

Introduction to Protein Structure Prediction

BMI/CS 776

www.biostat.wisc.edu/bmi776/

Mark Craven

craven@biostat.wisc.edu

Spring 2011

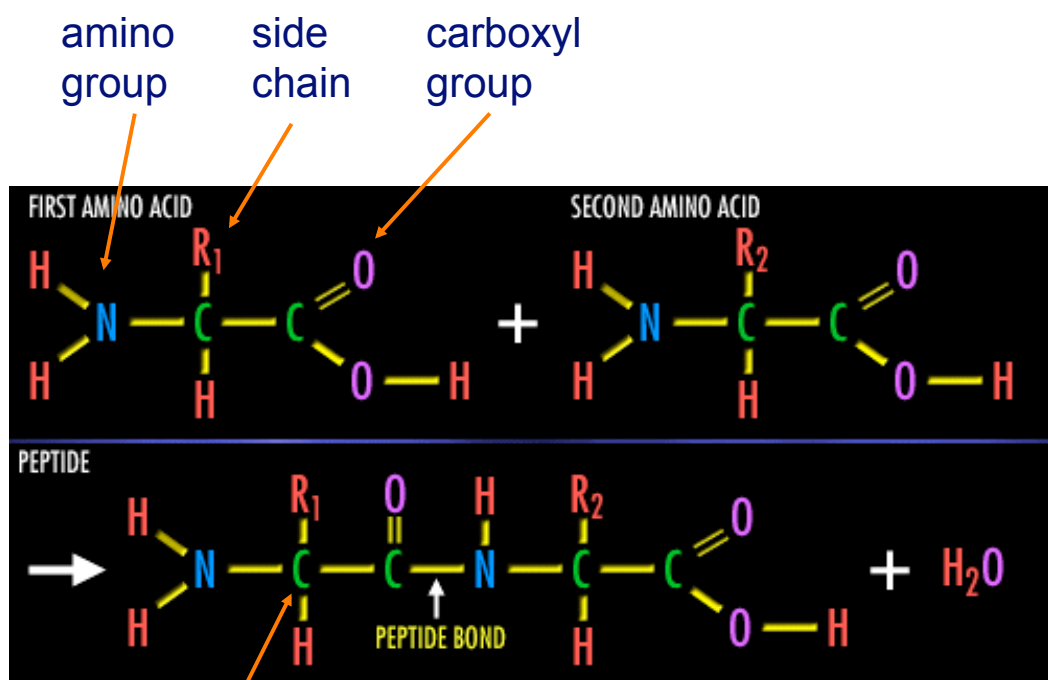
The Protein Folding Problem

- we know that the function of a protein is determined in large part by its 3D shape (*fold, conformation*)
- can we predict the 3D shape of a protein given only its amino-acid sequence?

Protein Architecture

- proteins are polymers consisting of amino acids linked by *peptide bonds*
- each amino acid consists of
 - a central carbon atom
 - an amino group, NH_2
 - a carboxyl group, COOH
 - a side chain
- differences in side chains distinguish different amino acids

Amino Acids and Peptide Bonds



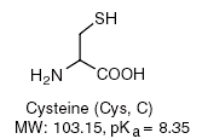
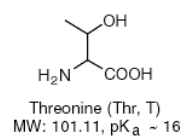
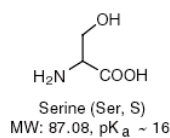
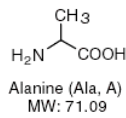
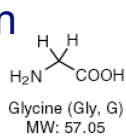
α carbon (common reference point for coordinates of a structure)

Amino Acid Side Chains

side chains vary in

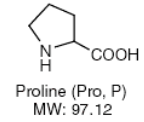
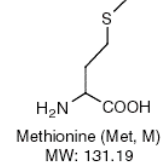
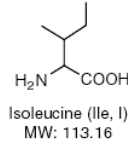
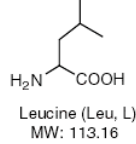
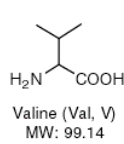
- shape
- size
- charge
- polarity

Small

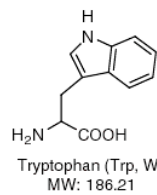
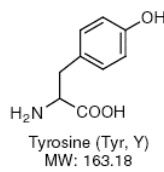
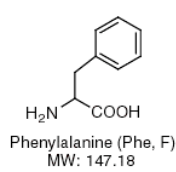


Nucleophilic

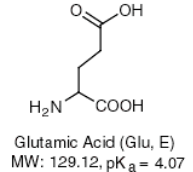
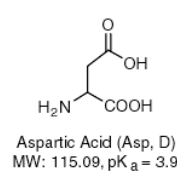
Hydrophobic



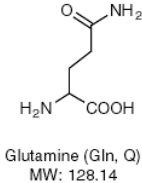
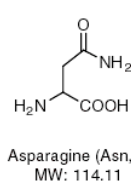
Aromatic



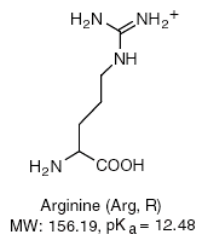
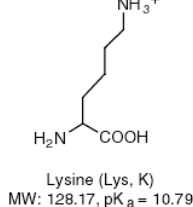
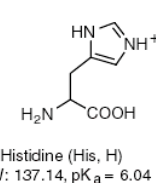
Acidic



Amide



Basic



What Determines Conformation?

- in general, the amino-acid sequence of a protein determines the 3D shape of a protein [Anfinsen et al., 1950s]
- but some qualifications
 - all proteins can be denatured
 - some proteins are inherently *disordered* (i.e. lack a regular structure)
 - some proteins get folding help from *chaperones*
 - there are various mechanisms through which the conformation of a protein can be changed in vivo
 - post-translational modifications such as *phosphorylation*
 - *prions*
 - etc.

What Determines Conformation?

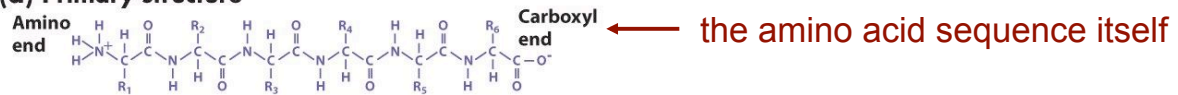
- Which physical properties of the protein determine its fold?
 - rigidity of the protein backbone
 - interactions among amino acids, including
 - electrostatic interactions
 - van der Waals forces
 - volume constraints
 - hydrogen, disulfide bonds
 - interactions of amino acids with water

Levels of Description

- protein structure is often described at four different scales
 - primary structure
 - secondary structure
 - tertiary structure
 - quaternary structure

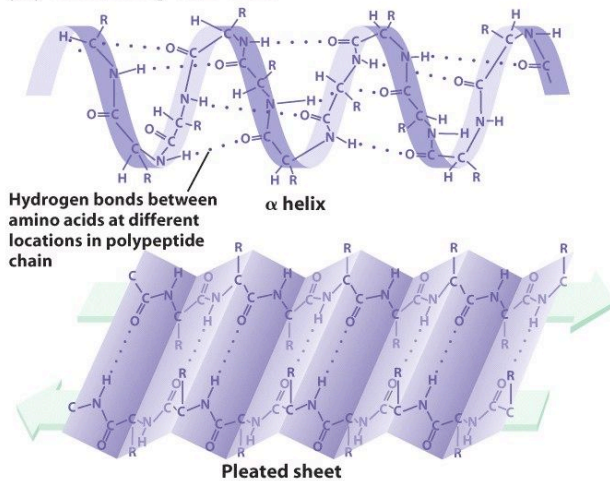
Levels of Description

(a) Primary structure



“local” description of structure:
describes it in terms of certain
common repeating elements

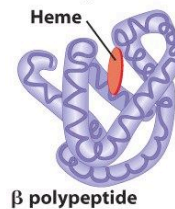
(b) Secondary structure



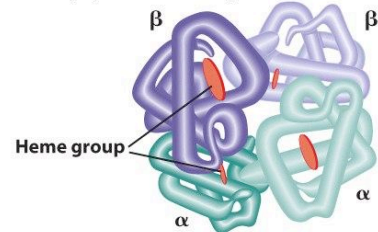
3D conformation
of a polypeptide

3D conformation
of a complex of
polypeptides

(c) Tertiary structure



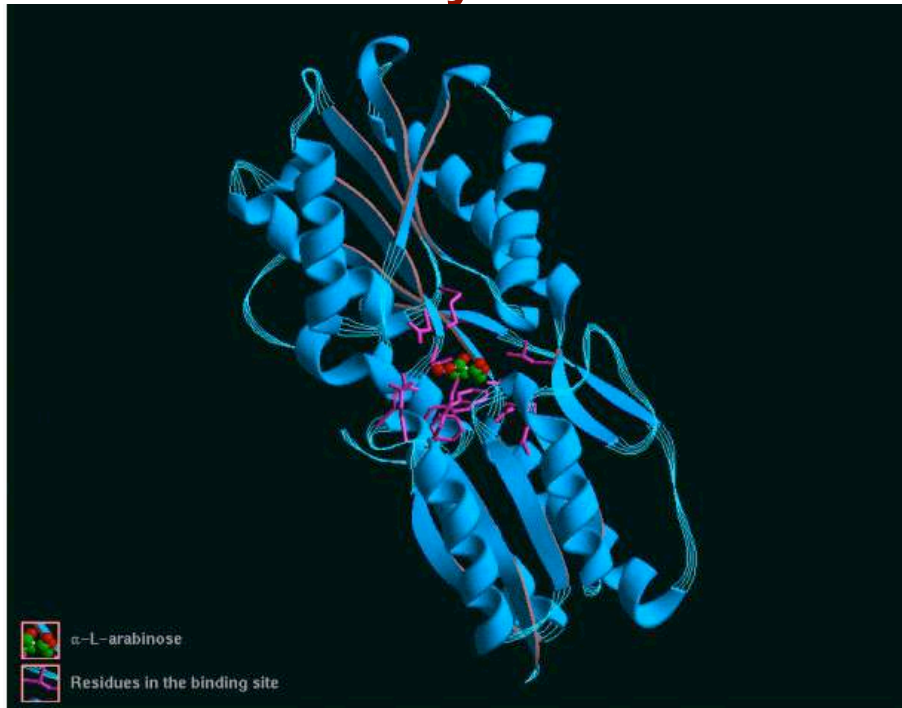
(d) Quaternary structure



Secondary Structure

- secondary structure refers to certain common repeating structures
- it is a “local” description of structure
- two common secondary structures
 - α helices
 - β strands/sheets
- a third category, called *coil* or *loop*, refers to everything else

Ribbon Diagram Showing Secondary Structures



Determining Protein Structures

- protein structures can be determined experimentally (in most cases) by
 - x-ray crystallography
 - nuclear magnetic resonance (NMR)
- but this is very expensive and time-consuming
- there is a large sequence-structure gap
 - \approx 500K protein sequences in SwissProt database
 - $<$ 67K protein structures in PDB database
- key question: can we predict structures by computational means instead?

Types of Protein Structure Predictions

- prediction in 1D
 - secondary structure
 - solvent accessibility (which residues are exposed to water, which are buried)
 - transmembrane helices (which residues span membranes)
- prediction in 2D
 - inter-residue/strand contacts
- prediction in 3D
 - homology modeling
 - fold recognition (e.g. via threading)
 - *ab initio* prediction (e.g. via molecular dynamics)

Prediction in 1D, 2D and 3D

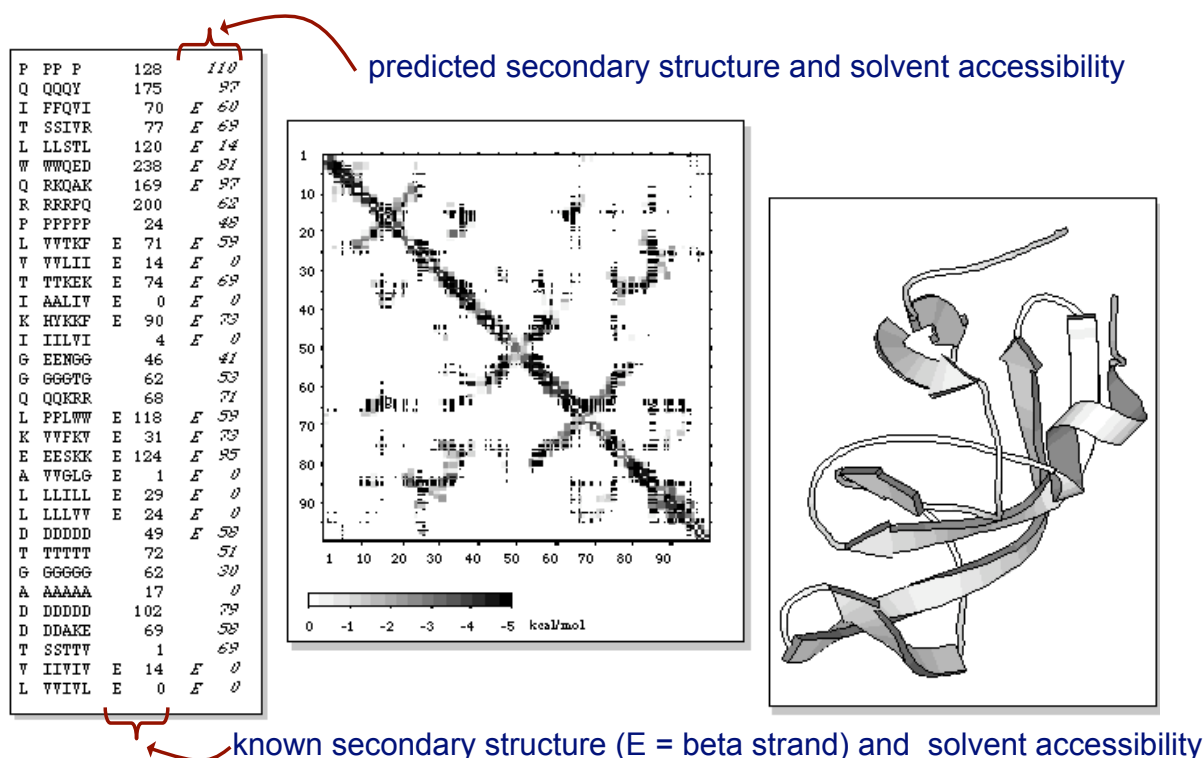


Figure from B. Rost, "Protein Structure in 1D, 2D, and 3D", *The Encyclopaedia of Computational Chemistry*, 1998

Prediction in 3D

- **homology modeling**

given: a query sequence Q , a database of protein structures
do:

- find protein P such that
 - structure of P is known
 - P has high sequence similarity to Q
- return P 's structure as an approximation to Q 's structure

- **fold recognition** (threading)

given: a query sequence Q , a database of known folds
do:

- find fold F such that Q can be aligned with F in a highly compatible manner
- return F as an approximation to Q 's structure

Prediction in 3D

- **“fragment assembly”** (Rosetta)

given: a query sequence Q , a database of structure fragments
do:

- find a set of fragments that Q can be aligned with in a highly compatible manner
- return fragment assembly as an approximation to Q 's structure

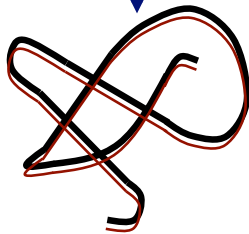
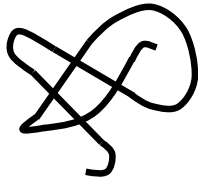
- **molecular dynamics**

given: a query sequence Q

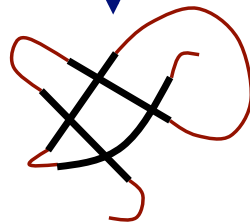
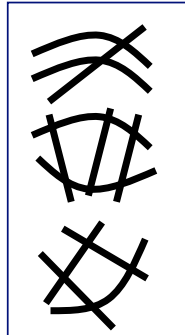
do: use laws of Physics to simulate folding of Q

Prediction in 3D

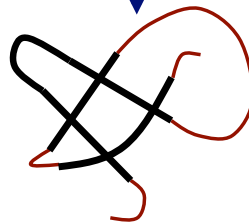
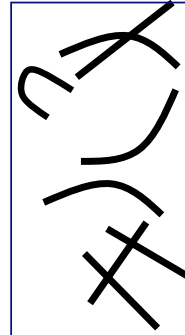
homology
modeling



threading



fragment assembly
(Rosetta)



molecular
dynamics

