomonas aeruginosa* the fucose-binding lectin II (PA-IIL) contributes to the pathogenic virulence of the bacterium. PA-IIL functions as a tetramer when binding fucose. Each monomer is comprised of a nonstranded, antiparallel beta-sandwich arrangement and contains two calcium cations that mediate the binding of fucose in a recognition mode unique among carbohydrate-protein interactions.

PubMed References

Structural basis for oligosaccharide-mediated adhesion of *Pseudomonas aeruginosa* in the lungs of cystic fibrosis patients. *Nat Struct Biol*. 2002 Dec 9;12(12):915-921

pfam07472 is a member of the superfamily cl06486.

Sequence Alignment

Format: Compact Hypertext Row Display: up to 10 Color Bits: 2.0 bit Type Selection: the most diverse members

10QX_A	6	FILPANTSGVIAFAGAAATQIVLVDS	VVK	ATFISGGISK	[1].LGS	[2].LNSGS	GAIK	63
q1_81658026	7	FILPARINFGVIVLVSSAATQIVVIFVDS	EPR	AAFSGVIGEN	[1].LGT	[2].INSGS	GNVK	64
q1_75465512	234	FQLPBNIKLSLAYSNTBQIVIVYIDD	QLV	DFLISGVNSV	LGF	[2].YSSST	GNVC	290
q1_123466540	14	FSIPFWDFRAIFFAGAAEQNIKLFIDD	SGE	[2].AYKLTTRDGP	[1].EAT	LNSGN	GNIR	71
q1_123570095	157	FSLPFWTFYGAIFYAGAADRQHLKLFIDD	APK	[2].ATFVNSKIDGV	[1].LFT	LNSGS	GNIR	244
q1_123585156	174	FHLPPNKFQVIALTSAANDQIVDIYIDD	NPX	[2].ATFKAGVQIQ	[1].LGT	[2].LDSGN	GNVK	233
2XRA_A	7	FHLPPNKFQVIALTSAANDQIVDIYIDD	DPK	[2].ATFKAGVQIQ	[1].LGT	[2].LDSGN	GNVK	66
2BOI_A	6	FILPARINFGVIVLVSSAATQIVVIFVDS	EPR	AAFSGVIGEN	[1].LGT	[2].INSGS	GNVK	63
q1_107102593	2	FILPANTSGVIAFAGAAATQIVLVDS	ETA	ATFSQGTNSA	[1].LGT	[2].LNSGS	[1].GNVQ	60
ZVNV_A	14	FSIPFWDFRAIFFAGAAEQNIKLFIDD	[2].EPA	AYKLTTRDGP	[1].EAT	LNSGN	GNIR	71

NCBI – Blast

Prohledávání databází známých aminokyselinových sekvencí

➤ celý protein

pfam12306: PlxA

Inclusion body protein
This family of proteins is found in bacteria. Proteins in this family are typically between 173 and 191 amino acids in length. PlxA is thought to be specifically produced in *Xenorhabdus nematophila*. It is an inclusion body protein.

Links
Statistics
Structure

PubMed References

Analysis of the PlxA inclusion body protein of *Xenorhabdus nematophila*. *J. Bacteriol.* 2006 Apr; 188(7):2706-2710

pfam12306 is classified as a model that may span more than one domain.
pfam12306 is not assigned to any domain superfamily.

Sequence Alignment

Reformat: Format: Compact Hypertext Row Display: up to 10 Color Bits: 2.0 bit Type Selection: the most diverse members

q1	123655921	2	- [2] -	NIWDLVLTIDVDT	IKL	[17] -	S	[2] -	PTQL	[4] -	SNG	[7] -	VHWVARD	[7] -	GSELAVALRQGD	84	
q1	123484895	13	- [2] -	QSIQILAVIDTDY	INK	[10] -	N		PTGI	[1] -	STA		LFVLSGHI	[8] -	TGRLGLKLSVGD	77	
q1	123180777	10	- [2] -	QDINIIVAVIDTEW	VWK	[10] -	A		PTGI	[1] -	SNG		QFLICTGA	[7] -	TADLETTAYPGD	73	
q1	53717990	9	- [2] -	QIKRVLVVIDIAY	IRS	[10] -	Q		PTGI	[1] -	SDE		QILLCTGS	[8] -	TGDLKFRANVGD	73	
q1	254245506	27	- [2] -	QQIDILAVIDTTEY	IKL	[10] -	L		PIAV	[1] -	SRA		VRLLYTGA	[8] -	PADPVLTLVPGD	81	
q1	170734880	2	- [2] -	VRCBALAVDAVT	LLS	[10] -	A		PIVI	[1] -	GRS		IVVLSPGD	[7] -	DSPLFAGLSPGD	85	
q1	53748592	18	- [2] -	LIINWVDTNVDV	ILA	[10] -	M		PTAI	[1] -	SAY		IVVLSDDP	[8] -	PGKILSNHVED	82	
q1	134279425	20	- [2] -	SRVLLVVIDSDY	VWK	[10] -	T		PIPV	[1] -	SRA		LFVICAGS	[8] -	SSEAICTAAVGD	82	
q1	170702239	10	- [2] -	QKITLLAVINDAK	[1] -	IKK	[10] -	R		PVQI	[1] -	SSE		QILLCDSP	[8] -	AKKIKFYAKCFD	75
q1	254224079	20	- [2] -	QIVWVDFLVDTIAY	IYA	[11] -	K		PKPI	[1] -	SNS		IVWACSFV	[7] -	TADLSFVWQGS	84	

InterPro protein sequence analysis & classification

InterPro is an integrated database of predictive protein signatures used for the classification and automatic annotation of proteins and genomes. InterPro classifies sequences at superfamily, family and subfamily levels, predicting the occurrence of functional domains, repeats and important sites. InterPro adds in-depth annotation, including GO terms, to the protein signatures.

European Bioinformatics Institute - <http://www.ebi.ac.uk/>

The screenshot displays the InterProScan Results page. At the top, there is a navigation bar with links for Research, Training, Industry, About Us, Help, Site Index, and RSS. Below this, the breadcrumb trail reads: EBI > Tools > Protein Functional Analysis > InterProScan Sequence Search. The main heading is "InterProScan Results", with tabs for Summary Table, Tool Output, Visual Output (selected), Submission Details, and Submit Another Job. A "Download in SVG format" button is visible. The main content area shows the following information:

- InterProScan (version: 4.8)**
- Sequence: Sequence_1
- Length: 288
- CRC64: 3FAE4C40C2498B64
- Launched Wed, May 16, 2012 at 17:31:03
- Finished Wed, May 16, 2012 at 17:35:39

The "InterPro Match" section shows a query sequence of length 288. The results are as follows:




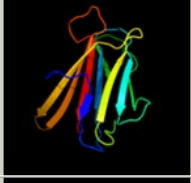


InterPro Match	Description	Signature
IPR010907	Calcium-mediated lectin	no description, PA-IIL, Calcium-mediated lectin
G3DSA:2.60.120.400		
PF07472		
SSF82026		
IPR021087	Uncharacterised protein family PixA/AidA	PixA
PF12306		

At the bottom, there is a legend for the database signatures used: PRODOM, HAMAP, PRINTS, PROSITE, PIR, SUPERFAMILY, PFAM, SIGNALP, SMART, TMHMM, TIGRFAMs, PANTHER, PROFILE, and GENE3D. The footer states: © European Bioinformatics Institute 2006-2012. EBI is an Outstation of the European Molecular Biology Laboratory.

Proč potřebujeme predikci domén

- Prohledávání sekvenčních databází bez predikce domén může být neúspěšné
- Automatická predikce struktury se zaměří jen na nejlépe „definovanou“ část
-

Phyre – whole protein http://www.sbg.bio.ic.ac.uk/phyre2/phyre2_output/a132b051273537c4/summary.html

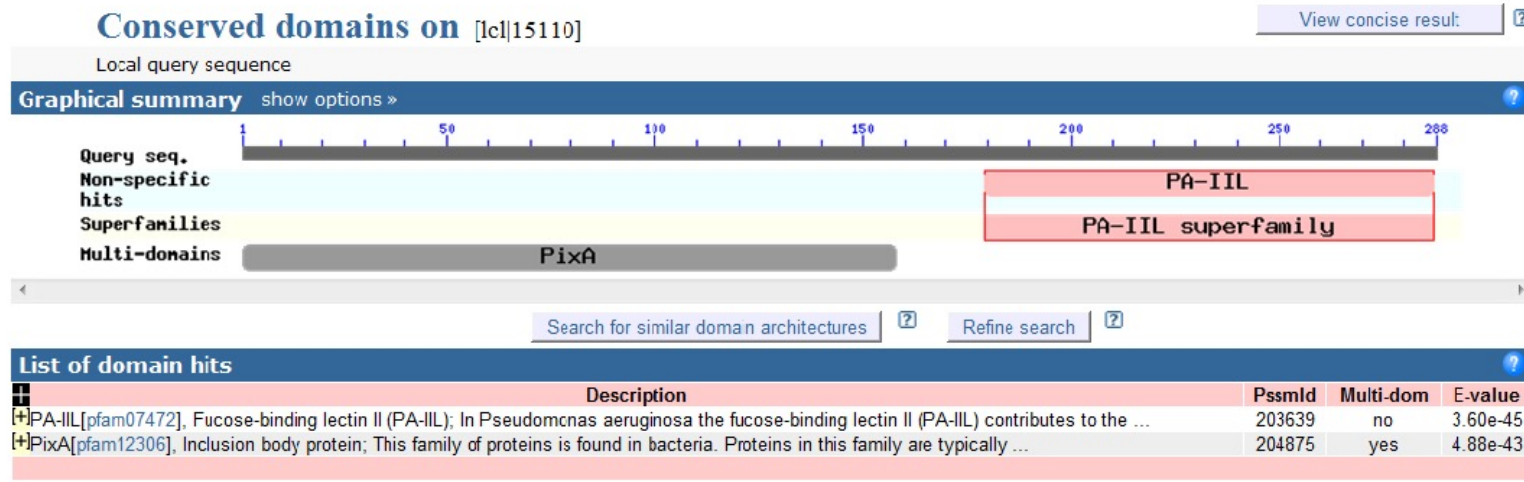
#	Template	Alignment Coverage	3D Model	Confidence	% I.d.	Template Information
1	c2vrvC <input type="radio"/> <input type="checkbox"/>	 <input type="button" value="Alignment"/>		100.0	60	PDB header: sugar-binding protein Chain: C; PDB Molecule: bcla; PDBTitle: crystal structure of bcla lectin from burkholderia2 cenocepacia in complex with alpha-methyl-mannoside at 1.73 angstrom resolution
2	c2xr4A <input type="radio"/> <input type="checkbox"/>	 <input type="button" value="Alignment"/>		100.0	43	PDB header: sugar binding protein Chain: A; PDB Molecule: lectin; PDBTitle: c-terminal domain of bc2l-c lectin from burkholderia cenocepacia
3	d2chha1 <input type="radio"/> <input type="checkbox"/>	 <input type="button" value="Alignment"/>		100.0	37	Fold: Calcium-mediated lectin Superfamily: Calcium-mediated lectin Family: Calcium-mediated lectin

Pokrocila bio

NCBI – Blast

Prohledávání databází známých aminokyselinových sekvencí

➤ celý protein




Phyre – C-term http://www.sbg.bio.ic.ac.uk/phyre2/phyre2_output/e332b1ecabb8d0a6/summary.html


#	Template	Alignment Coverage	3D Model	Confidence	% i.d.	Template Information
1	c2xr4A 	 <input type="button" value="Alignment"/>		100.0	44	PDB header: sugar binding protein Chain: A; PDB Molecule: lectin; PDBTitle: c-terminal domain of bc2l-c lectin from burkholderia cenocepacia
2	c2vnrC 	 <input type="button" value="Alignment"/>		100.0	62	PDB header: sugar-binding protein Chain: C; PDB Molecule: bcla; PDBTitle: crystal structure of bcla lectin from burkholderia2 cenocepacia in complex with alpha-methyl-mannoside at 1.73 angstrom resolution
3	d1uzva 	 <input type="button" value="Alignment"/>		100.0	30	Fold: Calcium-mediated lectin Superfamily: Calcium-mediated lectin Family: Calcium-mediated lectin

Phyre – n-term http://www.sbg.bio.ic.ac.uk/phyre2/phyre2_output/e332b1ecabb8d0a6/summary.html

#	Template	Alignment Coverage	3D Model	Confidence	% i.d.	Template Information
1	c1sddB 	 Alignment		83.7	9	PDB header: blood clotting Chain: B; PDB Molecule: coagulation factor v; PDBTitle: crystal structure of bovine factor vai
2	c3cdzB 	 Alignment		76.1	6	PDB header: blood clotting Chain: B; PDB Molecule: coagulation factor viii light chain; PDBTitle: crystal structure of human factor viii
3	d1kbva2 	 Alignment		68.0	13	Fold: Cupredoxin-like Superfamily: Cupredoxins Family: Multidomain cupredoxins


Swissprot – whole protein


 **BIOZENTRUM**
Universität Basel
The Center for Molecular Life Sciences


 **SWISS-MODEL Workspace**
Modelling Tools Repository Documentation


[myWorkspace] [login]

Workunit: P000007 - Overview


1  288


Print/Save this page as 


Model Summary 




Model information:
Modelled residue range: 169 to 288
Based on template: [2vnnD]* (1.7 Å)
Sequence Identity [%]: 56.35
Evalue: 0.00e-1

Quality information: [details] 
QMEAN Z-Score: -0.71

Quaternary structure information: [details] 
Template (2vnn): DIMER
Model built: SINGLE CHAIN

Ligand information: [details] 
Ligands in the template: CA: 3, MMA: 1, SO4: 1.
Ligands in the model: CA: 2

logs: [Templates]* [Alignment]* [Modelling]*
display model: as [pdb]* - as [DeepView project]* - in [AstexViewer]*
download model: as [pdb]* - as [Deepview project]* - as [text]*

Global Model Quality Estimation  [+/-]

http://swissmodel.expasy.org/workspace/index.phpuserid=michaw@chemi.muni.cz&key=0f449e99bc0176edfa75fba19b2d96e4&func=workspace_modelling&prjid=P000007



You don't have to be a scientist to do science.

By simply running a free program, you can help advance research in medicine, clean energy, and materials science.

Join Rosetta@home



HHMI
HOWARD HUGHES MEDICAL INSTITUTE



UNIVERSITY OF
WASHINGTON



Rosetta@home needs your help to determine the 3-dimensional shapes of proteins in research that may ultimately lead to finding cures for some major human diseases. By running the Rosetta program on your computer while you don't need it you will help us speed up and extend our research in ways we couldn't possibly attempt without your help. You will also be helping our efforts at designing new proteins to fight diseases such as HIV, Malaria, Cancer, and Alzheimer's. Please [join us](#) in our efforts!



The Science Behind Foldit

Foldit is a revolutionary crowdsourcing computer game enabling *you* to contribute to important scientific research. This page describes the science behind Foldit and how your playing can help.

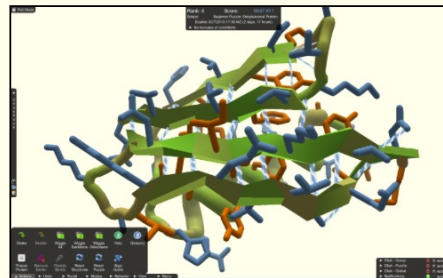
Page Contents:

- [What is protein folding?](#)
- [Why is this game important?](#)
- [Foldit Scientific Publications](#)
- [News Articles about Foldit](#)
- [News Articles about Rosetta](#)
- [Rosetta@Home Screensaver](#)
- [Community Rules](#)
- [Let's Foldit Podcast](#)
- [Instructions for Educators](#)
- [Terms of Service and Consent](#)
- [Credits](#)

<http://fold.it/portal/>

What is protein folding?

What is a protein? Proteins are the workhorses in every cell of every living thing. Your body is made up of trillions of cells, of all different kinds: muscle cells, brain cells, blood cells, and more. Inside those cells, proteins are allowing your body to do what it does: break down food to power your muscles, send signals through your brain that control the body, and transport nutrients through your blood. Proteins come in thousands of different varieties, but they all have a lot in common. For instance, they're made of the same



Folded up Streptococcal Protein Puzzle
(+) [Enlarge This Image](#)

GET STARTED: DOWNLOAD



Windows
(XP/Vista/7/8)



OSX
(10.7 or later)



Linux
(64-bit)

[Are you new to Foldit? Click here.](#)

[Are you a student? Click here.](#)

[Are you an educator? Click here.](#)

SEARCH

Only search fold.it

RECOMMEND FOLDIT

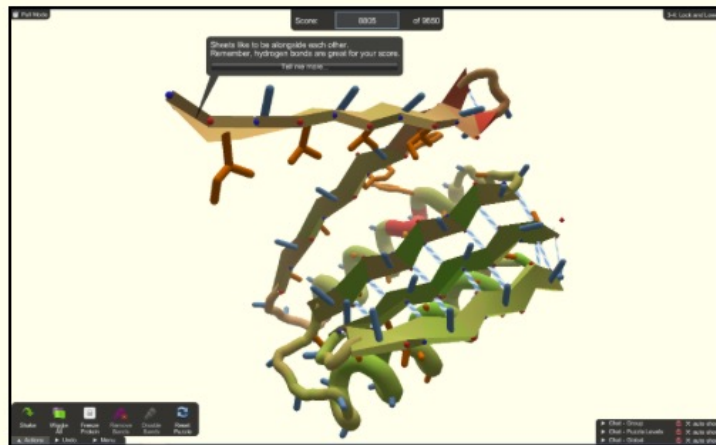
USER LOGIN

Username: *

Password: *

- [Create new account](#)
- [Request new password](#)

Just a game?



This is an example of a puzzle that a human can see the obvious answer to - fix the sheet that is sticking out!

[\(+\)](#) **Enlarge This Image**

proteins?

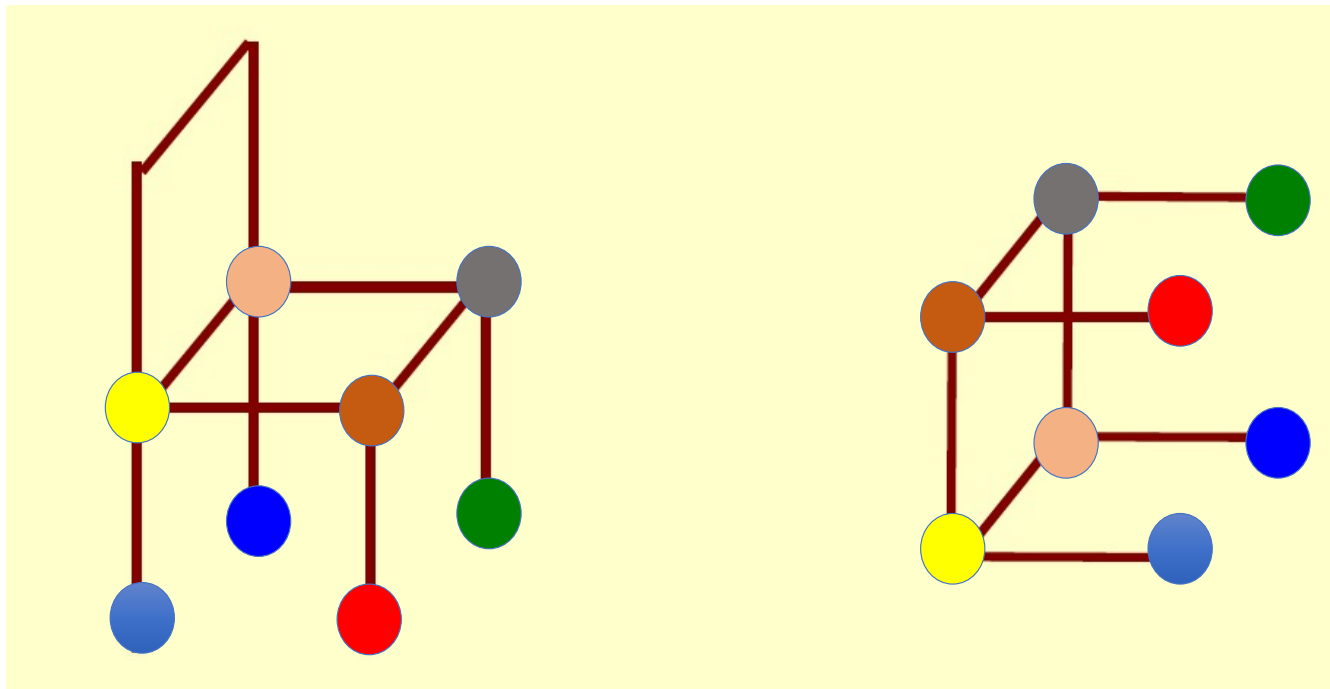
We're collecting data to find out if humans' pattern-recognition and puzzle-solving abilities make them more efficient than existing computer programs at pattern-folding tasks. If this turns out to be true, we can then teach human strategies to computers and fold proteins faster than ever!

What other good stuff am I contributing to by playing?

Proteins are found in all living things, including plants. Certain types of plants are grown and converted to biofuel, but the conversion process is not as fast and efficient as it could be. A critical step in turning plants into fuel is breaking down the plant material, which is currently done by microbial enzymes (proteins) called "cellulases". Perhaps we can find new proteins to do it better.

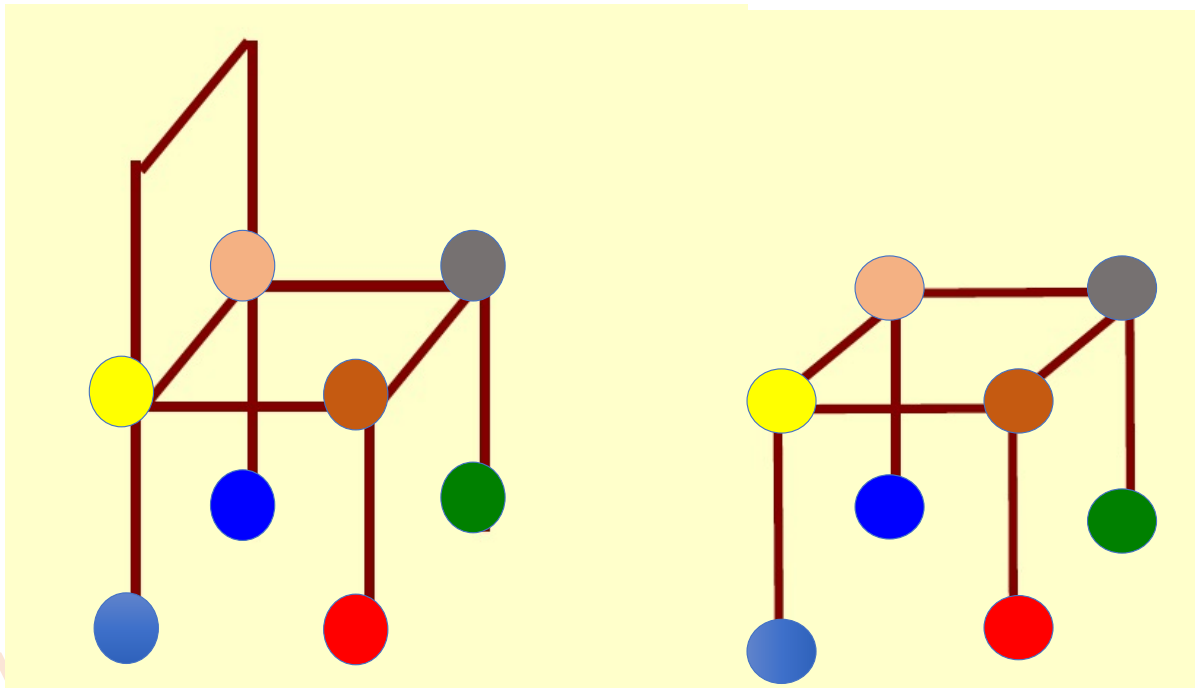
Can humans really help computers fold

Structure Superposition



The key is finding corresponding points between the two structures

Structure Superposition



The key is finding corresponding points between the two structures

Algorithms for Structure Superposition

Distance based methods:

DALI (Holm & Sander): Aligning scalar distance plots

SSAP (Orengo & Taylor): Dynamic programming using intra-molecular vector distances

MINAREA (Falicov and Cohen): Minimizing soap-bubble surface area

CE (Shindyalov & Bourne)

Vector based methods:

VAST (Bryant): Graph theory based secondary structure alignment

3D Search (Singh and Brutlag) & 3D Lookup (Holm and Sander): Fast secondary structure index lookup

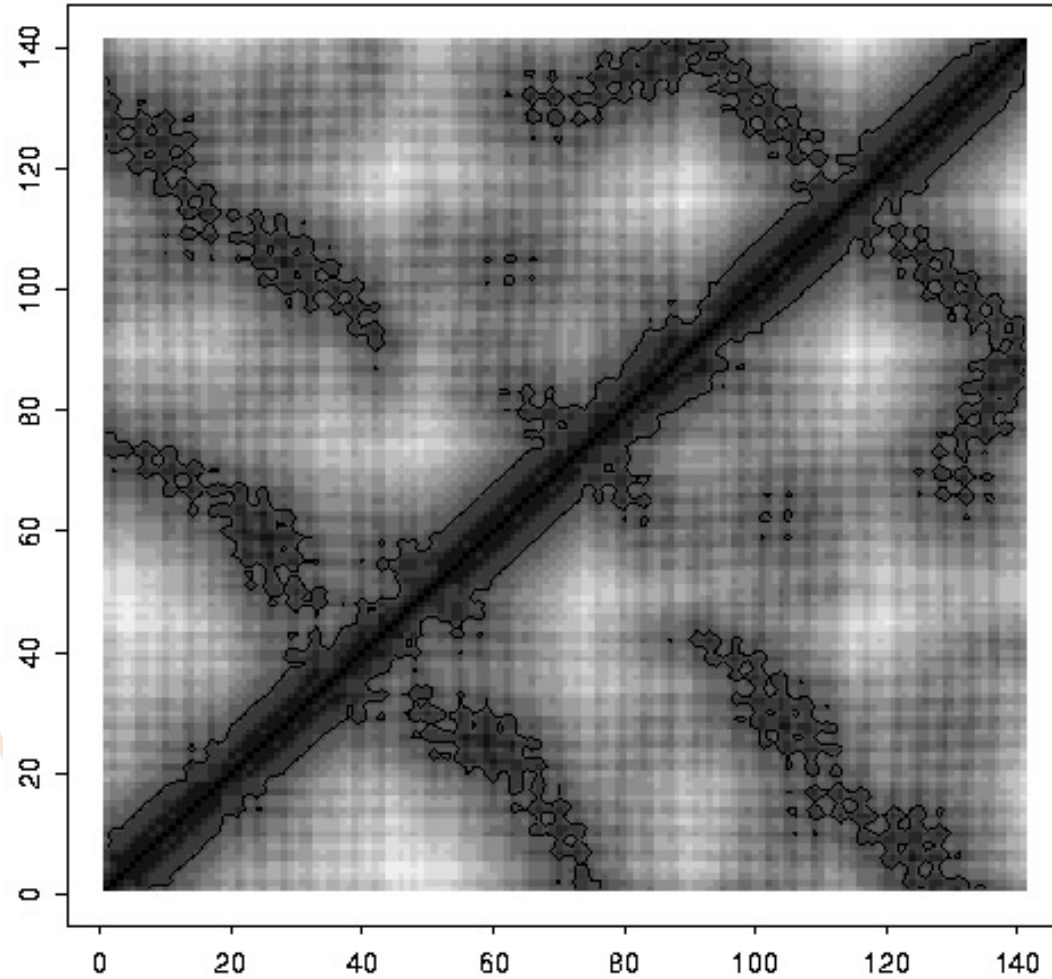
Both

LOCK (Singh & Brutlag) LOCK2 (Ebert & Brutlag): Hierarchically uses “Adaptive”

FATCAT(Flexible structure **A**lignment**T** by **C**haining **A**ligned fragment pairs allowing **T**wists, Ye & Godzik) – not further maintained?

<http://fatcat.godziklab.org/fatcat/>

DALI



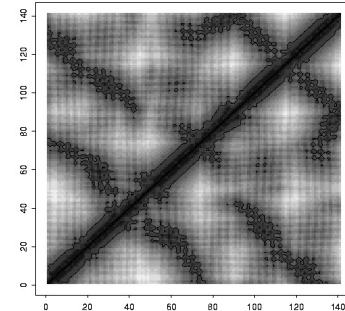
An intra-molecular distance plot for myoglobin

DALI

Based on aligning 2-D intra-molecular distance matrices

Computes the best subset of corresponding residues from the two proteins such that the similarity between the 2-D distance matrices is maximized

Searches through all possible alignments of residues using Monte-Carlo and Branch-and-Bound algorithms

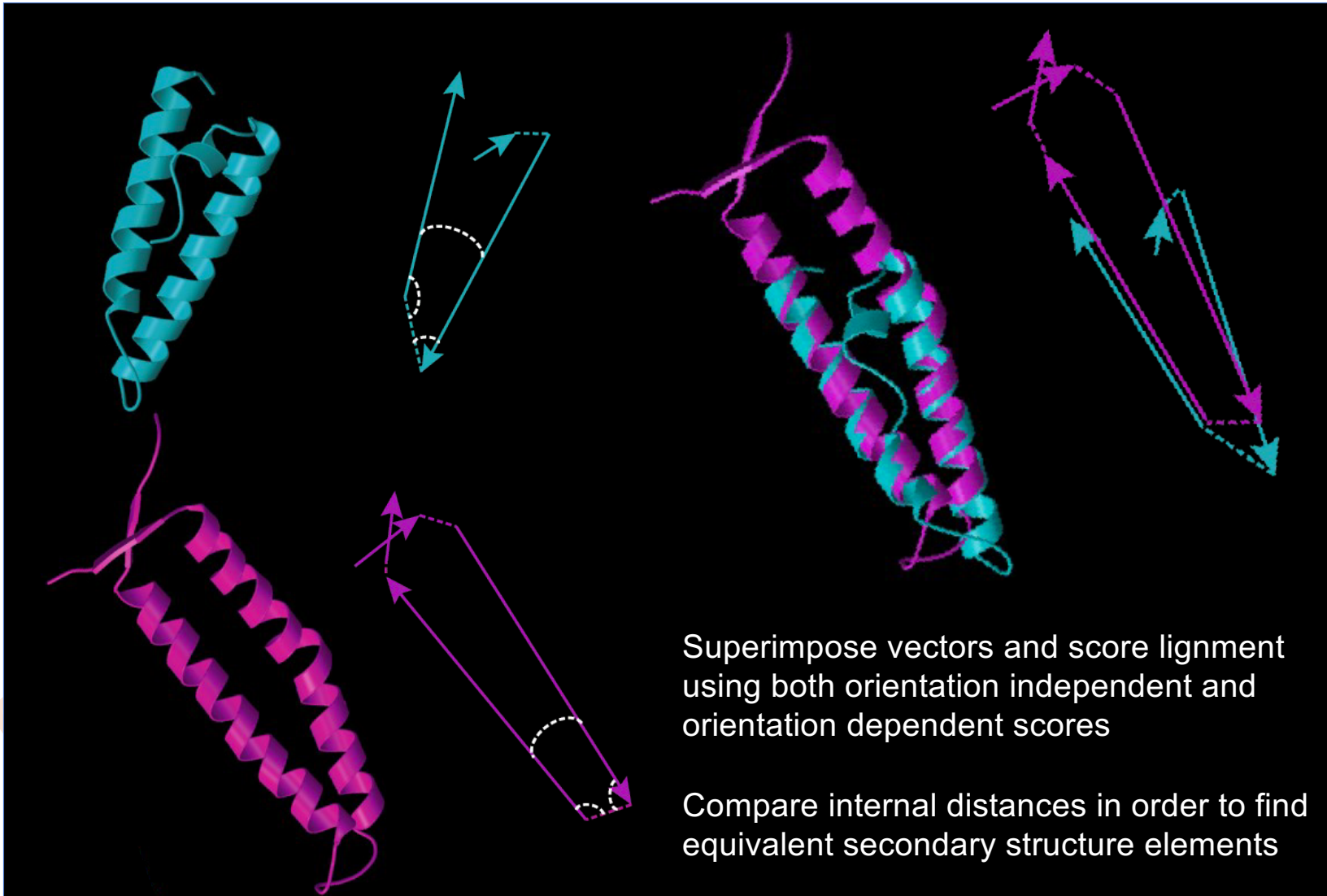


VAST – Vector Alignment Search Tool

Identifying similar structures by **purely geometric criteria** (and to identify distant homologs that cannot be recognized by sequence comparison). Find similarly shaped individual protein molecules or 3D domains (VAST+: similarly shaped macromolecular complexes)

- Aligns only secondary structure elements (SSE)
- Represents each SSE as a vector
- Finds all possible pairs of vectors from the two structures that are similar
- Uses a graph theory algorithm to find maximal subset of similar vector pairs
- Overall alignment score is based on the number of similar pairs of vectors between the two structures

LOCK2



FoldMiner: Structure Similarity Search Based on LOCK2 Alignment

FoldMiner aligns query structure with all database structures using LOCK2

FoldMiner up weights secondary structure elements in query that are aligned more often

FoldMiner outperforms CE and VAST in searches for structure similarity

The best to test as first:

Distance based methods

DALI

<http://ekhidna2.biocenter.helsinki.fi/dali/>

Vector and distance based method

FoldMiner (LOCK2) – local installation needed

“Adaptive”

FATCAT

<http://fatcat.godziklab.org/fatcat/>

Závěrem

- Struktura je klíčová pro správnou funkci proteinu
- Predikovat na základě sekvence (1D) lze 2D, 3D i 4D strukturu
- Vždy je nutné **kriticky kontrolovat** výstupy programů
- Ideální je využít více predikčních programů s různou metodologií a porovnat výsledky