

Parte II

Metodi per la predizione della struttura di proteine

Perchè predire la struttura terziaria?

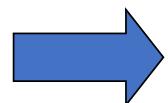
La cristallografia pone dei problemi

- E' necessario un campione di proteina molto puro
- Il campione di proteina deve essere in grado di formare cristalli (generalmente questo rappresenta il problema maggiore)
- Molte proteine non possono essere cristallizzate (i.e. proteine transmembrana)
- La determinazione di una struttura proteica ha in genere costi piuttosto elevati (~\$100K per struttura)

Perchè predire la struttura terziaria?

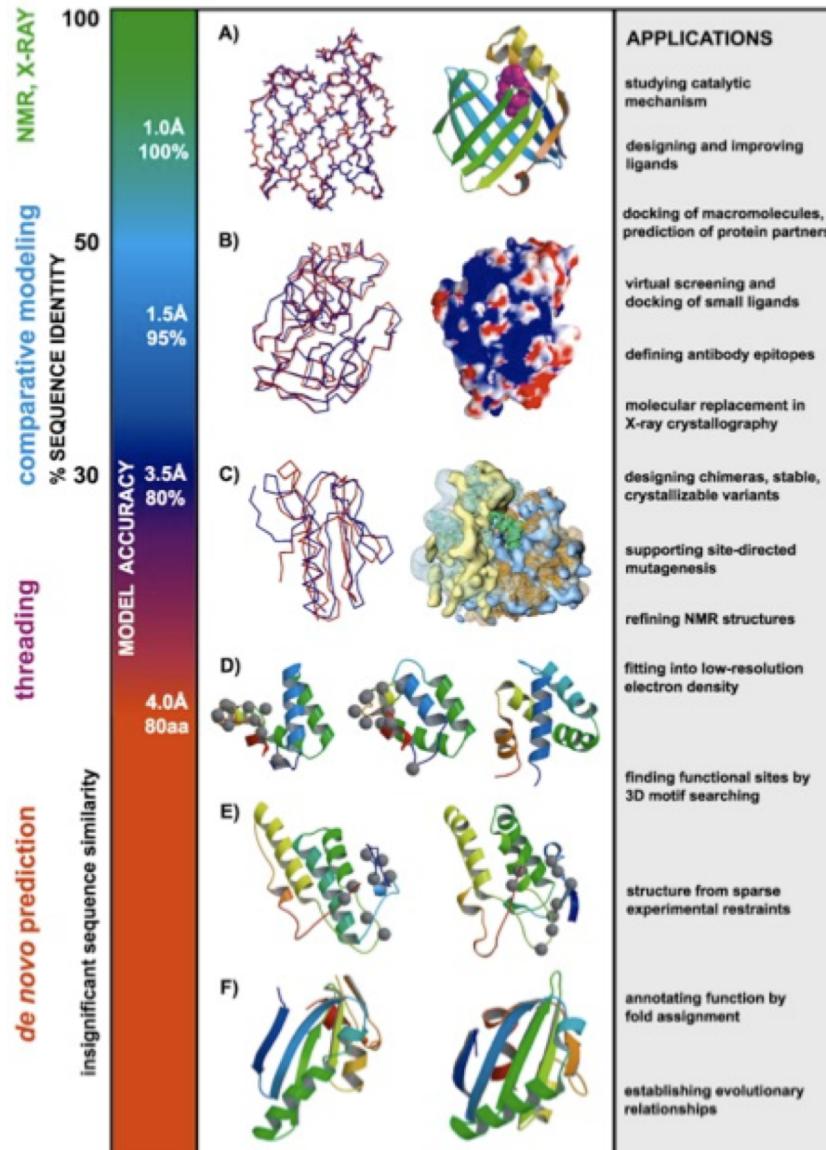
- **In cifre:**
 - 540.052 (SwissProt) + 33.995.348 (TrEMBL) May 2013
 - 90.424 (PDB), 23.887 uniche (<30% SeqId) May 2013
 - *La distanza tra il numero di sequenze e di strutture risolte sta aumentando.*
- **Metodi computazionali**
 - Veloci (minuti o ore), poco costosi (PC)
 - Soluzioni corrette ca. nel 60% dei casi.
 - Risoluzione più bassa, però spesso sufficiente per spiegare la funzione proteica
- Osservazione: La **sequenza** si evolve più rapidamente della **struttura** (*Chothia & Lesk, 1986*)
 - Numero limitato di fold

La sequenza di una proteina contiene tutta l'informazione necessaria perché essa possa ripiegarsi nella sua conformazione nativa



E' possibile modellare la struttura delle proteine a partire dalla loro sequenza?

Se sì, come si può utilizzare un modello?



Metodi predittivi

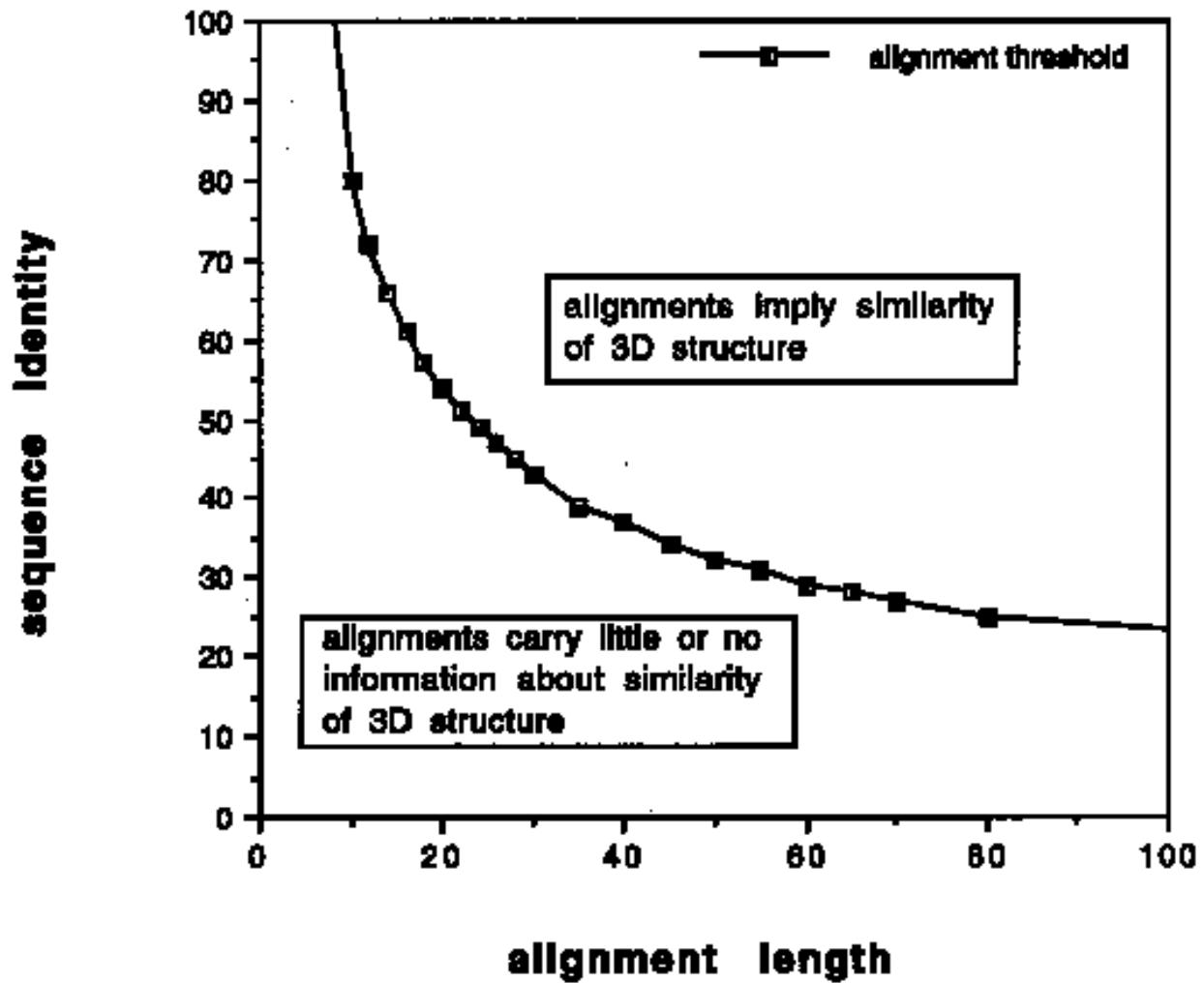
- » Comparative modeling
- » Threading/Fold recognition
- » Ab initio

> 30% similarità

0 – 30% similarità

nessun omologo

Threshold for structural homology



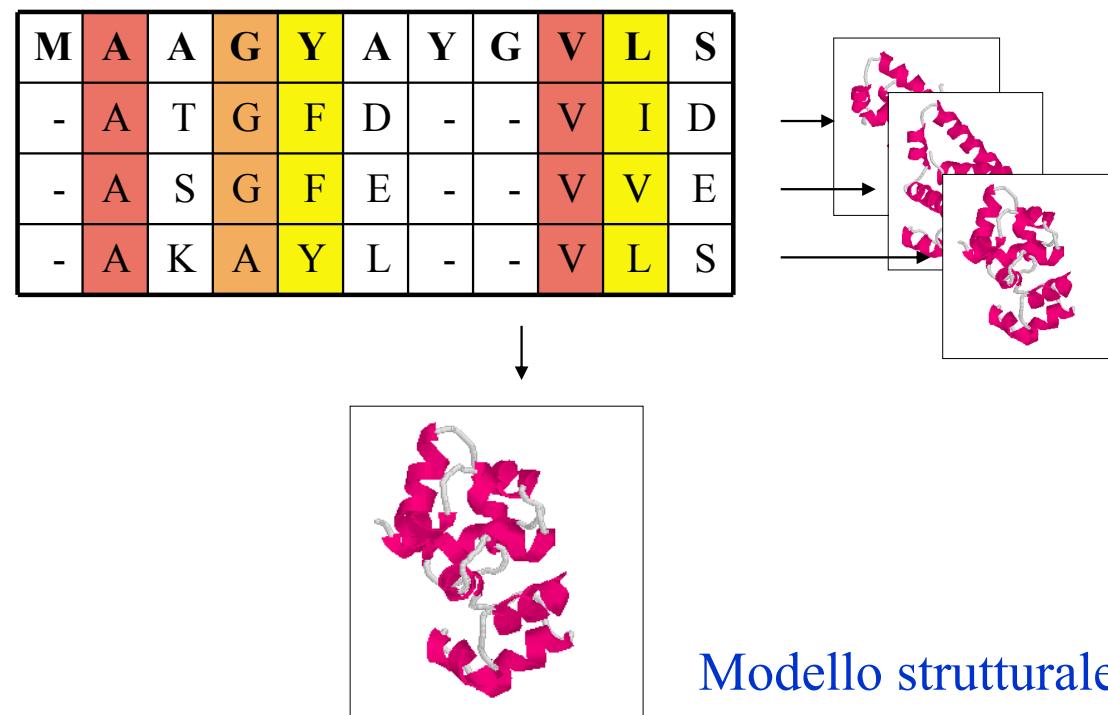
Qualità del modello comparativo

Identità di sequenza:

60-100%	Confrontabile con NMR media risoluzione Specificità di substrato
30-60%	Molecular replacement in cristallografia Partenza per site-directed mutagenesis
<30%	Gravi errori

Building by homology (Homology modelling)

Allineamento con proteine a struttura nota



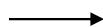
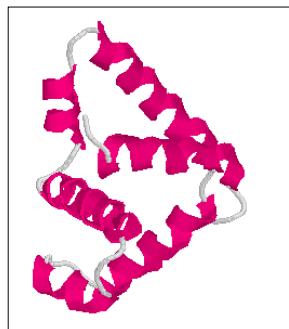
Fold recognition (Threading)

Sequenza:

M	A	A	G	Y	A	V	L	S
---	---	---	---	---	---	---	---	---

+

Motivi strutturali noti



Modello strutturale

Ab initio

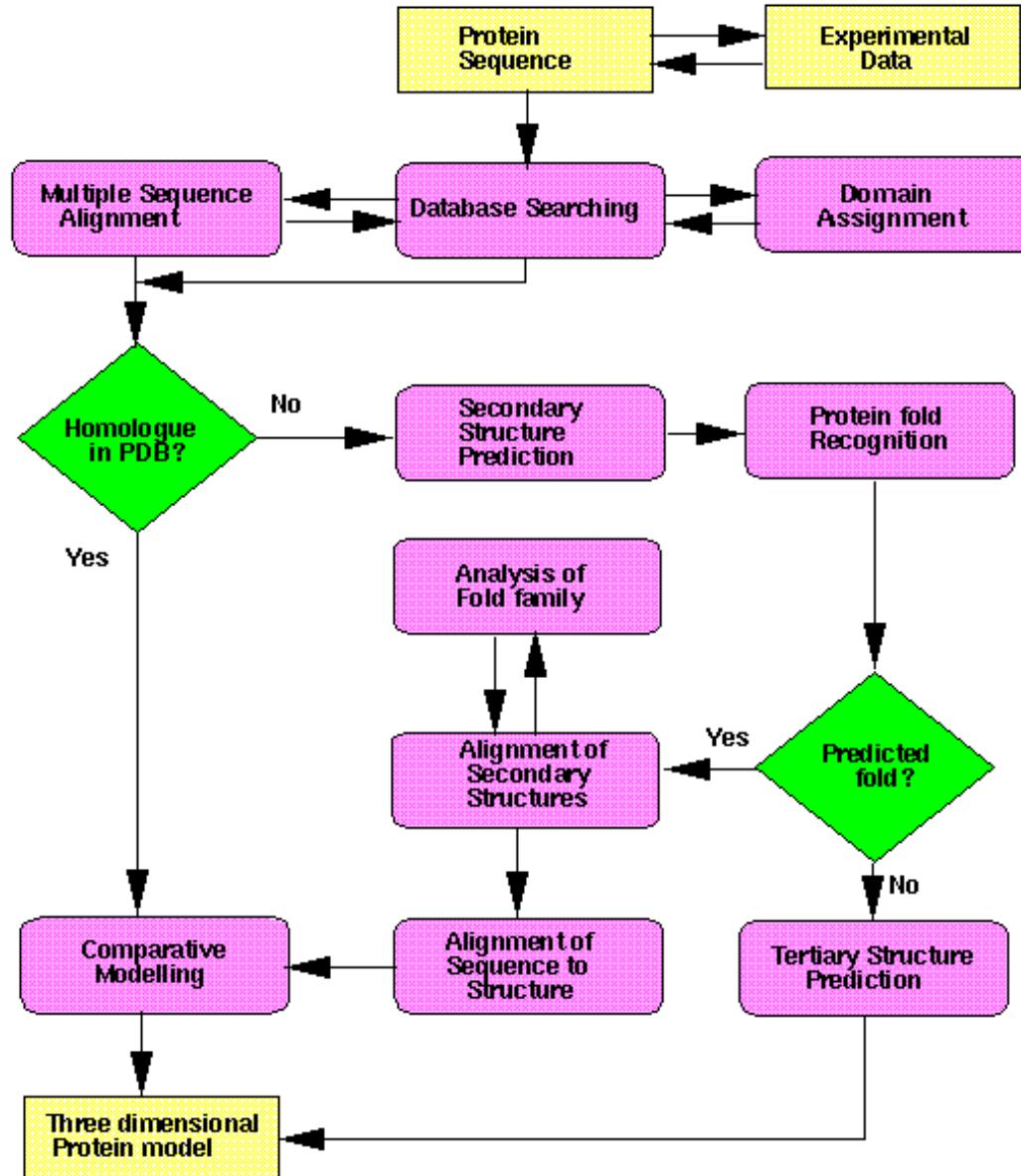
Sequenza

M	A	A	G	Y	A	V	L	S
---	---	---	---	---	---	---	---	---

 →

Modello strutturale

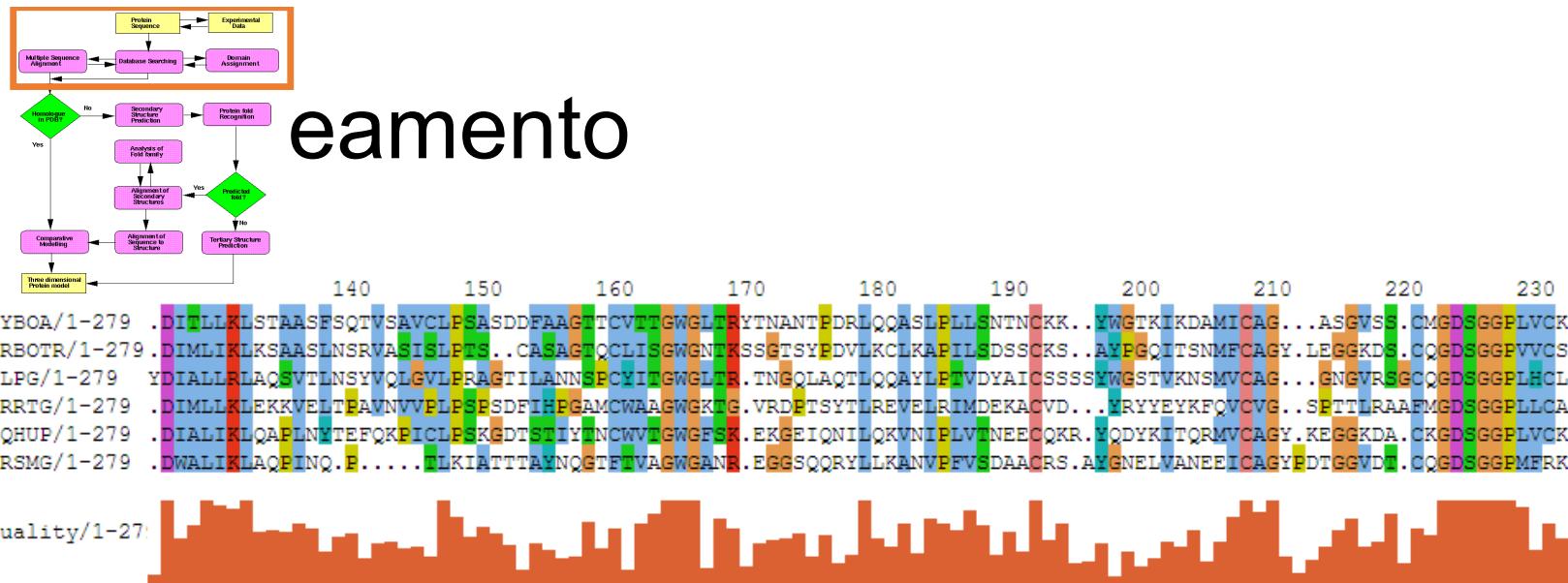
General Flowchart



Building by homology

Un numero grandissimo di polipeptidi si struttura in un numero finito (e relativamente piccolo) di folds

Almeno una proteina su due di quelle presenti nel database ha un omologo (identità > 30%) che quasi sempre ha lo stesso fold.



- Ricerca in database per trovare sequenze omologhe con struttura nota.
- Assegna le posizioni di residui equivalenti fra *target* e *template*. Determina *inserzioni* e *delezioni*.
- **L'allineamento determina la qualità del modello che si sta costruendo.**
- L'allineamento di sequenza non è sempre ottimale per costruire i modelli.
- Generalmente usato: *PSI-BLAST*
 - identifica sequenze omologhe anche molto remote utilizzando le *PSSM* (*position specific scoring matrix*).

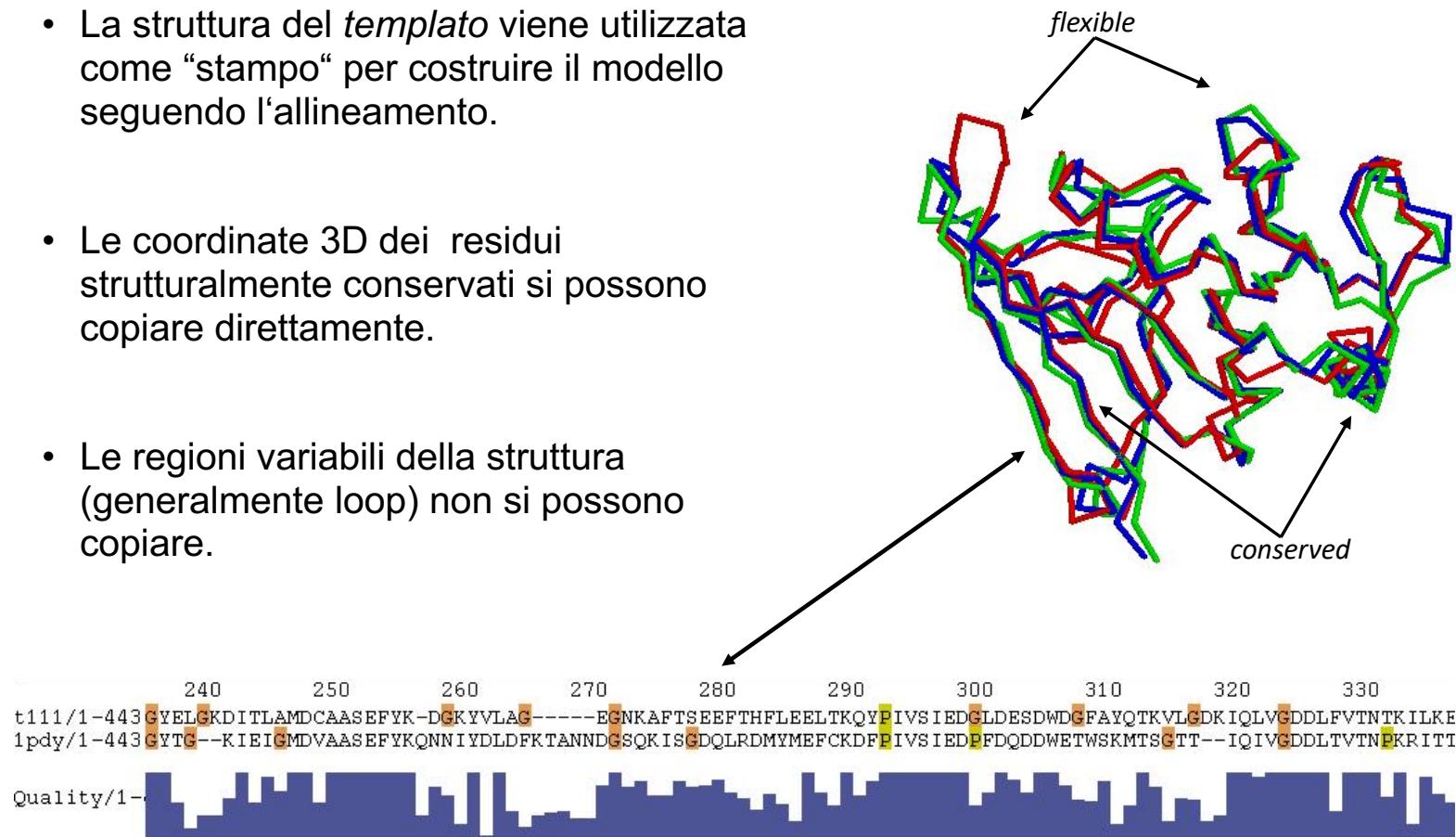
Costruire il modello comparativo

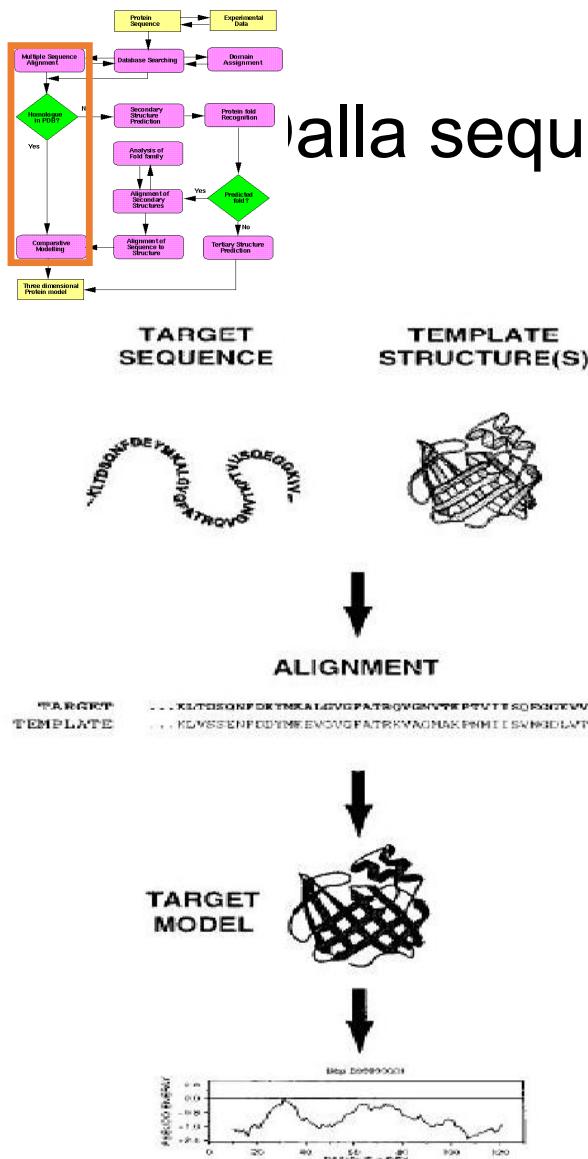
- 1) Cercare il massimo numero di omologhi che possiedano una entry nel PDB. Strumenti che utilizzano PSSM sono più sensibili. In questo caso vengono utilizzate sequenze senza struttura per costruire la PSSM.

- 2) Costruire un accurato allineamento multiplo tra la sequenza da modellare e tutte le entries che verranno utilizzate come templato.

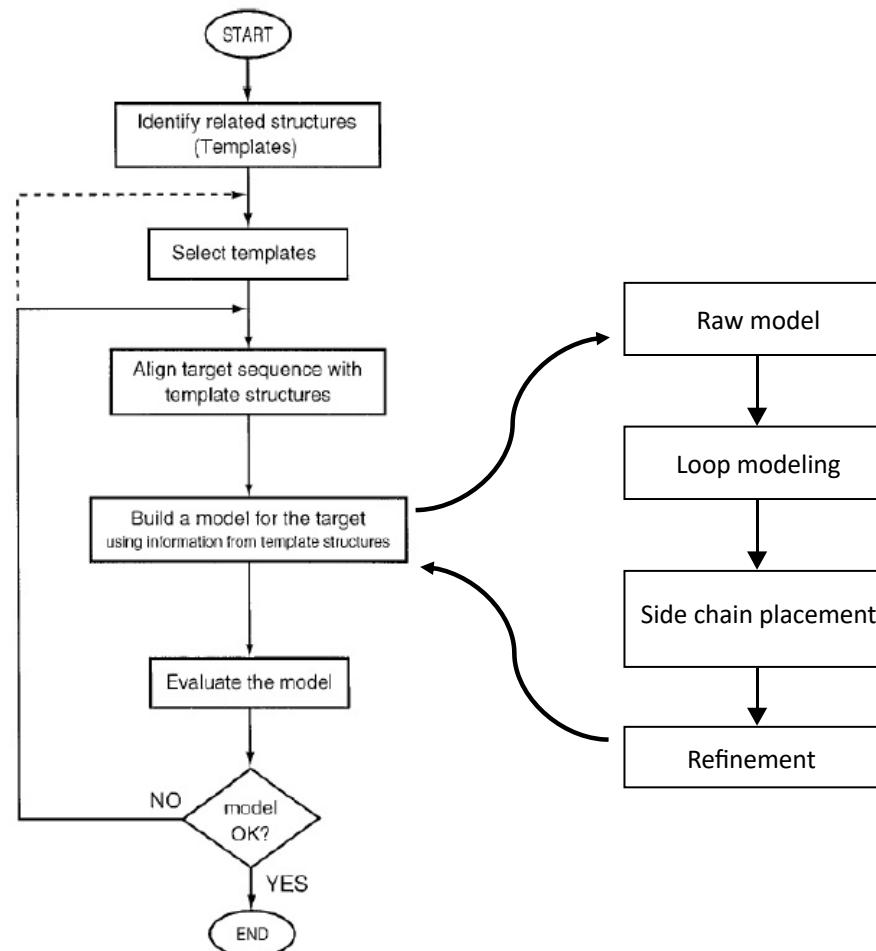
Costruzione del *pre*-modello

- La struttura del *template* viene utilizzata come “stampo” per costruire il modello seguendo l’allineamento.
- Le coordinate 3D dei residui strutturalmente conservati si possono copiare direttamente.
- Le regioni variabili della struttura (generalmente loop) non si possono copiare.





Proteina alla sequenza al modello



Costruire il modello stesso

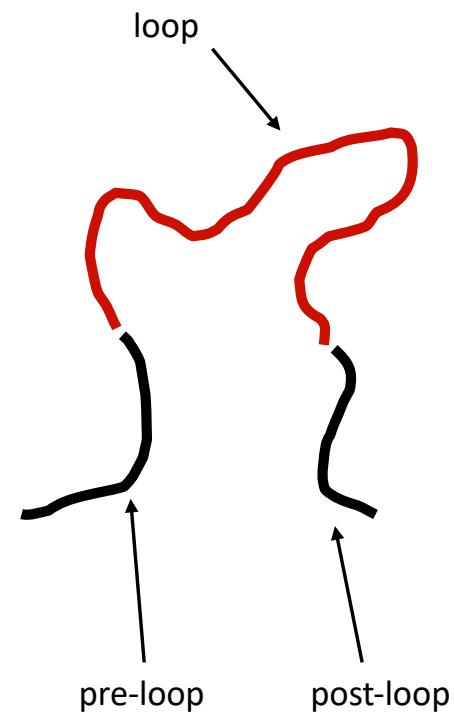
Determinare la struttura secondaria in base all'allineamento

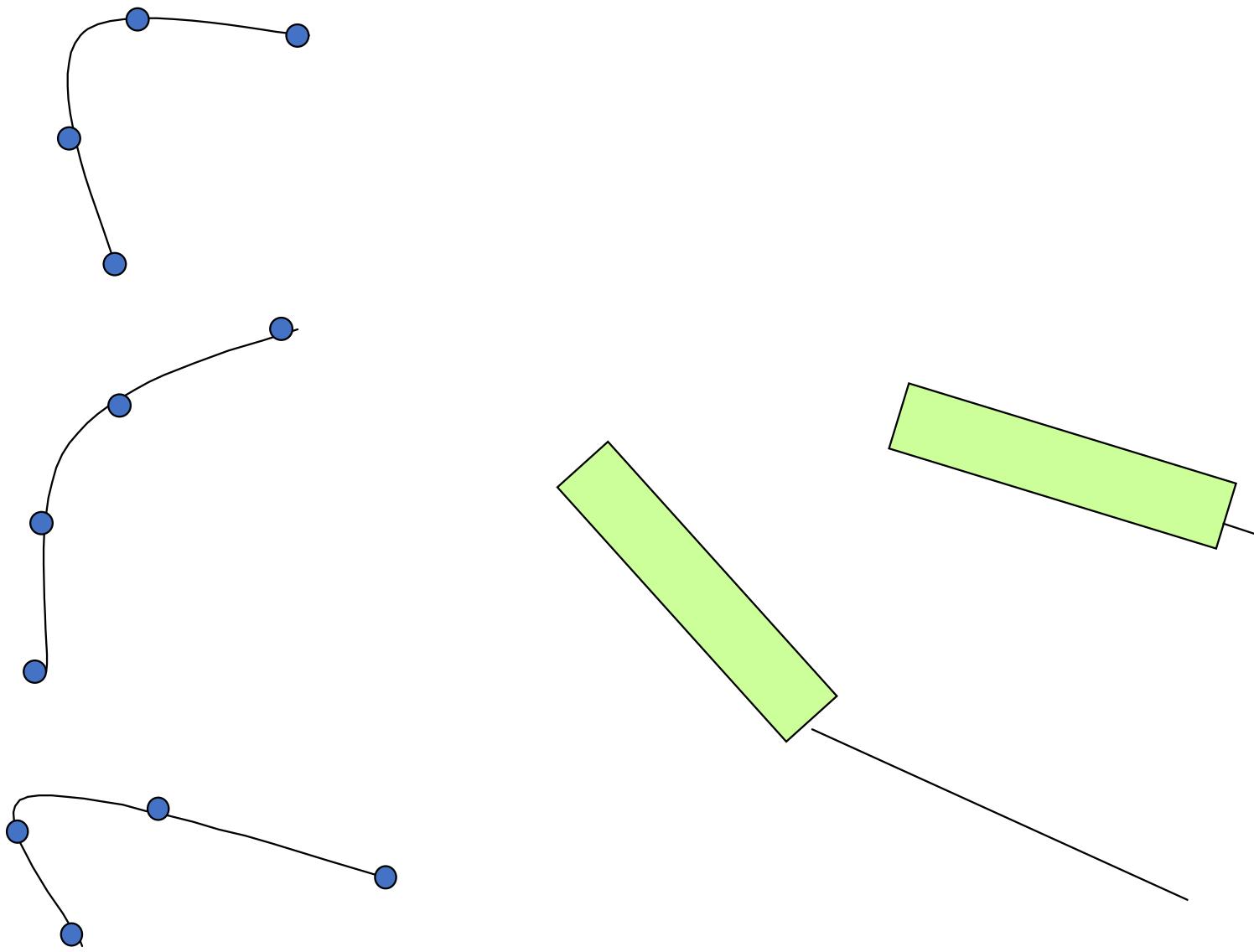
Costruire le regioni conservate. Per ciascuna regione possiamo prendere le coordinate del frammento con la maggior similarità di sequenza.

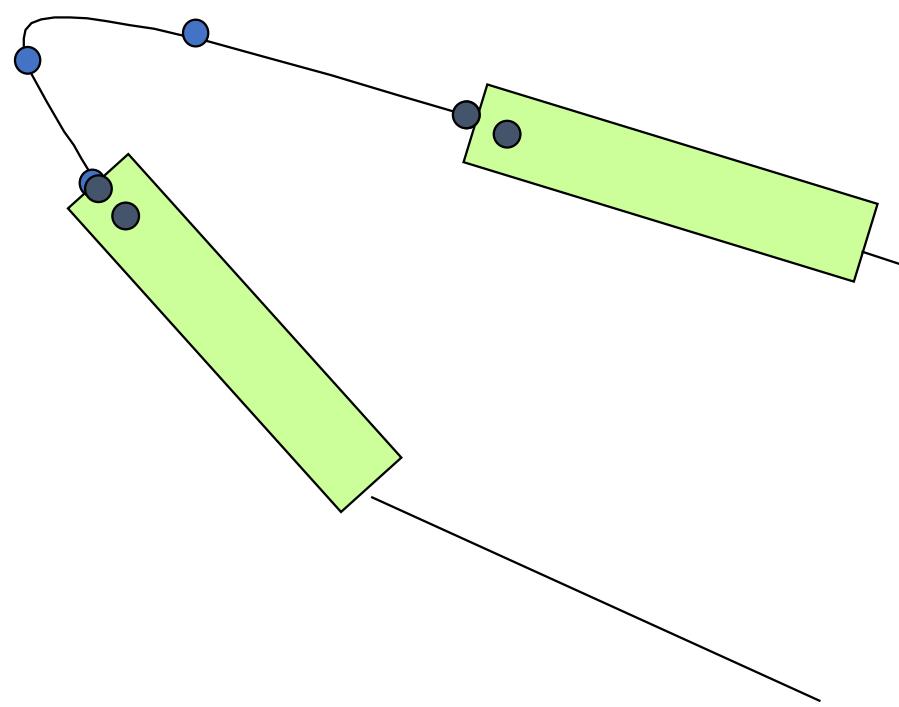
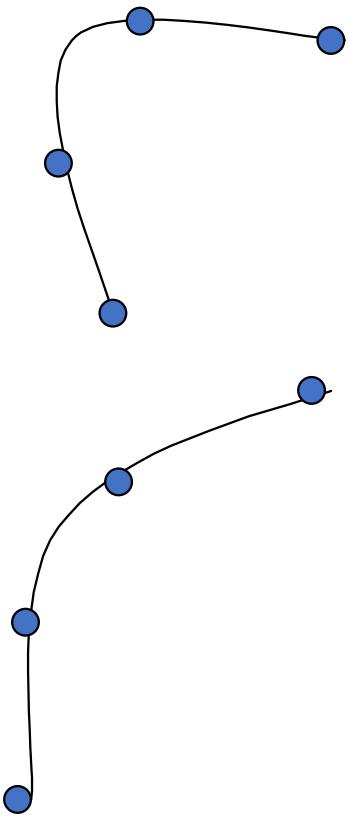
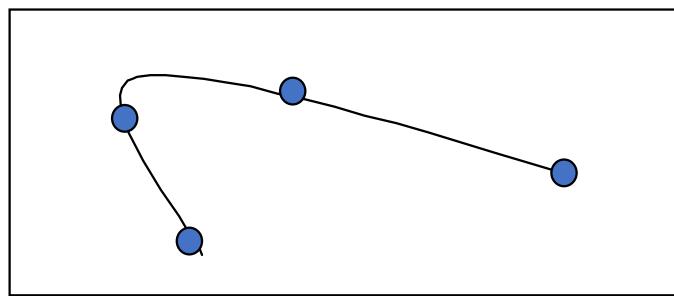
Costruire le regioni variabili, solitamente loops.

Loop modeling

- Al pre-modello possono mancare interi frammenti di catena principale
 - non conservati nella famiglia proteica
 - Inserzioni
 - Delezioni
- Descrizione del problema:
 - Si cerca un fold che colleghi il frammento N-terminale (pre-loop) con quello C-terminale (post-loop) tramite k residui
 - (f, y) sono gli unici parametri liberi
- Metodi di database
 - Estrarre frammenti di loop dal PDB. Scegli il frammento che rispetta meglio i vincoli geometrici.
- Metodi *ab initio*
 - Genera molti frammenti alternativi basati sui vincoli geometrici (angoli torsionali). Seleziona il frammento “migliore”.

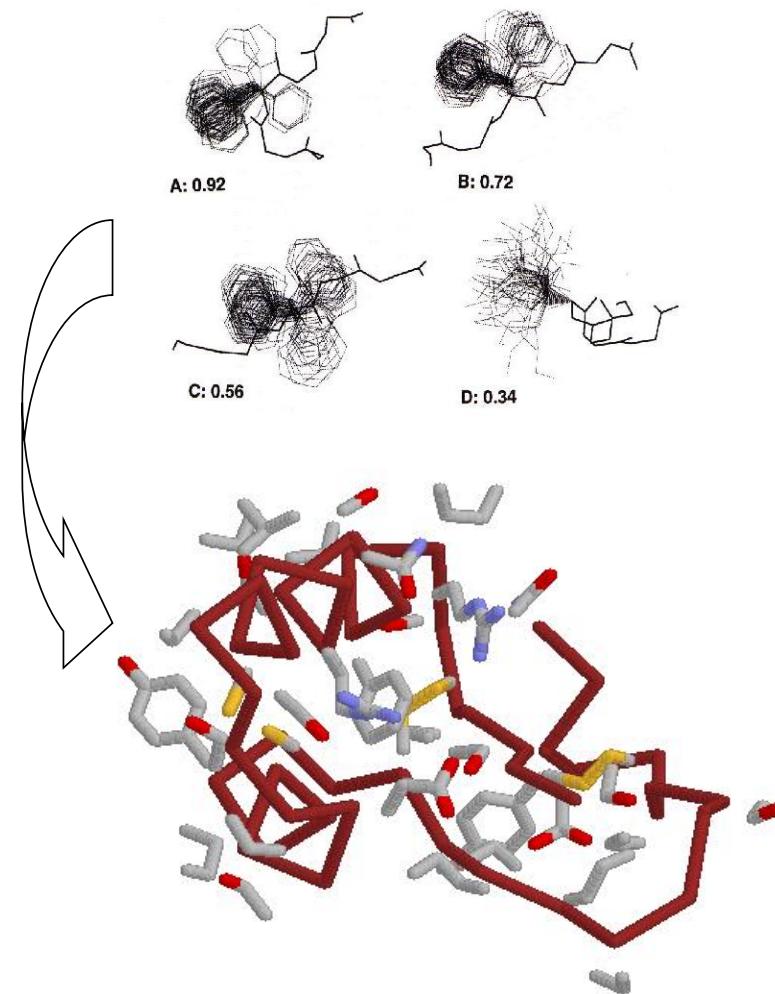






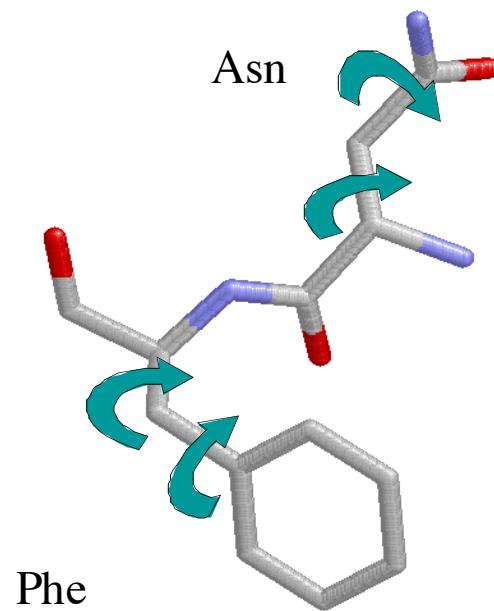
Catene laterali

- **Problema:** Applicando le coordinate del *templato* sulla sequenza del *target* cambiano tipo, dimensione e posizione delle catene laterali.
- Rotameri
 - 3 posizioni per angolo torsionale
 - Interdipendenza, effetto domino
- L'RMSD cambia relativamente poco, però possono cambiare le conformazioni di residui importanti (p.es. del sito attivo)
- Dove possibile è meglio mantenere le conformazioni delle catene laterali del *templato*.
- Esistono metodi standard per risolvere questo problema.



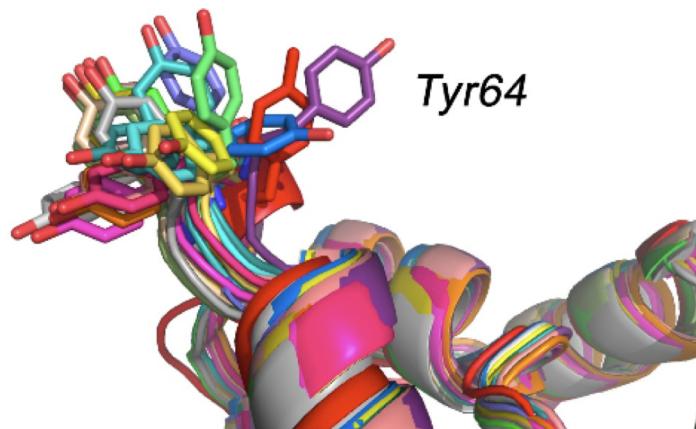
Conformation - a given set of dihedral angle which defines a structure.

Rotamer - energetically favourable conformation.



Side Chain Conformation Search

- ❑ With two rotatable backbone bonds per residue, it's very difficult to find the best conformation of a side chain
- ❑ In addition, side chains of many residues have one or more degree of freedom.
- ❑ Side chain conformational search in loop regions is much more difficult



Selection Of Good Rotamers

- ❑ Fortunately, statistical studies show side chain adopt only a small number of many possible conformations
- ❑ The correct rotamer of a particular residue is mainly determined by local environment
- ❑ Side chain generally adopt conformations where they are closely packed

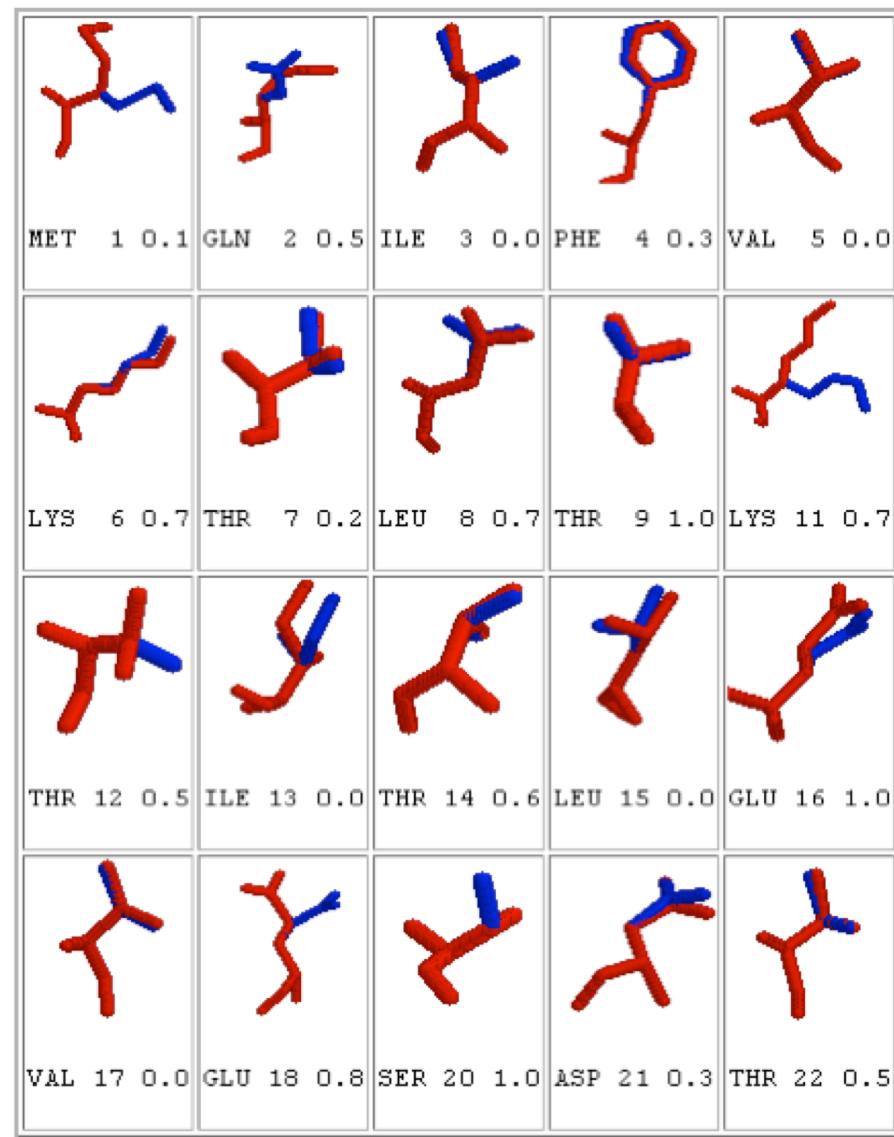
It has been observed that:

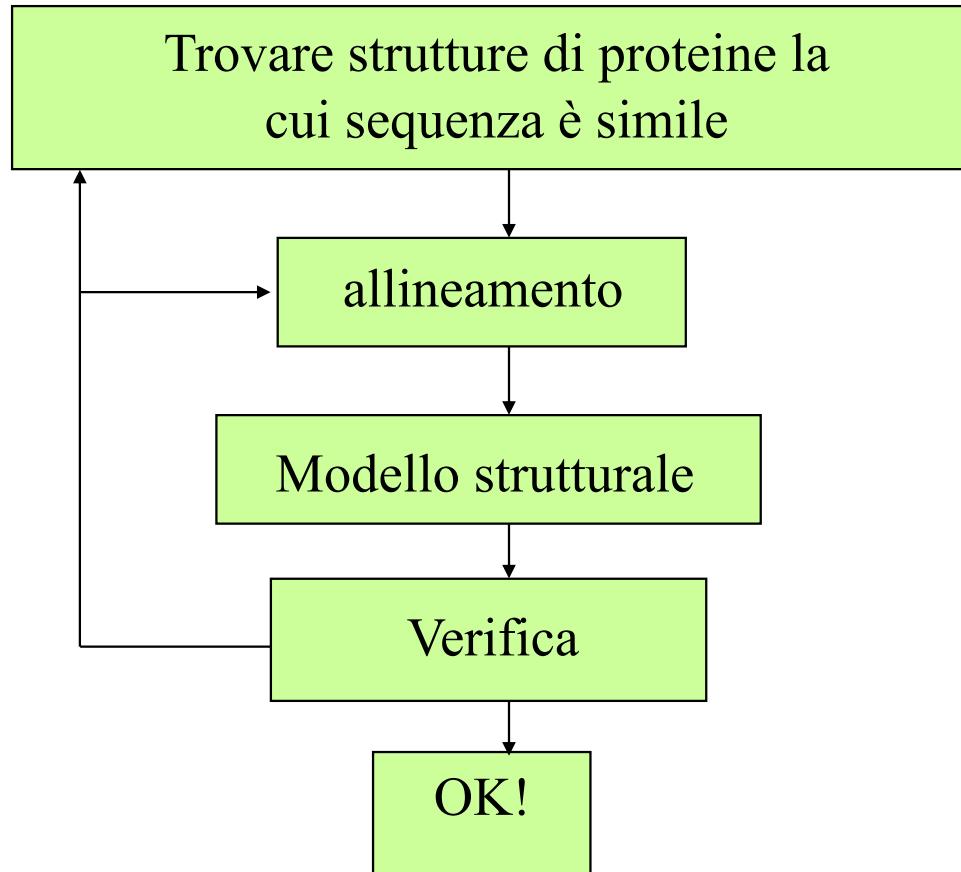
- In homologous proteins, corresponding residues virtually retain the same rotameric state ([Ponder and Richards 1987](#), [Benedetti et al. 1983](#))
- Certain rotamers are almost always associated with certain secondary structure([McGregor et al. 1987](#)).

Esempio di libreria di rotameri

SER	59.6	41.0
SER	-62.5	26.4
SER	179.6	32.6

TYR	63.6	90.5	21.0
TYR	68.5	-89.6	16.4
TYR	170.7	97.8	13.3
TYR	-175.0	-100.7	20.0
TYR	-60.1	96.6	10.0
TYR	-63.0	-101.6	19.3





Refinement

- Per ridurre tutti quei piccoli errori che si accumulano durante il processo di modelling si può ricorrere ai campi di forza (p.es. *CHARMM* o *AMBER*) per minimizzare l'energia del modello.
- Riducono le collisioni molecolari e rendono il modello “più bello”.
- Non modificano significativamente il modello e richiedono relativamente tanto tempo di calcolo.
- In caso di eccesso possono incrementare l'RMSD complessiva del modello.

- Every model contains errors.
- Two main reasons:
 - % sequence identity between reference and model
 - The number of errors in templates
- Hence it is essential to check the correctness of overall fold/structure, errors of localized regions and stereochemical parameters: bond lengths, angles, geometries

Evaluation of model accuracy

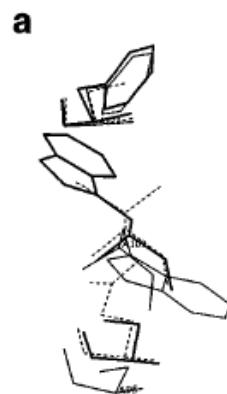


“... a model must be wrong, in some respects -- else it would be the thing itself. The trick is to see ... where it is right.”

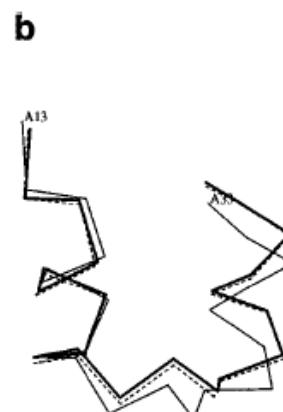
*Henry A. Bent
"Uses (and Abuses) of Models in Teaching Chemistry,"
J. Chem. Ed. **1984** 61, 774.*

SIV Model based on: 1BL3 (C) HIV-1 Integrase core domain
Experimental structure: 1C6V (C) SIV Integrase core domain
Seq. Identity: 61 %

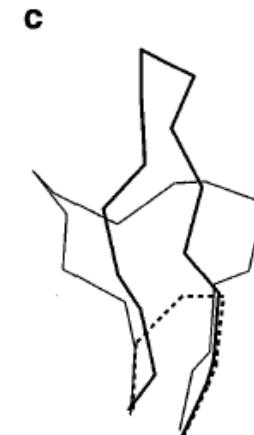
Errori tipici



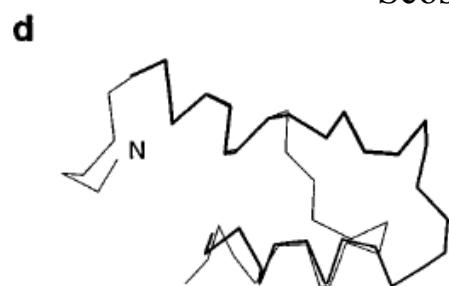
Catene laterali



Scostamento



Loops

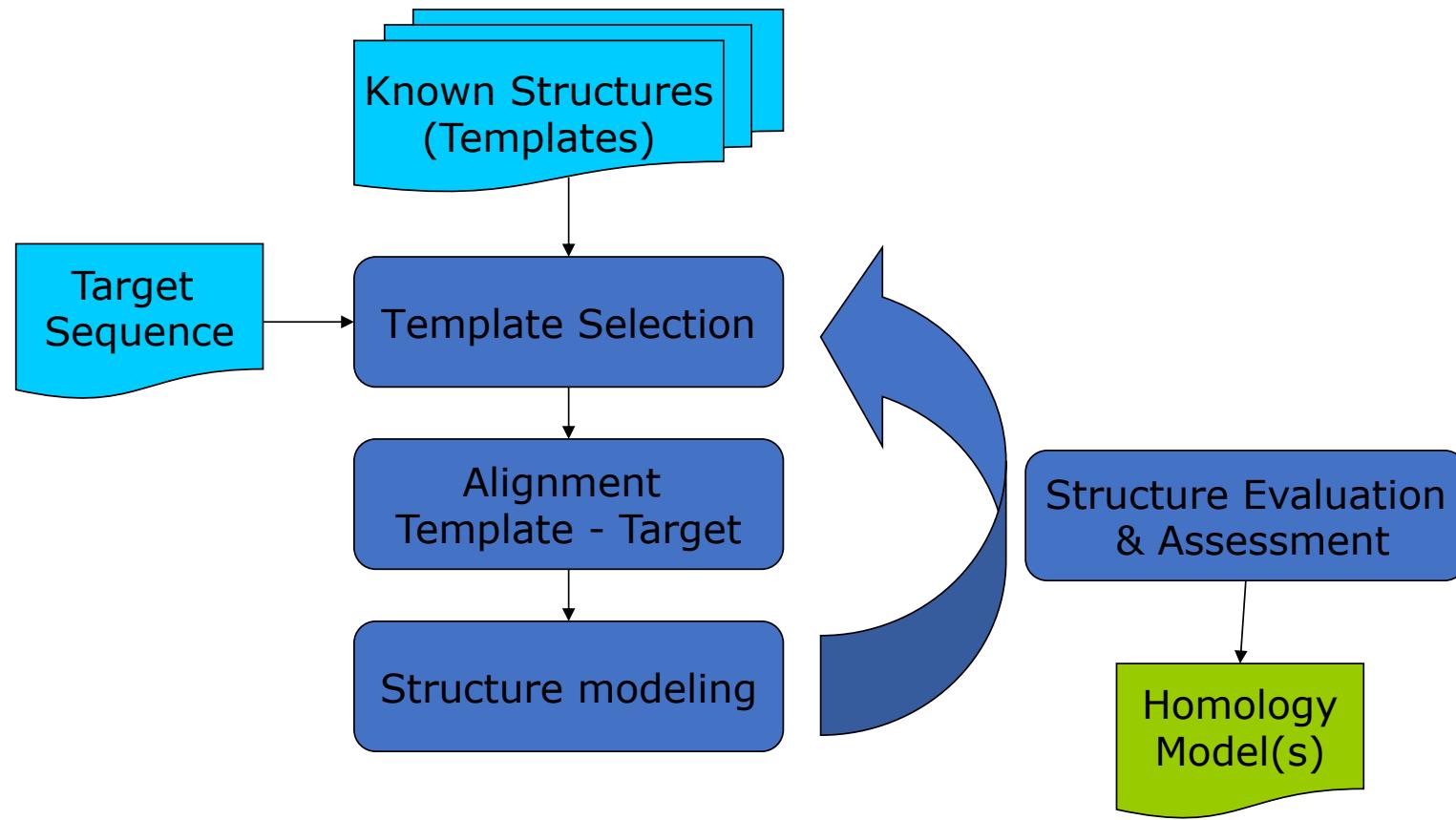


EDN ---KPPQFTWAQWFETOHNMTSQOCTNAMO
7RSA KETAAAKFERQHMDSSSTAASSSNYCNQMMK
aaaaaaaaaaaa aaaaaaaaaa

Allineamento errato



Templato errato



CASP4 Target T0111

1. Protein Name

enolase

2. Organism Name

Escherichia coli

3. Number of amino acids (approx)

431

4. Accession number

P08324

5. Sequence Database

Swiss-prot

6. Amino acid sequence

```
SKIVKIIGREIIDSRGNPTVEAEVHLEGGFVGMAAPSGASTGSREALEL  
RDGDKSRFLGKGVTKAVAAVNGPIAQALIGDKADQAGIDKIMIDLDGTE  
NKSFKGANAILAVSLANAKAAAAAKGMPLYEHIAELNGTPGKYSMPVPMM  
NIINGGEHADNNVDIQEFMIQPVGARTVKEARIMGSEVFHHLAKVLKARG  
MNTAVGDEGGYAPNLGSNAEALAVIAEAVKAAGYELGKDITLAMDCAASE  
FYKDGGKVVLAGEGNKAKTSEEFTFLEELTKQYPIVSIEDGLDSDWDGF  
AYQIKVLCDDKIQLVGGDLFVTNTIKLREGIEKGIANSILIKFNQIGSLTE  
TLAAIRMAKDAGYTAVISHRSGETEDATIADLAvgTAAGQIKTGSMSSRSD  
RVAKYNQLRIEEALGEKAPYNGRKEIKGQA
```

7. Additional Information

oligomerization state: dimer in the presence of magnesium by dynamic light scattering and small angle x-ray solution scattering and in the recently solved crystal structure.

8. Homologous Sequence of known structure

yes

9. Current state of the experimental work

Protein supply: overexpressed in E. coli
crystals: grown at 20 C from PEG 3550
diffraction quality: strong data to 2.5 Å with good redundancy

Structure solved by molecular replacement. Currently, the refinement to 2.5 Å resolution is near completion.
Current Rfree 27 % ; R 22 %

10. Interpretable map?

yes

11. Estimated date of chain tracing completion

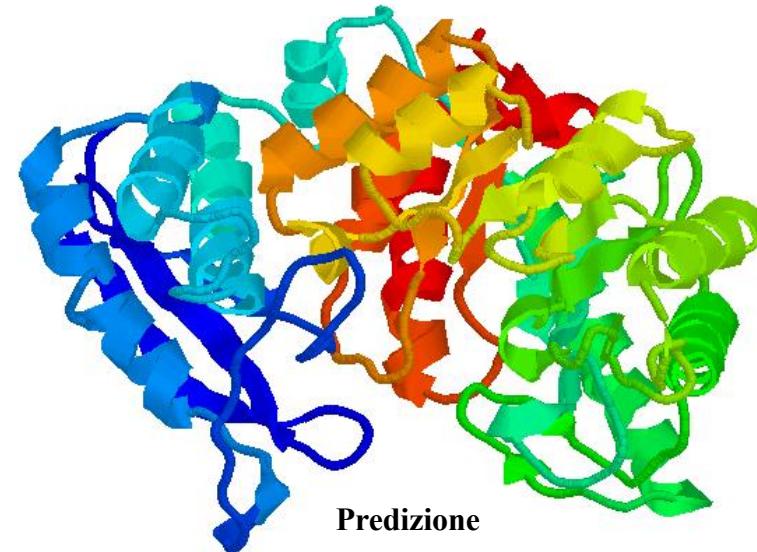
complete

12. Estimated date of public release of structure

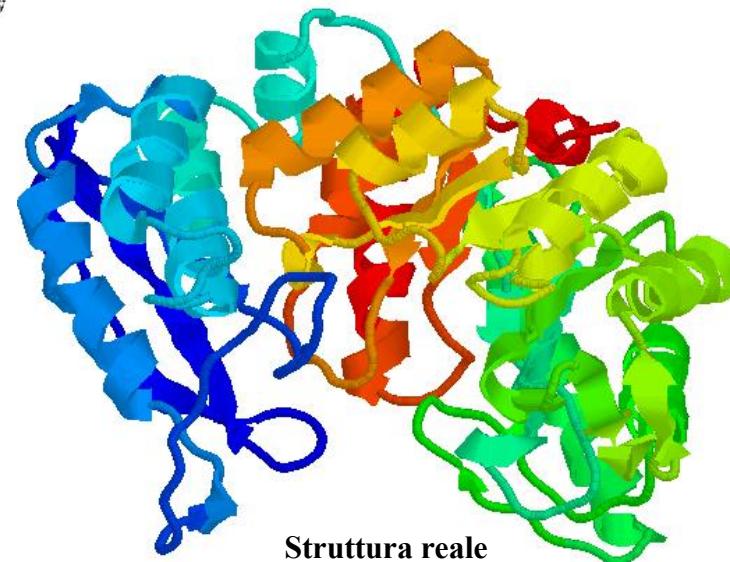
Dec 2000

13. Name

Unavailable until after public release of structure



Predizione



Struttura reale

Alcuni siti web di homology modeling

COMPOSER – felix.bioccam.ac.uk/soft-base.html

MODELLER – guitar.rockefeller.edu/modeller/modeller.html

WHAT IF – www.sander.embl-heidelberg.de/whatif/

SWISS-MODEL – www.expasy.ch/SWISS-MODEL.html

Swiss-Model

<http://www.expasy.ch/swissmod/SWISS-MODEL.html>



SWISS-MODEL

An Automated Comparative Protein Modelling Server

[SWISS-MODEL](#) is a fully automated protein structure homology-modeling server, accessible via the [ExPASy](#) web server, or from the program [DeepView](#) (Swiss Pdb-Viewer). The purpose of this server is to make Protein Modelling accessible to all biochemists and molecular biologists World Wide.

The present version of the server is 3.5 and is under constant improvement and debugging. In order to help us refine the sequence analysis and modelling algorithms, please [report](#) of possible bugs and problems with the modelling procedure.

SWISS-MODEL was initiated in 1993 by Manuel Peitsch, and is now being further developed within the [SIB](#) in collaboration between [GlaxoSmithKline R&D](#) (Geneva) and the [Structural Bioinformatics Group](#) at the Biozentrum (University of Basel). The computational resources for the SWISS-MODEL server are provided by collaboration with the [Advanced Biomedical Computing Center](#) (NCI Frederick, USA).

Methods and Programs used by SWISS-MODEL

- Sequence Alignment:

- BLAST:
Altschul S.F., Gish W., Miller W., Myers E.W., Lipman D.J. Basic local alignment search tool. *J. Mol. Biol.* 215:403-410. (1990)
 - SIM:
Huang, X., and Miller, M. A time-efficient, linear-space local similarity algorithm. *Adv. Appl. Math.* 12,337-367. (1991)
 - ProModII:
Guex, N., and Peitsch, M.C. Structurally corrected multiple alignments.
-

- Comparative Protein Modelling:

- ProMod / ProModII:
[Several publications](#).
-

- Energy Minimisation:

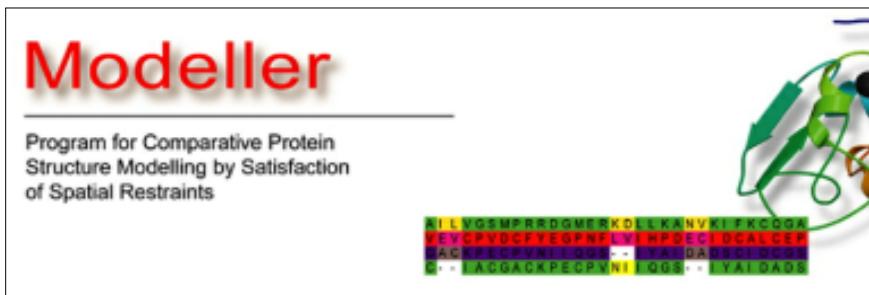
- Gromos96:
Information on this force field can be obtained from the [ETH](#) in Zürich.
-

- Model Evaluation:

- [Swiss-PdbViewer](#):
Provides all necessary tools to evaluate the quality of a model. This feature is thus no longer provided by the SWISS-MODEL server.

Modeller

http://guitar.rockefeller.edu/modeller/about_modeller.shtml



Advanced program for homology modeling

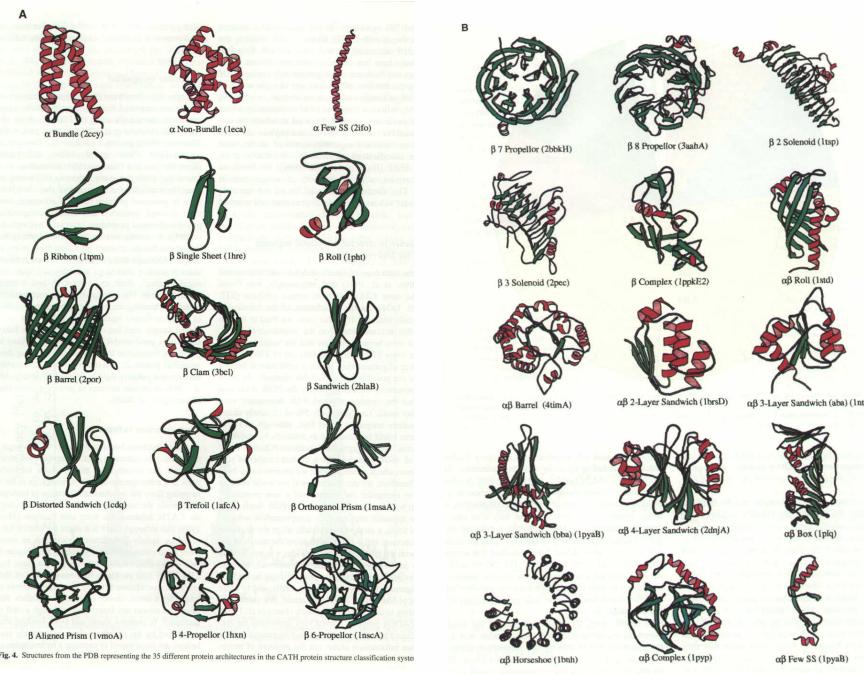
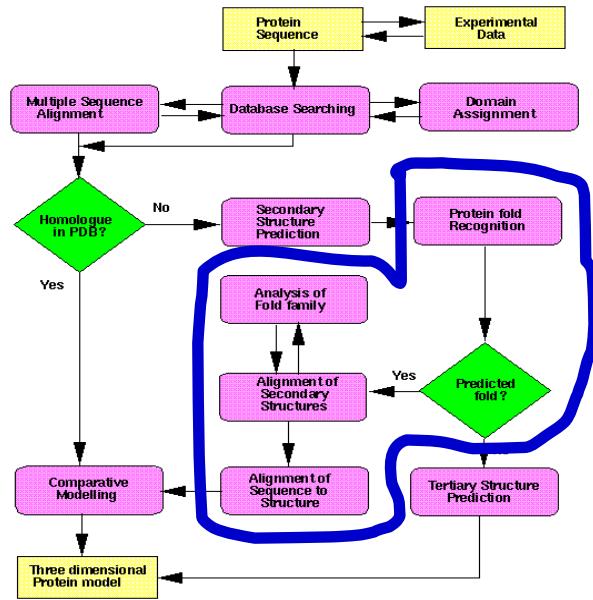
Based on distance constraints

Implemented in several popular modelling packages
such as InsightII

The source is available for unix platforms at the above URL

Fold Recognition

Protein threading is based on two basic observations: that the number of different folds in nature is fairly small (approximately 1300); and that 90% of the new structures submitted to the PDB in the past three years have similar structural folds to ones already in the PDB



- Predizione di sequenza con poca o nessuna similarità con strutture note.
- Osservazione: La natura utilizza solamente un numero limitato di fold diversi (< 1000 ?)
- **Idea del fold recognition:** Cerca di rappresentare la struttura ignota con dei fold conosciuti, valuta quale potrebbe essere quello “giusto”.

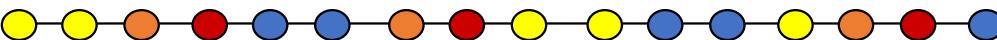
Threading (fold recognition)

La sequenza di input viene confrontata con una libreria di folds noti

Si calcola un punteggio che esprima la compatibilità tra la sequenza e ciascun fold considerato

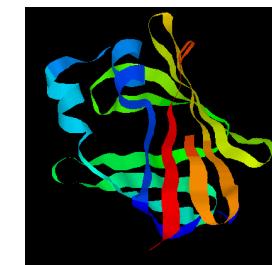
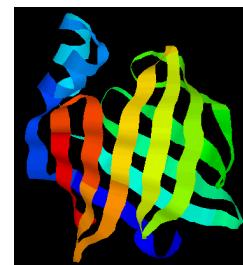
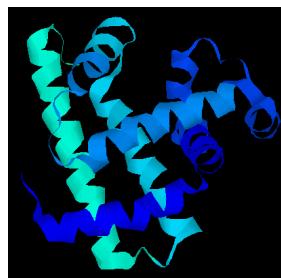
Punteggi statisticamente significativi indicano che la sequenza ha una certa probabilità di assumere la stessa struttura 3D del fold considerato

Input:

 Sequenza

- Donatore H
- Accettore H
- Gly
- Idrofobico

Collezione di folds di proteine note

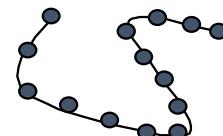
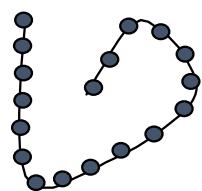
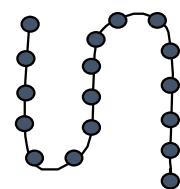


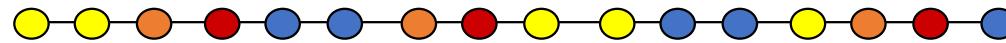
Input:

 Sequenza

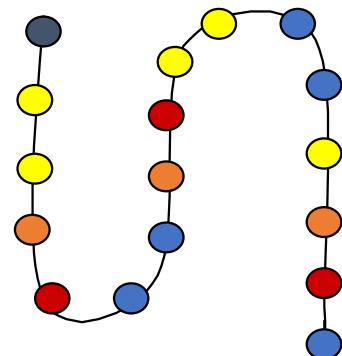
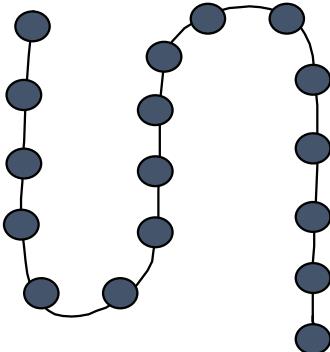
- Donatore H
- Accettore H
- Gly
- Idrofobico

Collezione di folds di proteine note

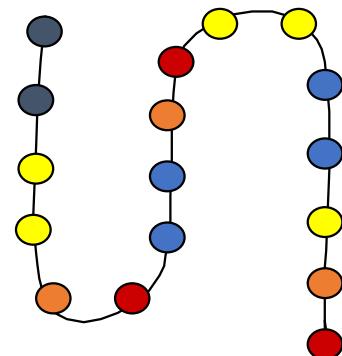




- Donatore H
- Accettore H
- Gly
- Idrofobico

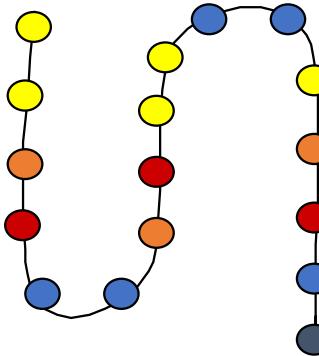


$S=-2$
 $Z= -1$

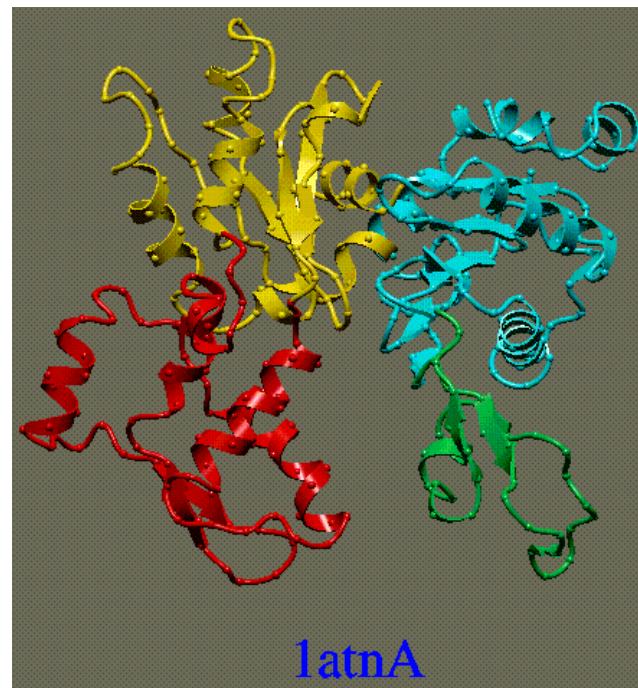
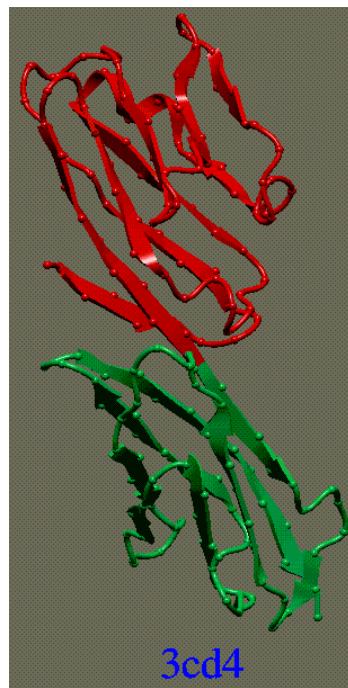


$S=5$
 $Z= 1.5$

$S=20$
 $Z=5$



Chain/Domain Library



Scoring functions for fold recognition

- Ci sono due metodi per valutare la compatibilità sequenza-struttura (1D-3D)
- Nei metodi basati su profili strutturali, per ciascun fold è costruito un profilo basato sulle caratteristiche strutturali del fold e sulla compatibilità di ciascun aminoacido in ciascuna posizione.
- Questa compatibilità è determinata in funzione di struttura secondaria, accessibilità al solvente e caratteristiche di idrofobicità dell'ambiente locale
- Il profilo ha la forma di una funzione matematica adatta al confronto a coppie ed alla programmazione dinamica.

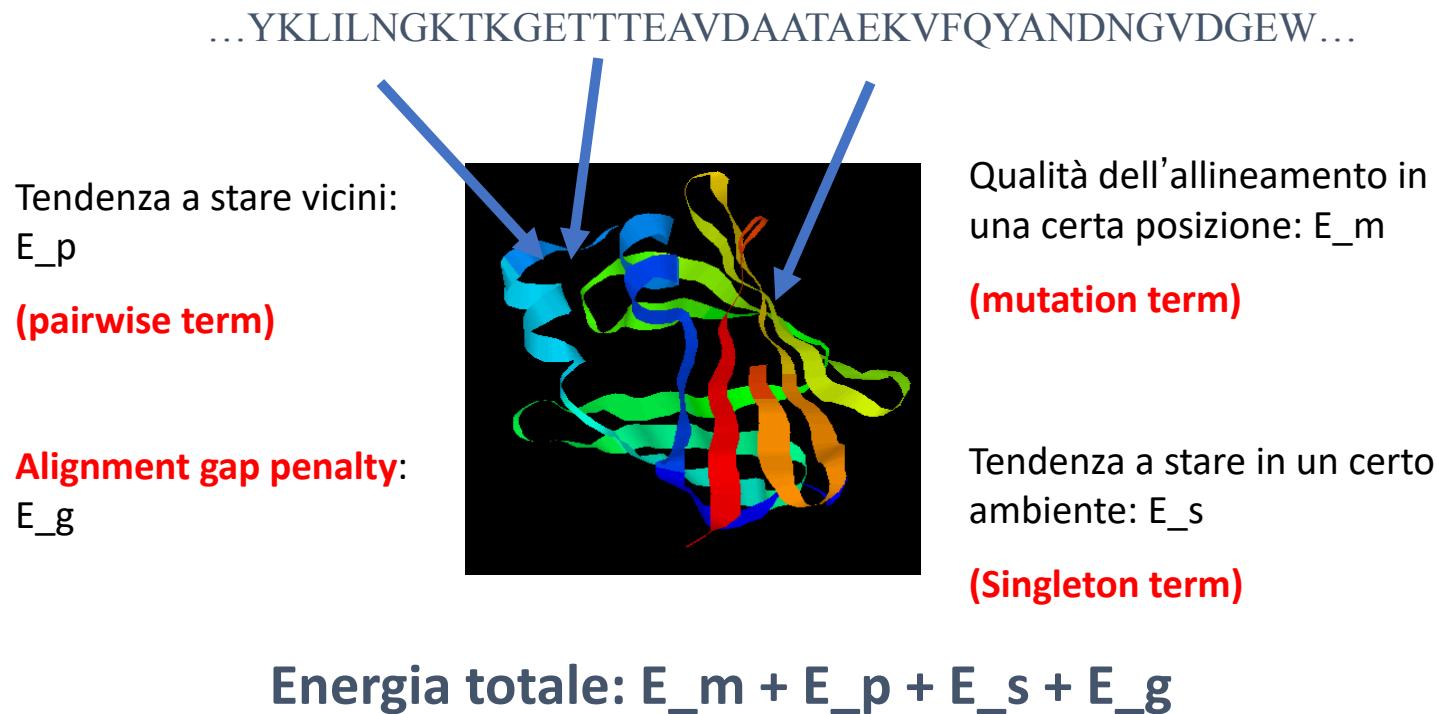
Potenziali di contatto

Basato su tabelle che descrivono punteggi pseudo-energetici per ciascuna interazione tra coppie di aminoacidi.

Rappresenta diversi fold in termini di matrici di distanze.

Somma delle energie sulle coppie di residui in contatto. La somma totale indica la qualità del fit tra sequenza e struttura del fold.

Scoring Function

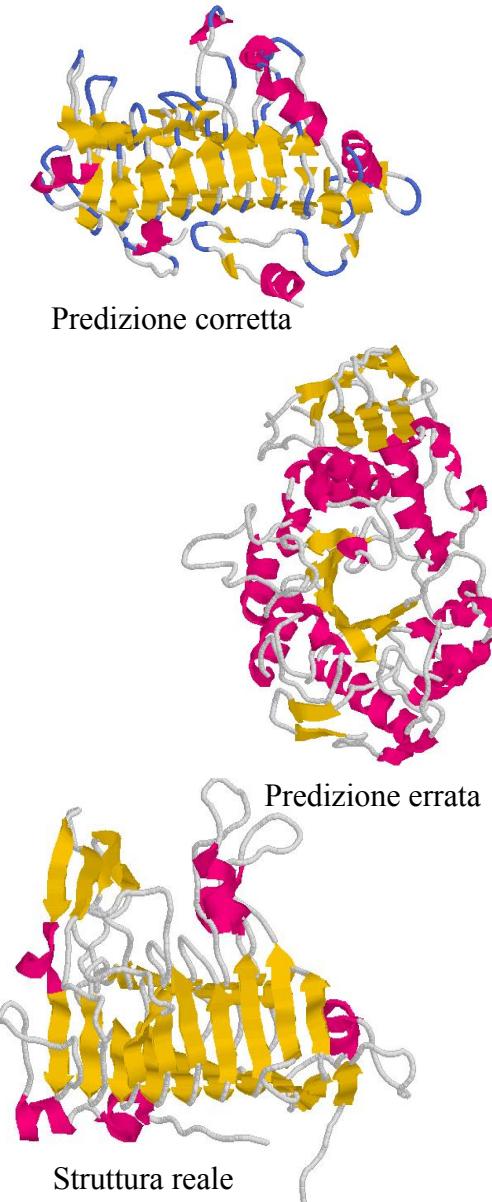


Describe quanto la sequenza assomiglia al tempiato

Dopo che il templato a cui la sequenza si adatta meglio e' stato
selezionato, il modello viene costruito sulla base
dell'allineamento tra il target ed il templato

Cosa si può ottenere dalla fold recognition?

- Predizione del *fold* corretto (media su più metodi) nel 60-70% ca. di casi senza omologia chiara.
 - Stima a priori della qualità del risultato difficile.
 - Riconoscimento di *novel folds* (casi senza soluzione) spesso impossibile.
- I metodi automatici fino ad oggi producono modelli decisamente inferiori a quelli prodotti con l'aiuto di esperti.
 - Alexey Murzin e i dati contenuti in letteratura
 - I server automatici stanno migliorando
- I server *consensus* che combinano diverse predizioni funzionano mediamente meglio dei singoli metodi e danno maggiore affidabilità alle predizioni.
 - Meta Server (<http://bioinfo.pl/meta/>)



Web sites for fold recognition

Profiles:

3D-PSSM - <http://www.bmm.icnet.uk/~3dpssm>

Libra I - http://www.ddbj.nig.ac.jp/htmls/E-mail/libra/LIBRA_I.html

UCLA DOE - <http://www.doe-mbi.ucla.edu/people/frsver/frsver.html>

Contact potentials

123D - <http://www-Immb.ncifcrf.gov/~nicka/123D.html>

Profit - <http://lore.came.sbg.ac.at/home.html>

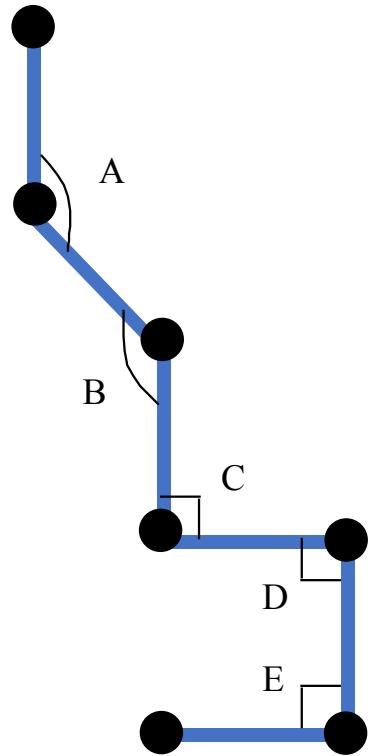
Ab initio methods for modelling

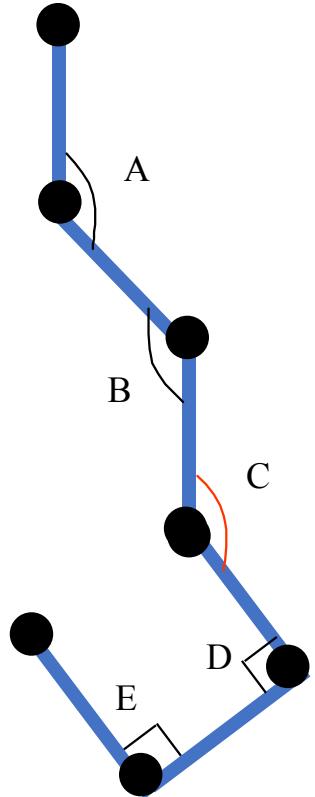
NO allineamento

NO struttura nota

Costruire una funzione empirica che descriva le forze
di interazione

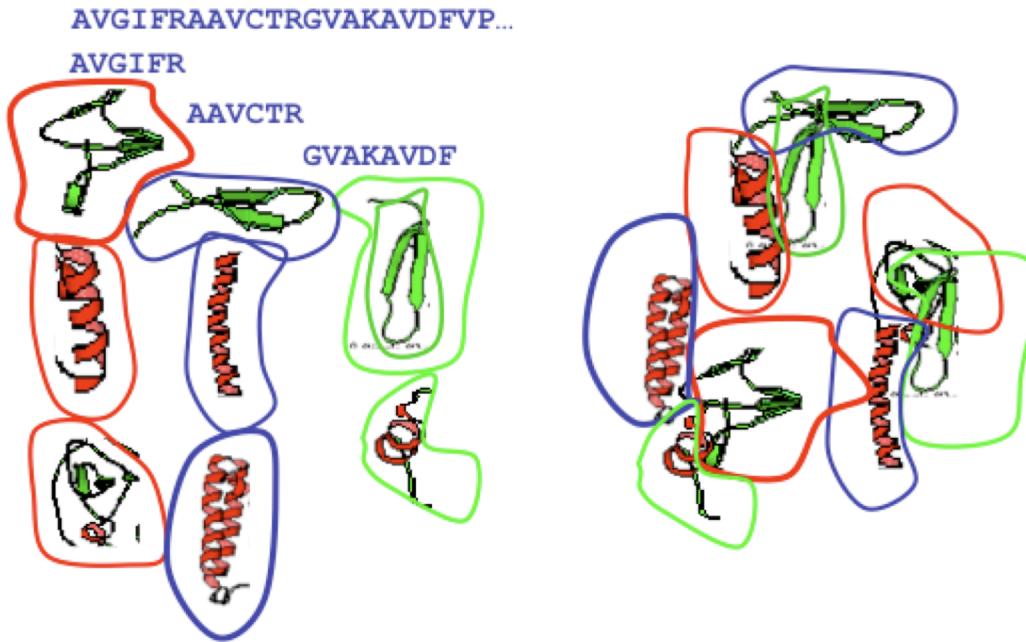
Esplorare lo spazio conformazionale per massimizzare
funzione di merito





Rosetta – David Baker

- Based on the assumption that the distribution of conformations sampled by a local segment of the polypeptide chain is reasonably approximated by the distribution of structures adopted by that sequence and closely related sequences in known protein structures.
- Fragment libraries for all possible three and nine residue segments of the chain are extracted from PDB by profile methods



Bystroff and Baker, JMB, 1998

AVGIFRAAVCTRGVAKAVDFV...
AVGIFR

AAVCTR

GVAKAVDF



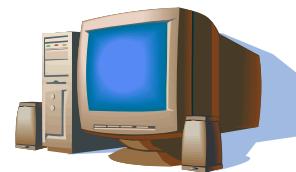
Optimize and score

CASP/CAFASP

- CASP: Critical Assessment of Structure Prediction
- CAFASP: Critical Assessment of Fully Automated Structure Prediction

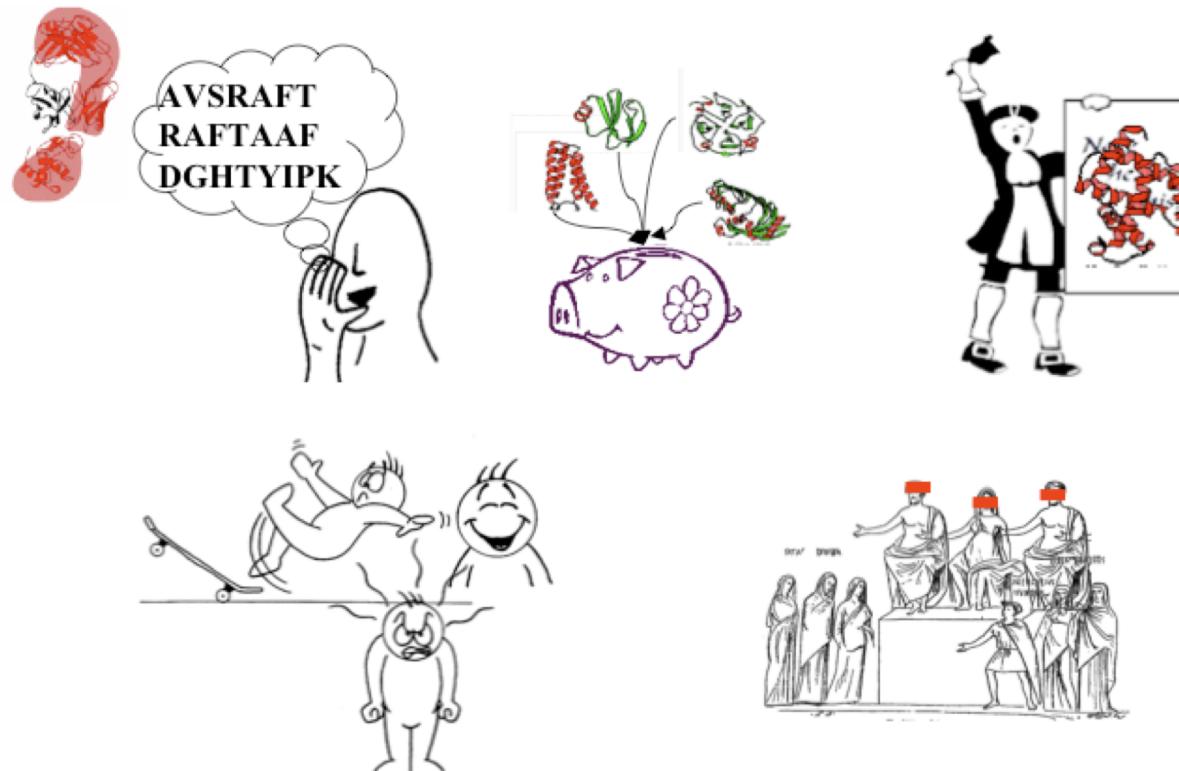


CASP
Predictor

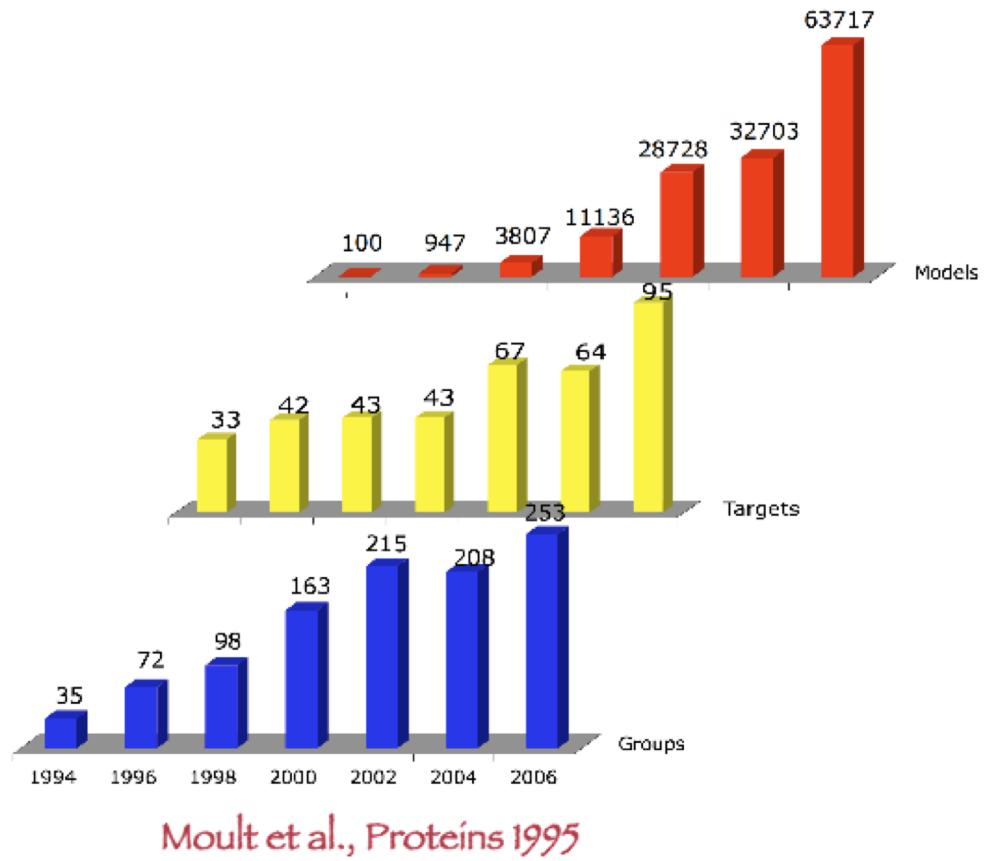


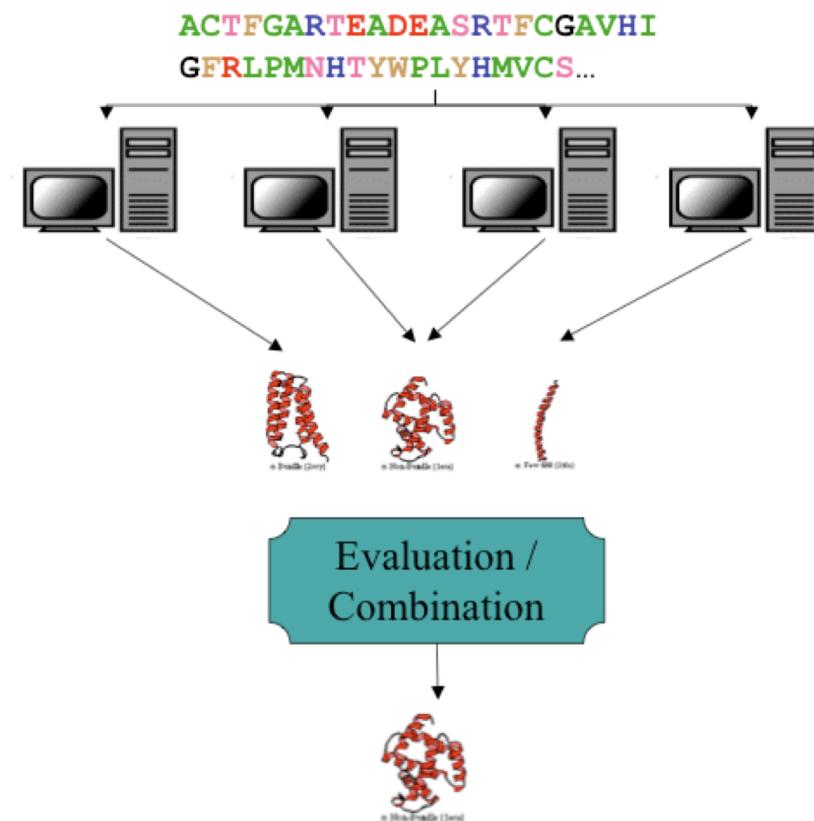
CAFASP
Predictor

1. Won't get tired
2. High-throughput



Moult et al., Proteins 1995







- Il problema del folding delle proteine è stato “risolto” ???
- Dichiarazioni contrastanti fino a circa dieci anni fa.
- *Critical Assessment of Techniques for Protein Structure Prediction*
 - “blind test” che coinvolge tutti i principali gruppi, ripetuto ogni 2 anni
 - CASP-5 (e CAFASP-3) nel 2002
 - Oltre 250 gruppi di predittori, 65 targets
- Cerca di misurare lo stato dell’arte ed i miglioramenti in tutti i maggiori settori della predizione di strutture proteiche
 - (Stabilisce un ranking dei migliori gruppi)

C
A
S
P
5



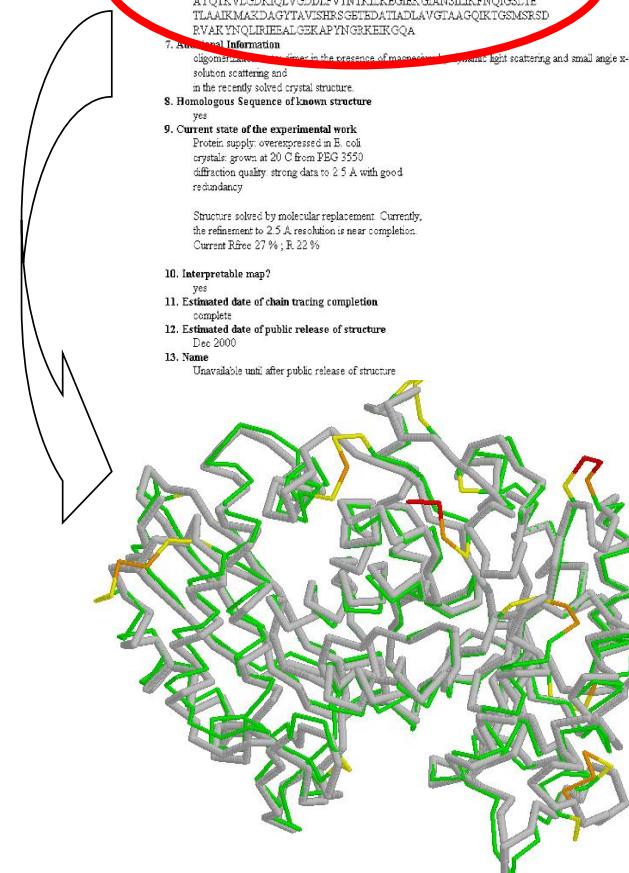
5

Le principali categorie del CASP:

- *Homology modelling*
- *Fold recognition*
- *Ab initio / novel folds*
- Struttura secondaria

CASP4 Target T0111

1. Protein Name
endo-
2. Organism Name
Escherichia coli
3. Number of amino acids (approx)
431
4. Accession number
P08324
5. Sequence Database
Swiss-Prot
6. Amino acid sequence
SKIVKIGREIISGRNPTVEAEVHLEGGFVGMAAAPSGASTGSREALEL
RDGEKSRFLGRGSGTKAVAAVNGPQIAQALIGKDARDQAGIDKIMIDLDGTE
NKSFKFGANALLAVSLANAKAAAAAKGMPLYEHIAELNGTPGKYSMPVPM
NINNGGEHADNNVJQEFMIPVYGAKTKEARMGSEVFHLAATVRAKG
MTAVGDEGGYAPNLGSNABALVALAVALAVALAAGVYELGKDITLAMDCAASE
FYRKDRKVYLAGEGNKAFTSEEFITIEFLELITCQYIVSIEDQGLDESWDGF
AYQTQVLCDKIQVLCDDLFVINTKILRKGHEKGLANSILKFNQGSCTE
TLAIRKMAKDAGYTAVISHRSGETEDATIAIDLAVGTAAGQIKTGMSRSRD
EVAKYNQURIEEALGEKAPYNGRKEKGQA
7. Ab initio Information
cigmentation, colour in the presence of macrocyclic ligands, light scattering and small angle x-ray
solution scattering and
in the recently solved crystal structure
8. Homologous Sequence of known structure
yes
9. Current state of the experimental work
Protein supply: overexpressed in E. coli
crystals grown at 20 °C from PEG 3550
diffraction quality: strong data to 2.5 Å with good
redundancy
Structure solved by molecular replacement: Currently,
the refinement to 2.5 Å resolution is near completion.
Current Rfree 27 %; R 22 %
10. Interpretable map?
yes
11. Estimated date of chain tracing completion
complete
12. Estimated date of public release of structure
Dec 2000
13. Name
Unavailable until after public release of structure



- **Target:**
 - Sequenza di cui si cerca la struttura
- **Template:**
 - Sequenza con struttura nota, “stampo” per il modello
- ***Comparative o homology modeling***
 - Ricerca in database
 - Modello costruito da struttura omologa
- ***Fold recognition (Threading)***
 - Tenta di riconoscere omologie remote
 - Approcci differenti che utilizzano struttura secondaria, profili di sequenza, funzioni energetiche specializzate, ...

