# Protein Shape Descriptors

## Patrice Koehl
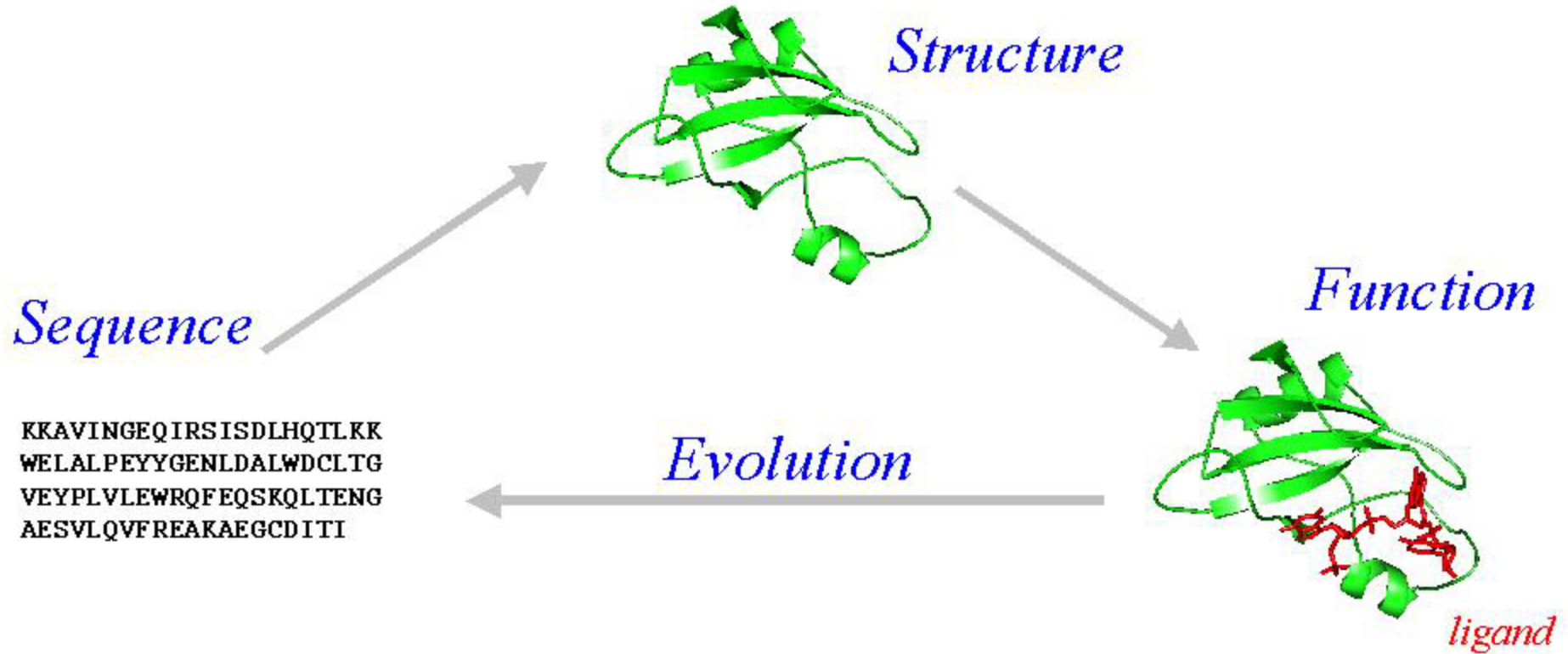
### Stanford University
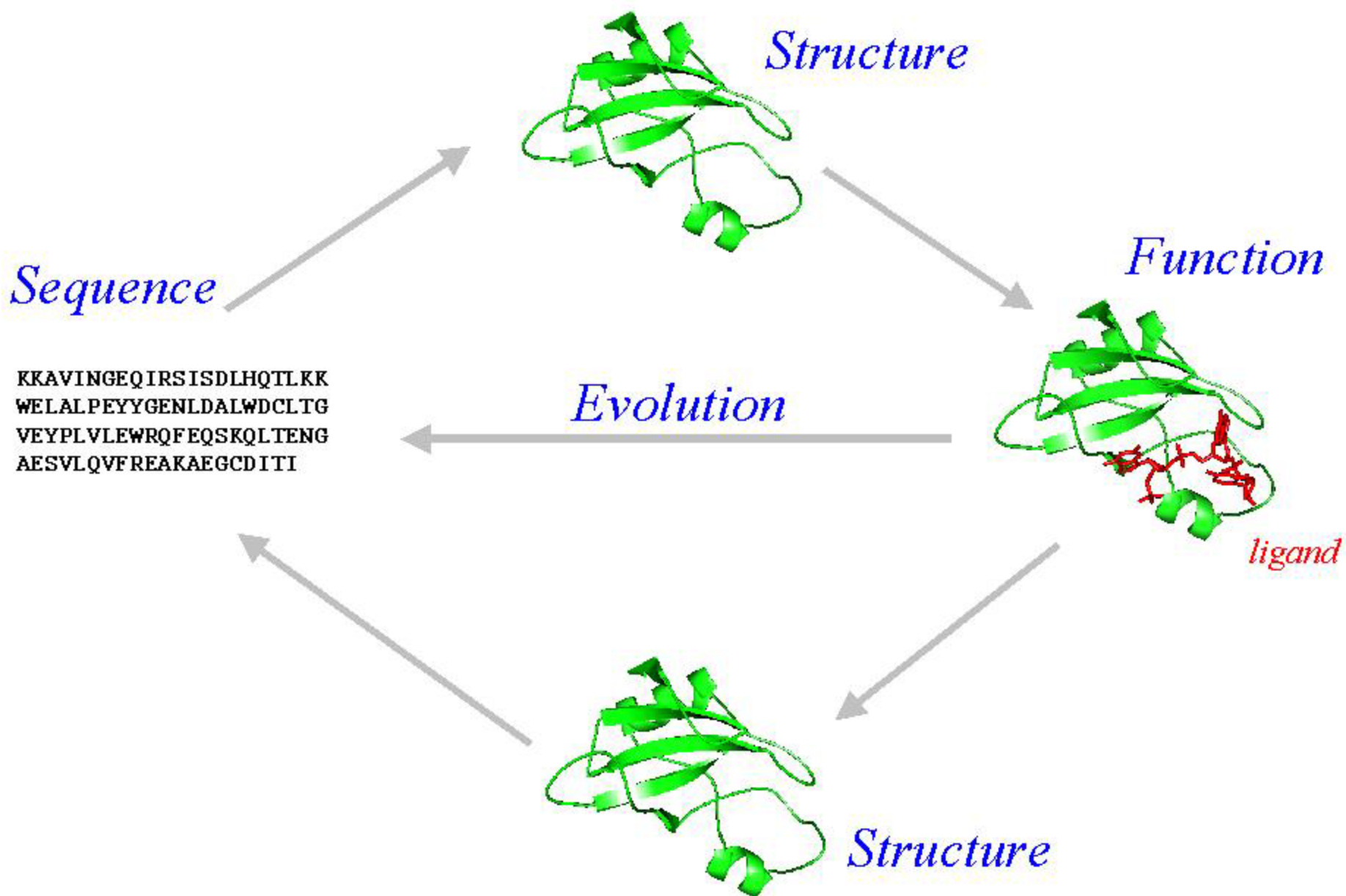
http://csb.stanford.edu/koehl/

# From Sequence to Function and Back…



*Structure*

*Function*

*Sequence*

KKAVINGEQIRSISDLHQTLKK
WELALPEYYGENLDALWDCLTG
VEYPLVLEWRQFEQSKQLTENG
AESVLQVFREAKAEGCDITI

*ligand*

# From Sequence to Function and Back…



*Structure*

*Function*

*Sequence*

KKAVINGEQIRSISDLHQTLKK
WELALPEYYGENLDALWDCLTG
VEYPLVLEWRQFEQSKQLTENG
AESVLQVFREAKAEGCDITI

*Evolution*

*ligand*

# From Sequence to Function and Back…



*Structure*

*Function*

*Sequence*

KKAVINGEQIRSISDLHQTLKK
WELALPEYYGENLDALWDCLTG
VEYPLVLEWRQFEQSKQLTENG
AESVLQVFREAKAEGCDITI

*Evolution*

*ligand*

*Structure*

# From Sequence to Function and Back…



Structure

*Sequence*

*Function*

*ligand*

# Outline

•Introduction

> *What is a Protein?*

•Protein Energy Functions

> *Computational Geometry Tools*

•Classifying Proteins

> *The Shapes of Protein Structures*

# Outline

- Introduction

  *What is a Protein?*

- Protein Energy Functions

  *Computational Geometry Tools*
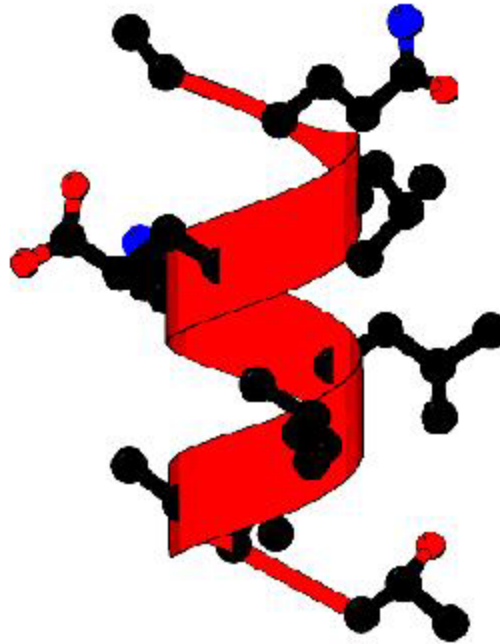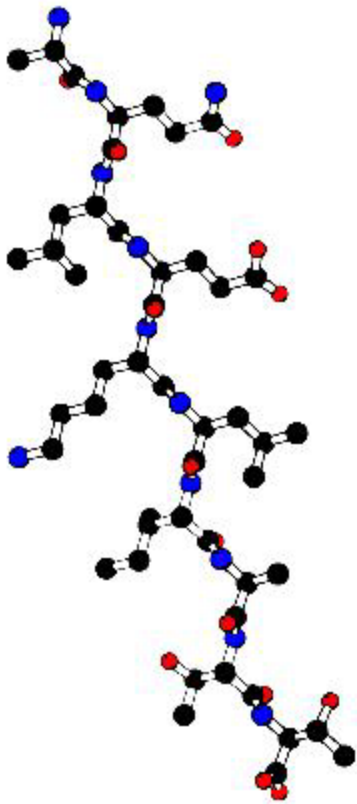
- Classifying Proteins

  *The Shapes of Protein Structures*

# What is a Protein ?

*Primary Structure*   *Secondary Structure*   *Tertiary Structure*



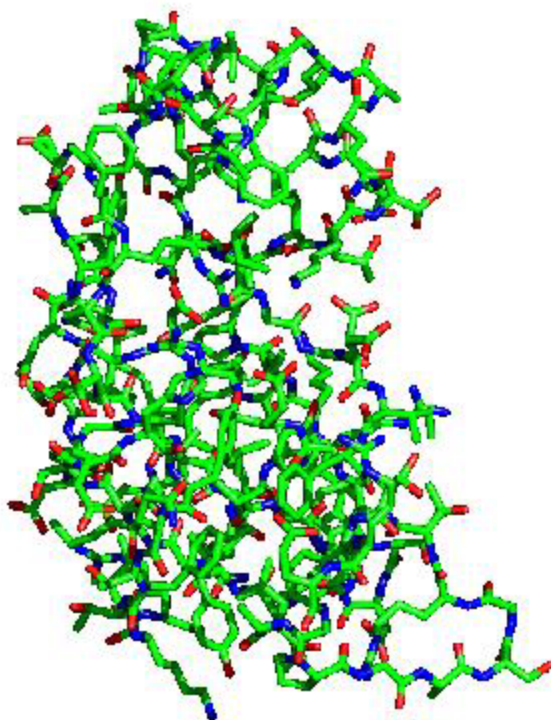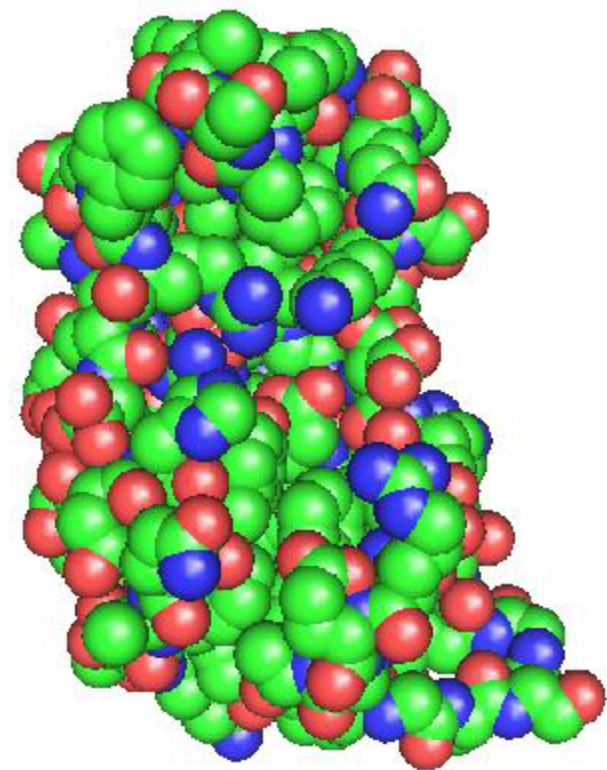Sequence of amino acids                    Native protein

# Protein Representations



*Cartoon*                    *Stick*              *Space-filling Model*

# Outline

- Introduction

  *What is a Protein?*

- Protein Energy Functions

  *Computational Geometry Tools*

- Classifying Proteins

  *The Shapes of Protein Structures*

# Energy of a Protein

## *Bonded Interactions*
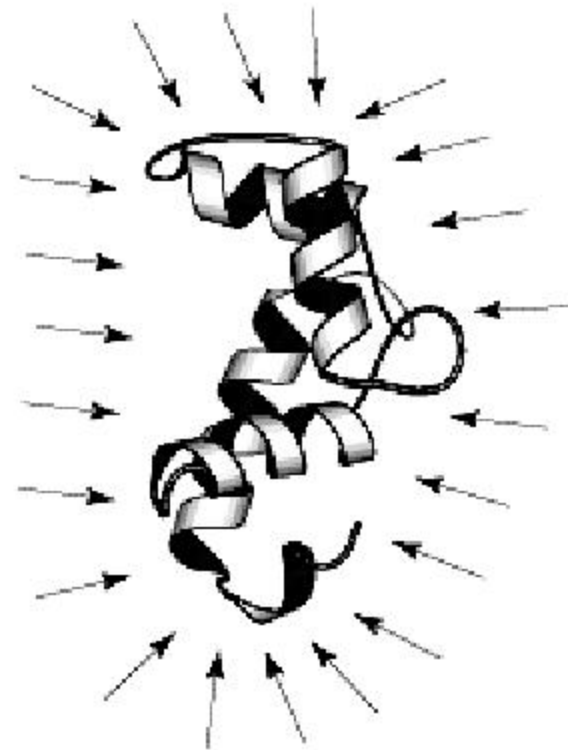
Bonds, Angles, Dihedral angles

## *Non Bonded Interactions*

van der Waals interactions, Electrostatics
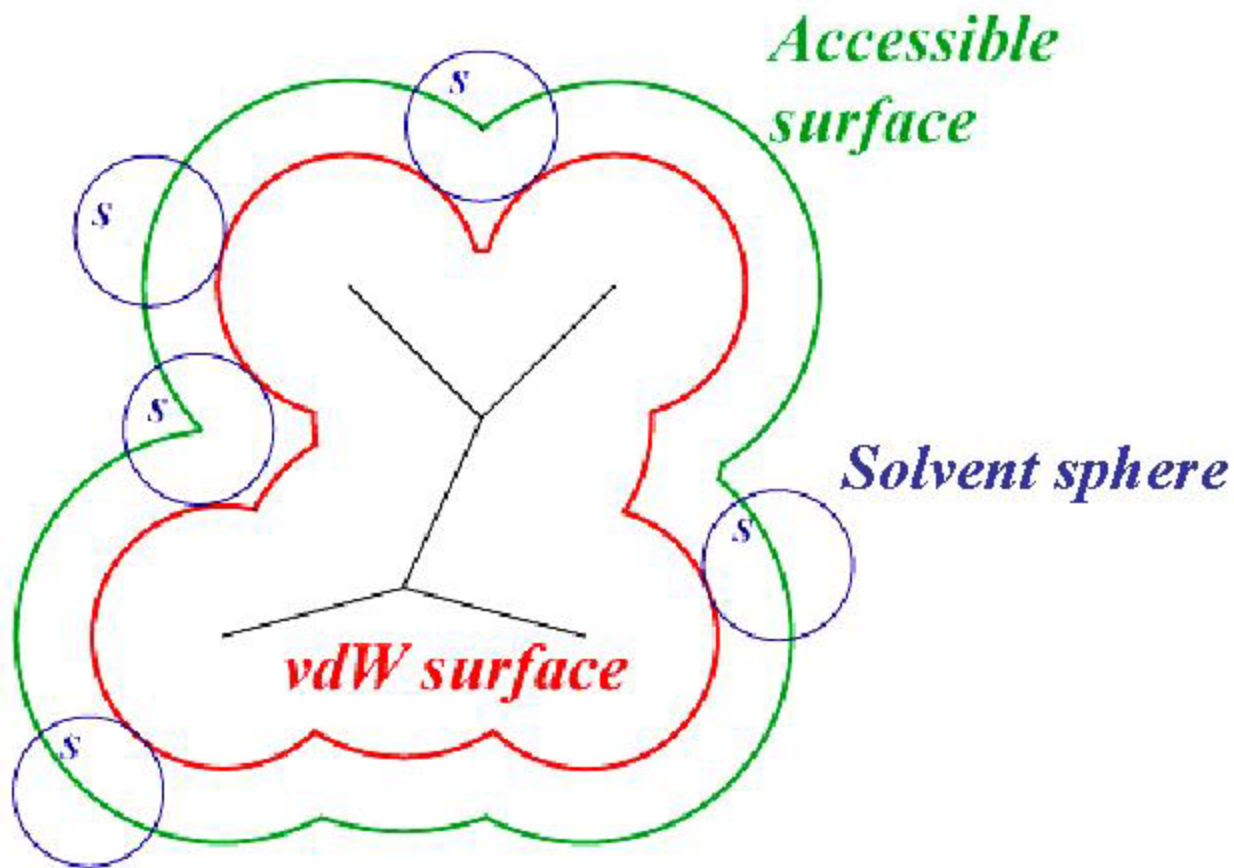
## *Solvent*

Most difficult

# Solvent

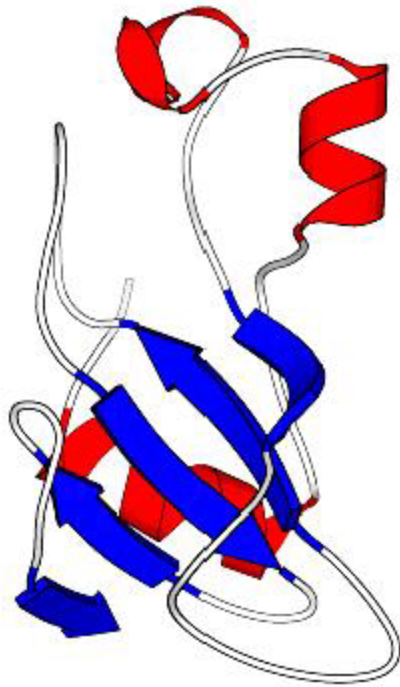*Explicit or Implicit ?*



$$G_{sol} = G_{pol} + G_{cav} + G_{vdW}$$

# Solvation Potential


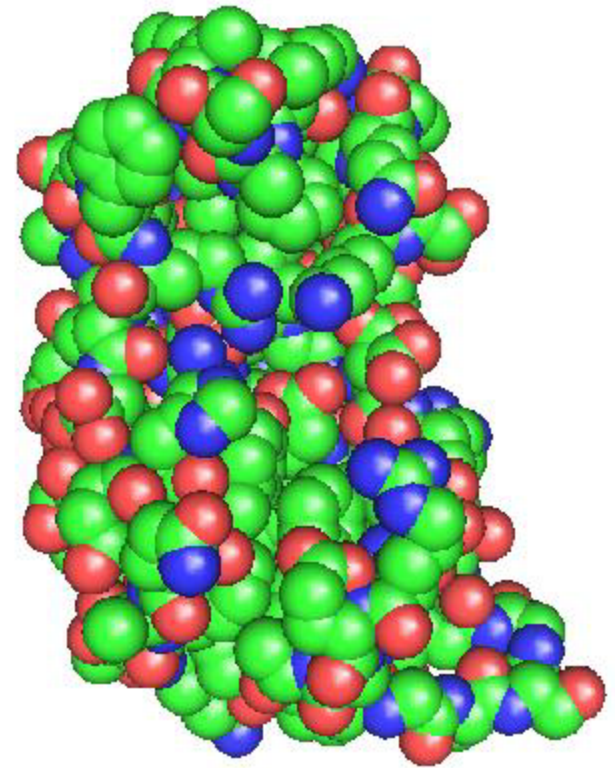
Accessible surface

Solvent sphere

vdW surface

$$G_{cav} + G_{vdW} = \sum_{k=1}^{N} \sigma_k SA_k$$

Need Surface and Volume
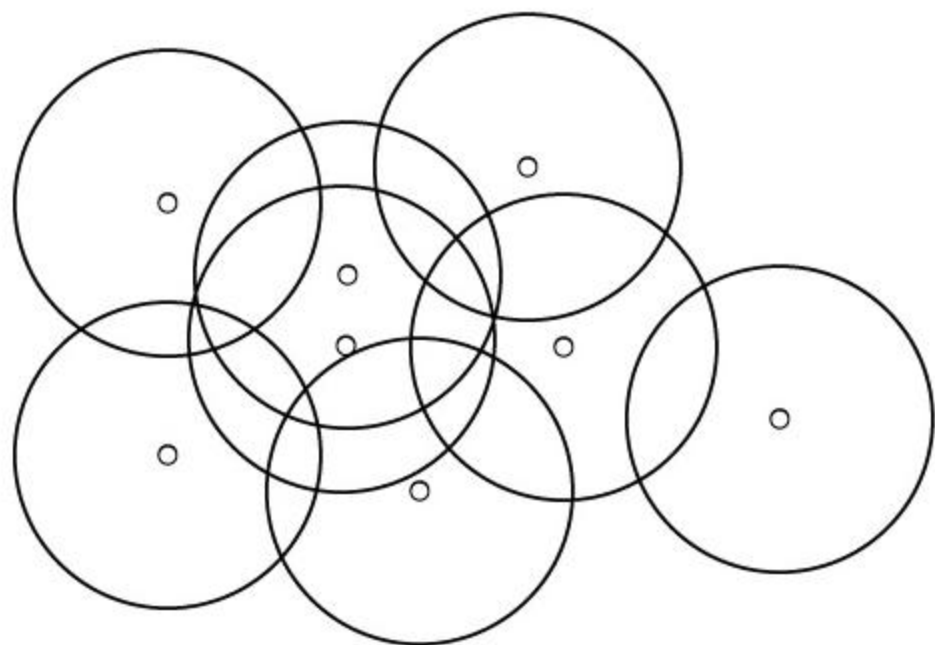
# Geometry of Protein Structure
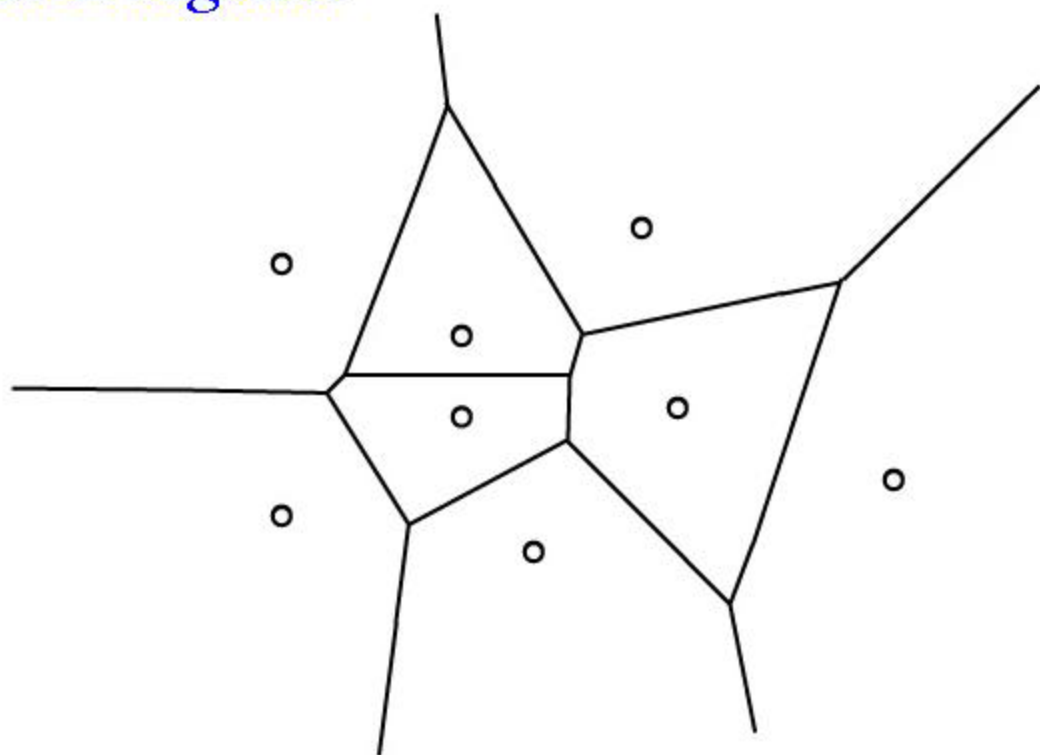


*Cartoon*

*Space-filling Model*

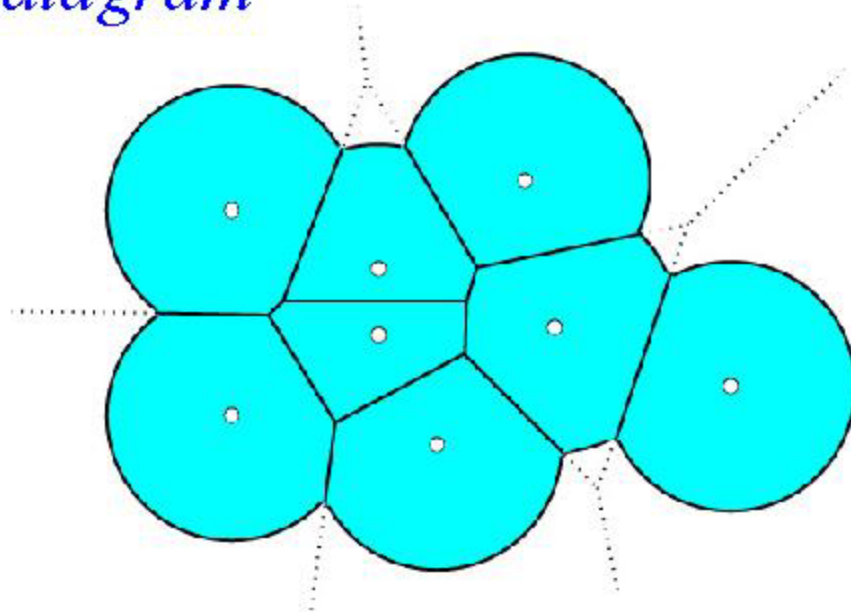# Computing the Surface Area
# and Volume of a Union of Balls

# Computing the Surface Area
# and Volume of a Union of Balls
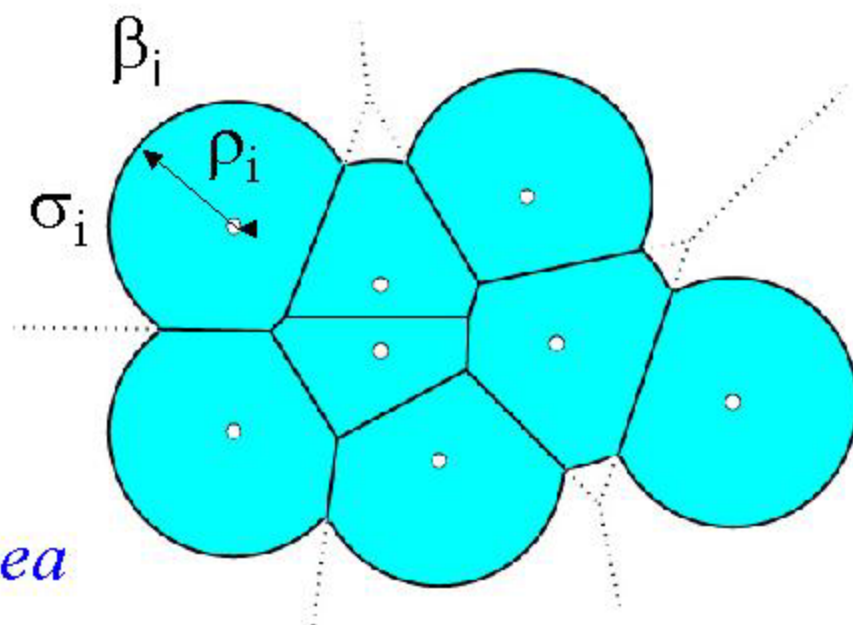
*Power Diagram*

# Computing the Surface Area
# and Volume of a Union of Balls
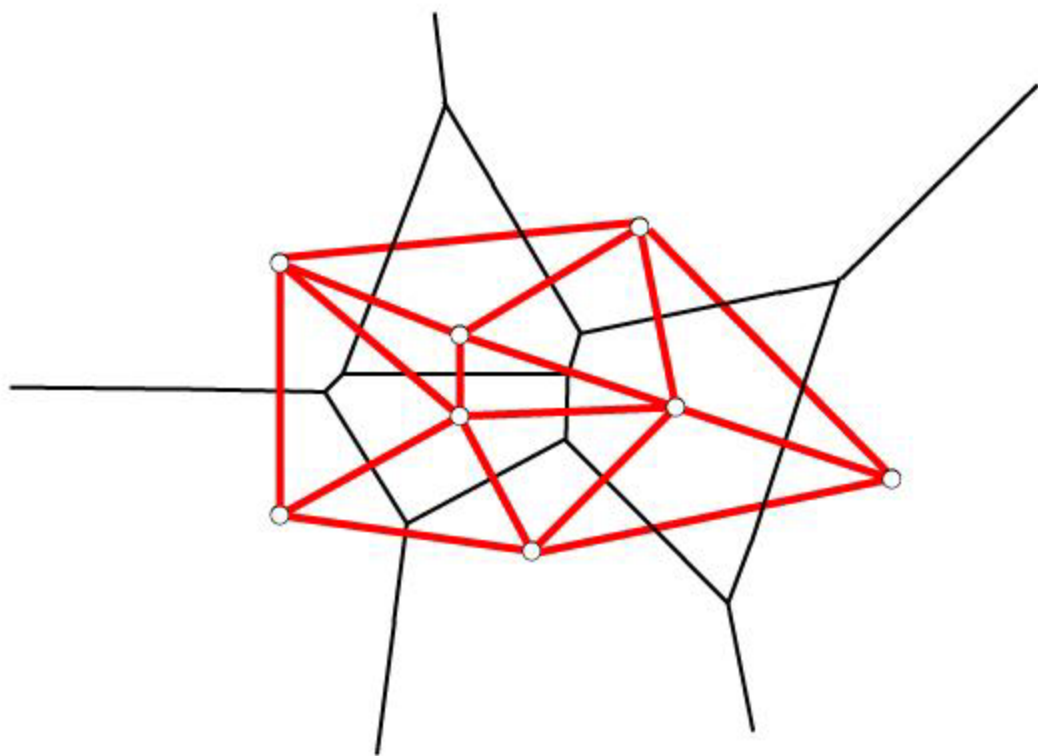
*Decomposition of the*
*Space-filling diagram*

# Computing the Surface Area
# and Volume of a Union of Balls



*Surface Area*

$$A = 4\pi \sum_{i=1}^{N} \rho_i^2 \sigma_i$$
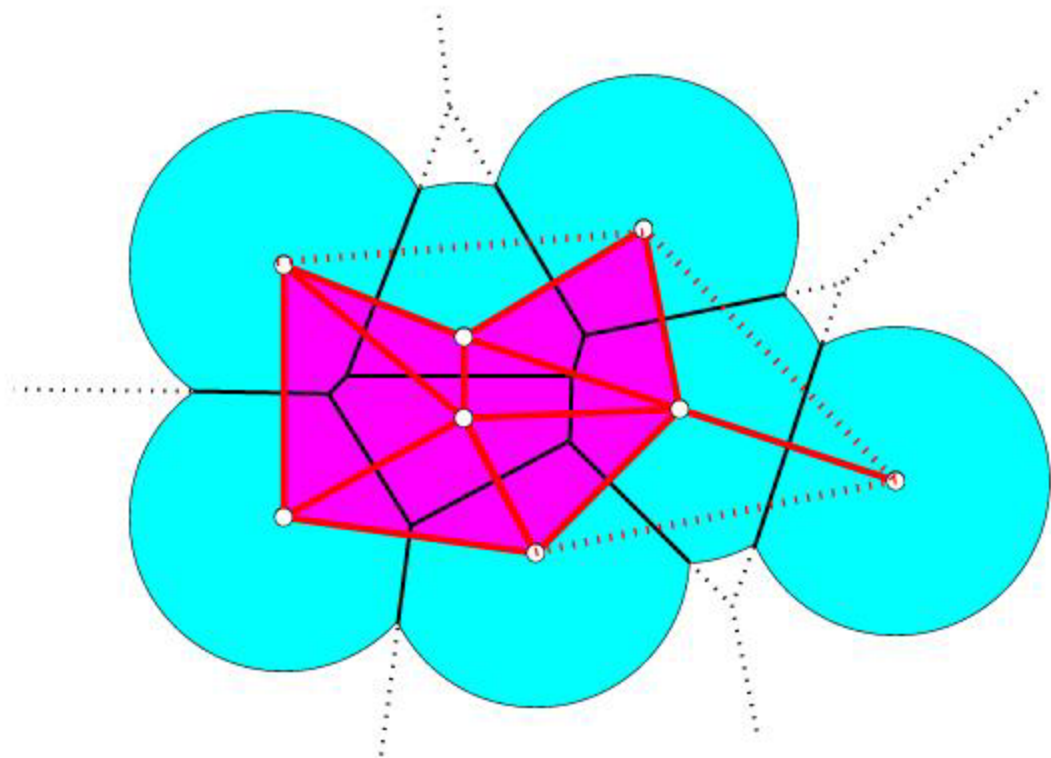
*Volume*

$$V = \frac{4\pi}{3} \sum_{i=1}^{N} \rho_i^3 \beta_i$$

# Computing the Surface Area
# and Volume of a Union of Balls



*The weighted Delaunay triangulation is the dual of the power diagram*
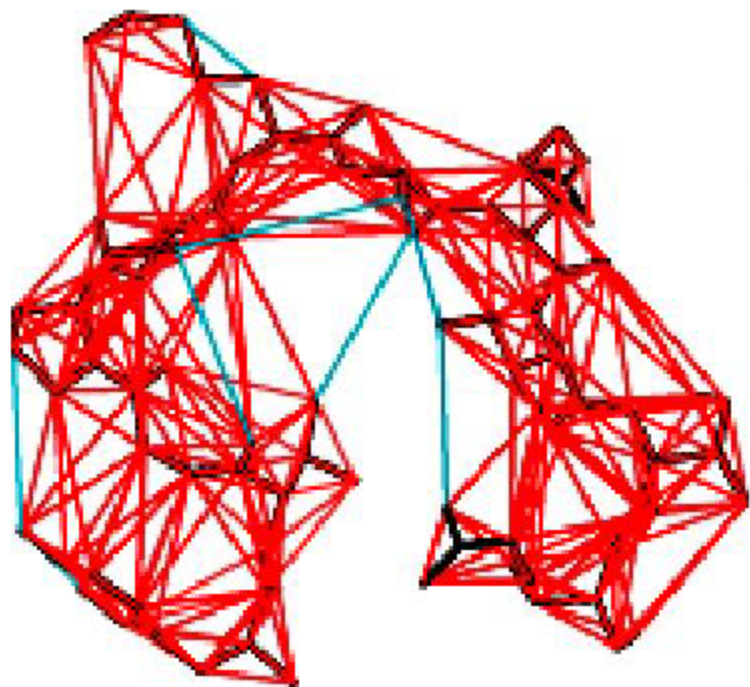
# Computing the Surface Area and Volume of a Union of Balls



*The dual complex K is the dual of the decomposition of the space-filling diagram*

# Computing the Surface Area
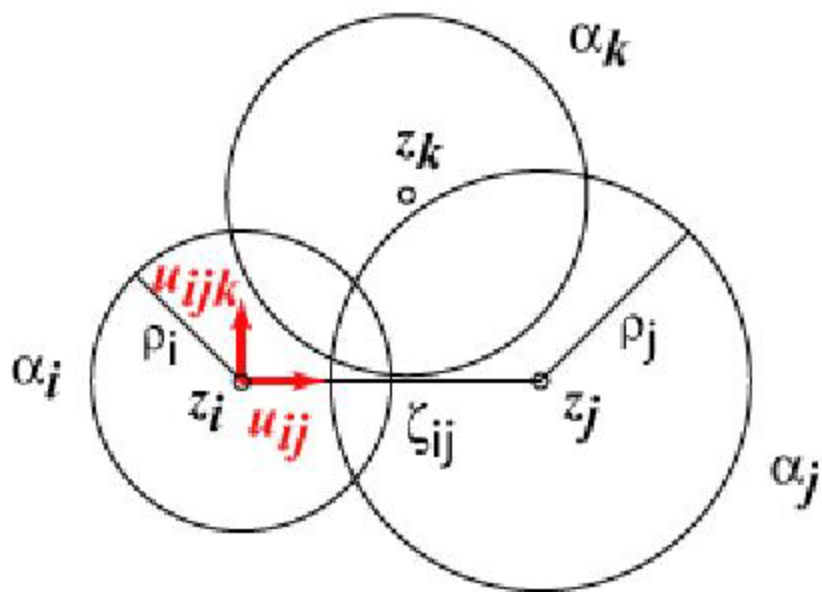# and Volume of a Union of Balls

*Inclusion-exclusion formulas for surface and volume*

$$Area\left(\bigcup B\right) = \sum_{X \in K} (-1)^{\dim X} Area\left(\bigcap X\right)$$

$$Vol\left(\bigcup B\right) = \sum_{X \in K} (-1)^{\dim X} Vol\left(\bigcap X\right)$$
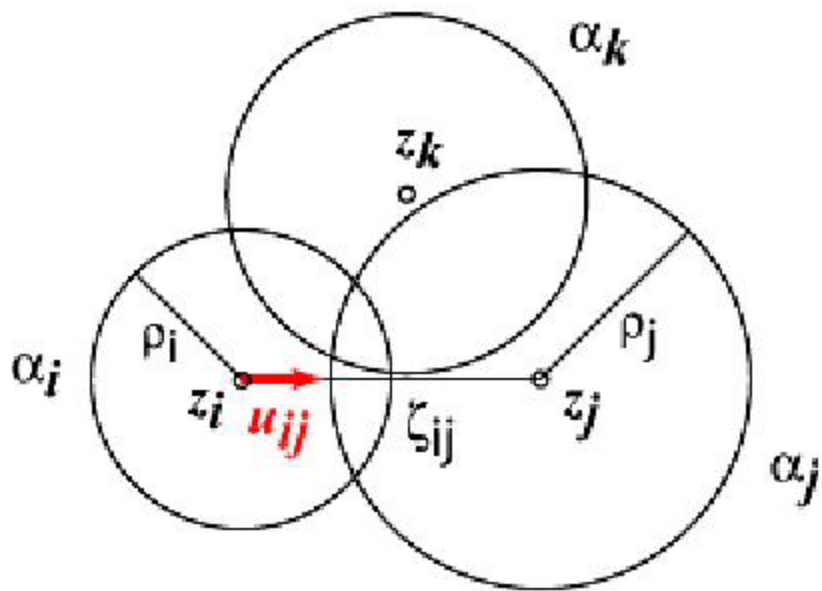
# Weighted Area Derivative



$$\begin{bmatrix} e_{3i-2} \\ e_{3i-1} \\ e_{3i} \end{bmatrix} = \sum_j \left( \sigma_{ij} \bullet e_{ij} + \sum_k \beta_{ijk} \bullet e_{ijk} \right)$$

$$e_{ij} = \pi \left[ \left( \alpha_i \rho_i + \alpha_j \rho_j \right) - \left( \alpha_i \rho_i - \alpha_j \rho_j \right) \frac{\rho_i^2 - \rho_j^2}{\zeta_{ij}^2} \right] u_{ij}$$

$$e_{ijk} = 2\rho_{ijk} \frac{\alpha_i \rho_i - \alpha_j \rho_j}{\zeta_{ij}} u_{ijk}$$

# Weighted Volume Derivative



$$\begin{bmatrix} w_{3i-2} \\ w_{3i-1} \\ w_{3i} \end{bmatrix} = \sum_{j} S_{ij} \left( w_{ij} \bullet u_{ij} + x_{ij} \bullet v_{ij} \right)$$
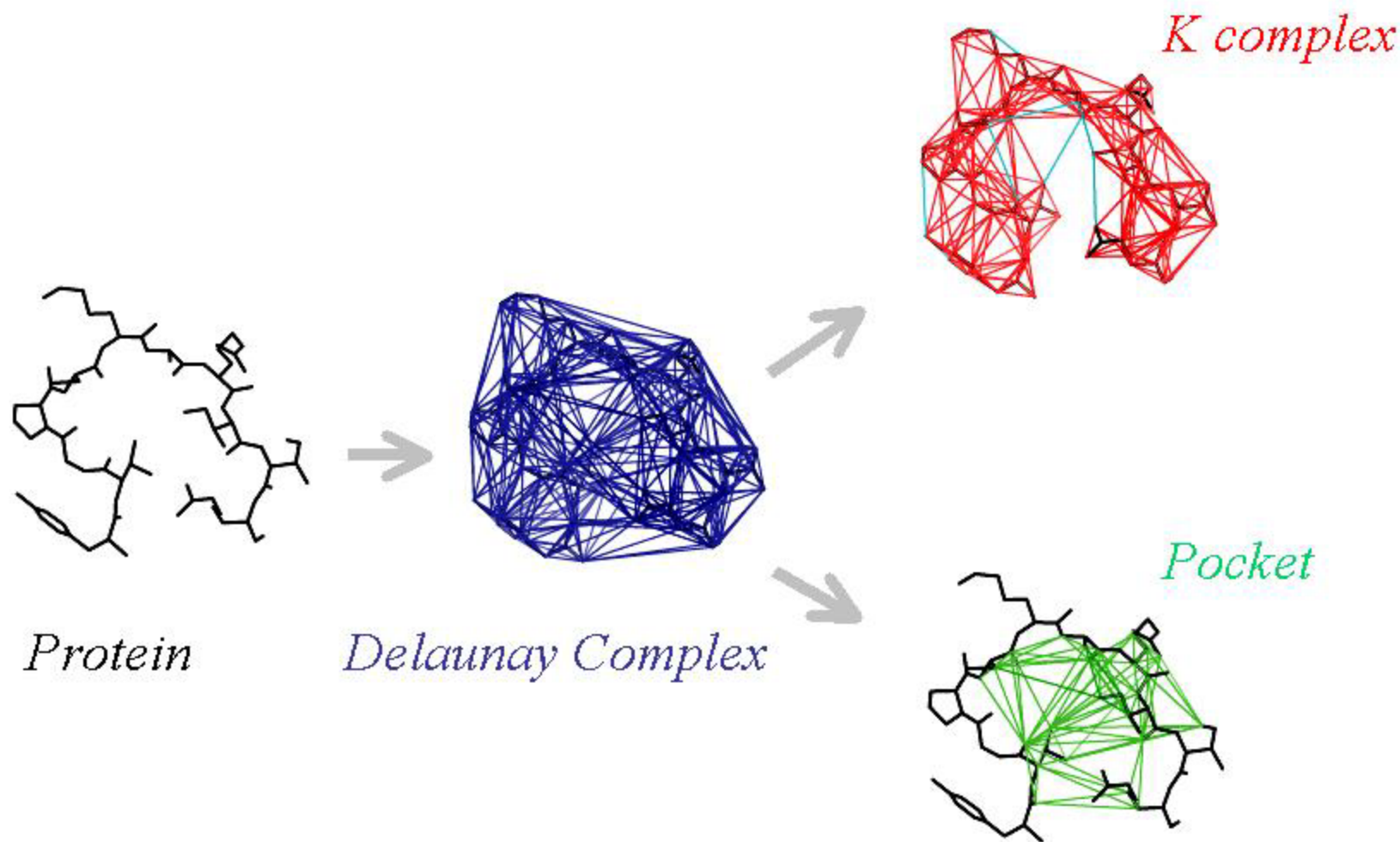
$$w_{ij} = \left[ \frac{\alpha_i + \alpha_j}{2} + \frac{\left(\alpha_i - \alpha_j\right)\left(\rho_i^2 - \rho_j^2\right)}{2\zeta_{ij}^2} \right]$$

$$x_{ij} = \frac{2\left(\alpha_i - \alpha_j\right)}{3\zeta_{ij}}$$

$S_{ij}$: *Area of facet (i,j)*

$V_{ij}$: *Average vector from center to boundary of facet (i,j)*
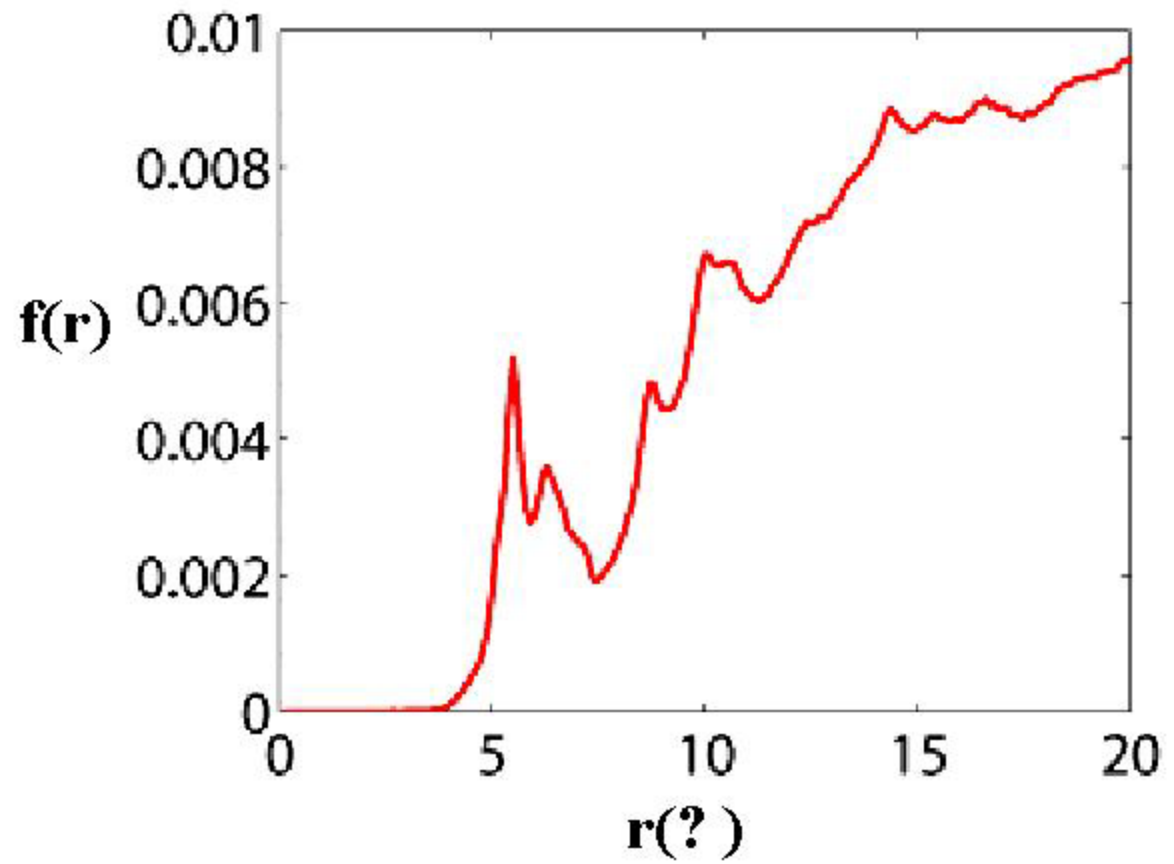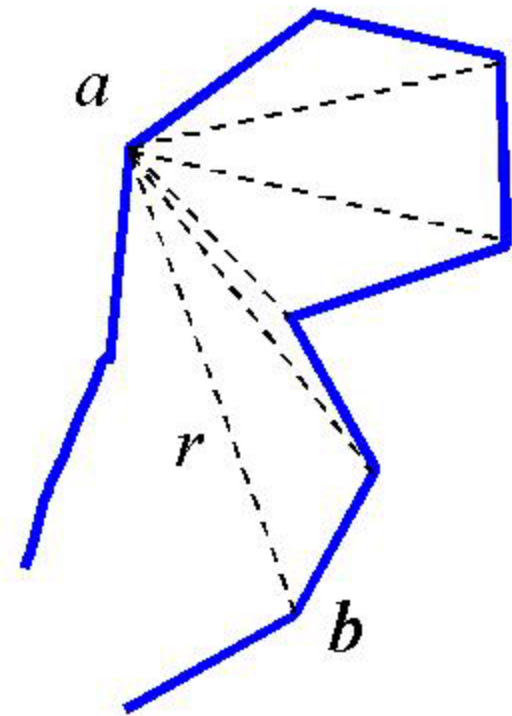
# Computing the Surface Area
# and Volume of a Protein



K complex

Pocket

Protein

*Delaunay Complex*

# Software: ProGeom

|  | Regular Triangulation | Dual Complex | Inclusion-Exclusion |
|---|---|---|---|
| Surface Area | 0.19 | 0.06 | 0.05 |
| Surface Area, Derivatives | 0.19 | 0.06 | 0.08 |
| Volume | 0.19 | 0.06 | 0.18 |
| Volume, Derivatives | 0.19 | 0.06 | 0.29 |

*System: 3740 balls (corresponding to a protein with 492 residues)*
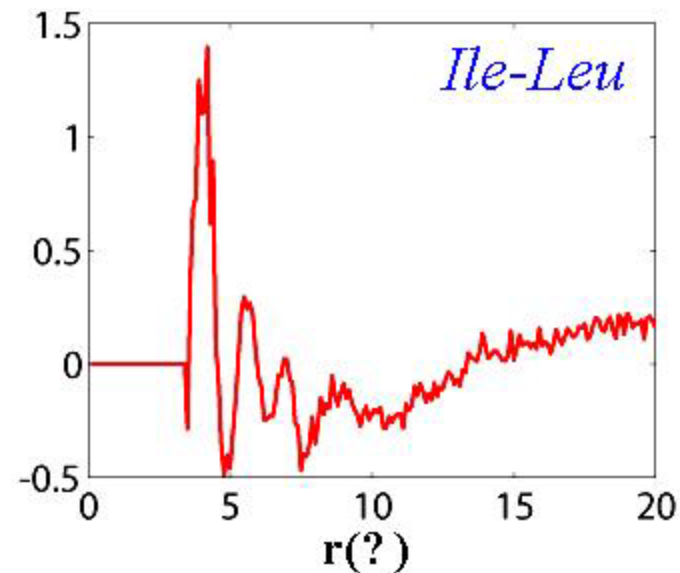*Computing time is seconds, on a Athlon 1.8 Ghz PC computer)*

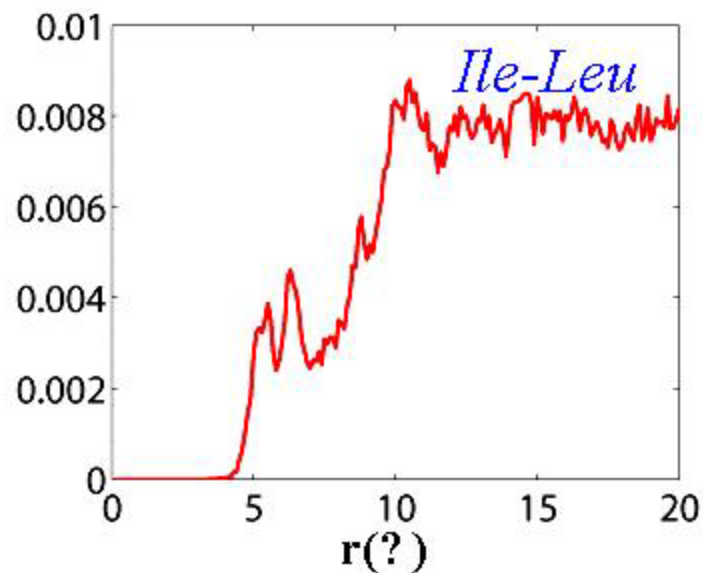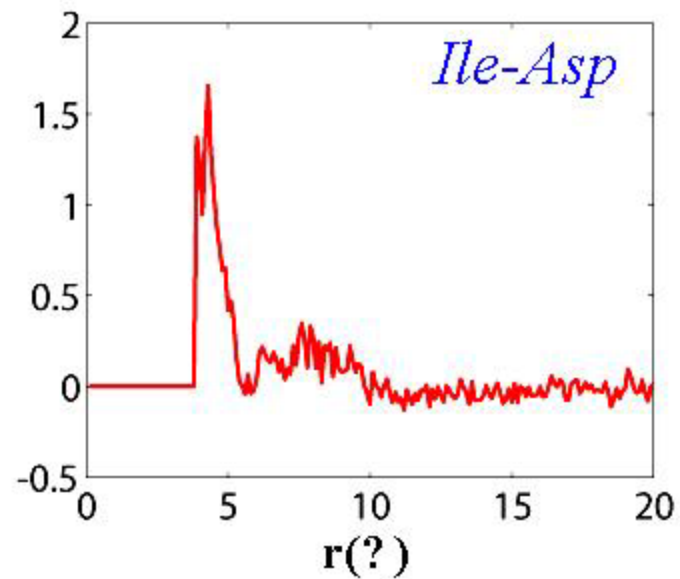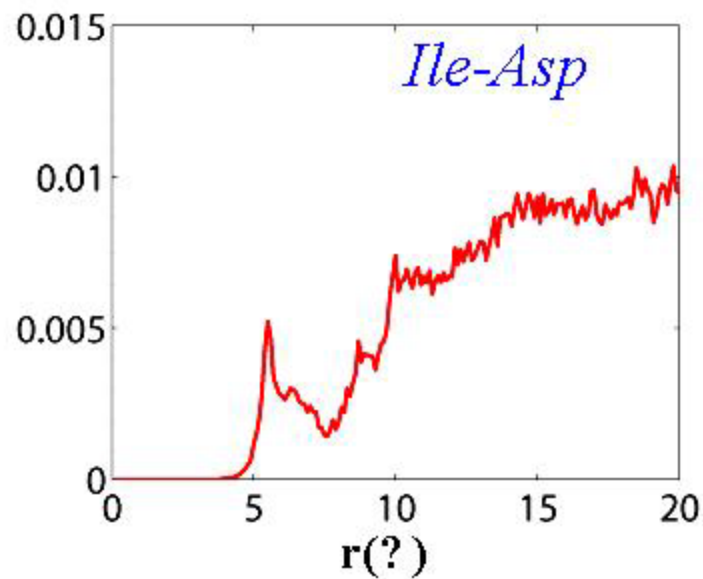*(http://csb.stanford.edu/koehl/ProShape/download.php)*

# Statistical Potentials

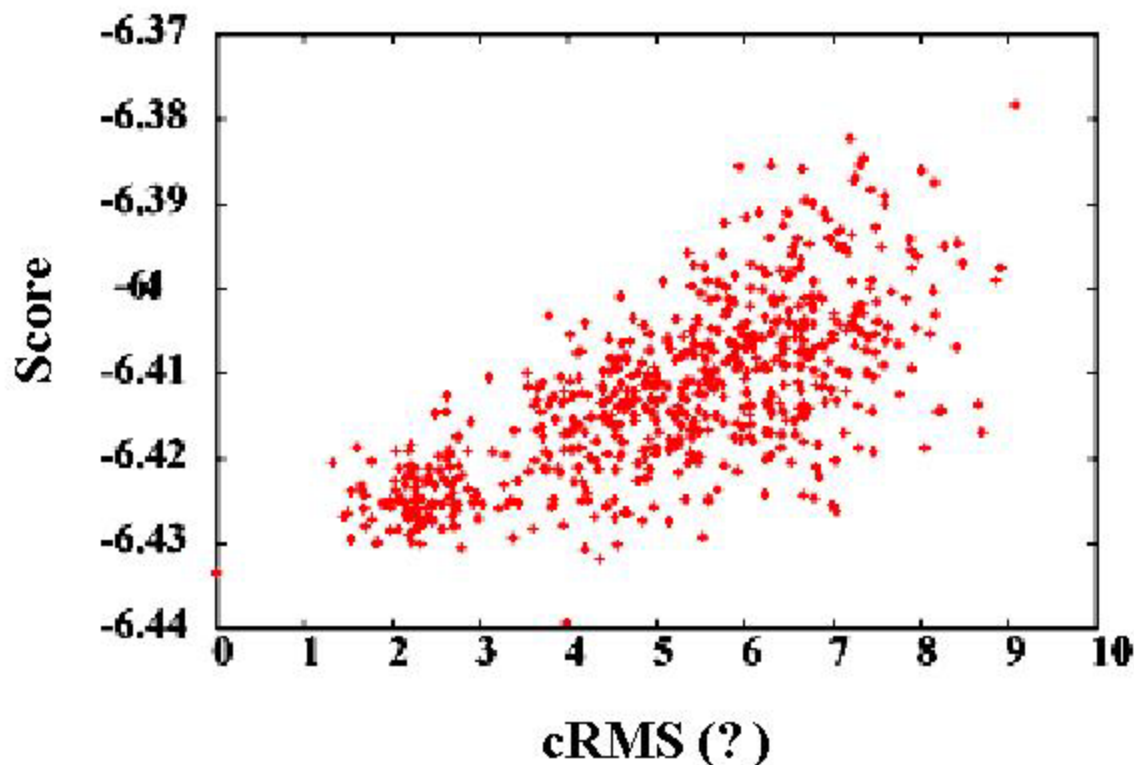$$E(a,b,r) = -\ln\left(\frac{P_{(a,b)}(r)}{P(r)}\right)$$
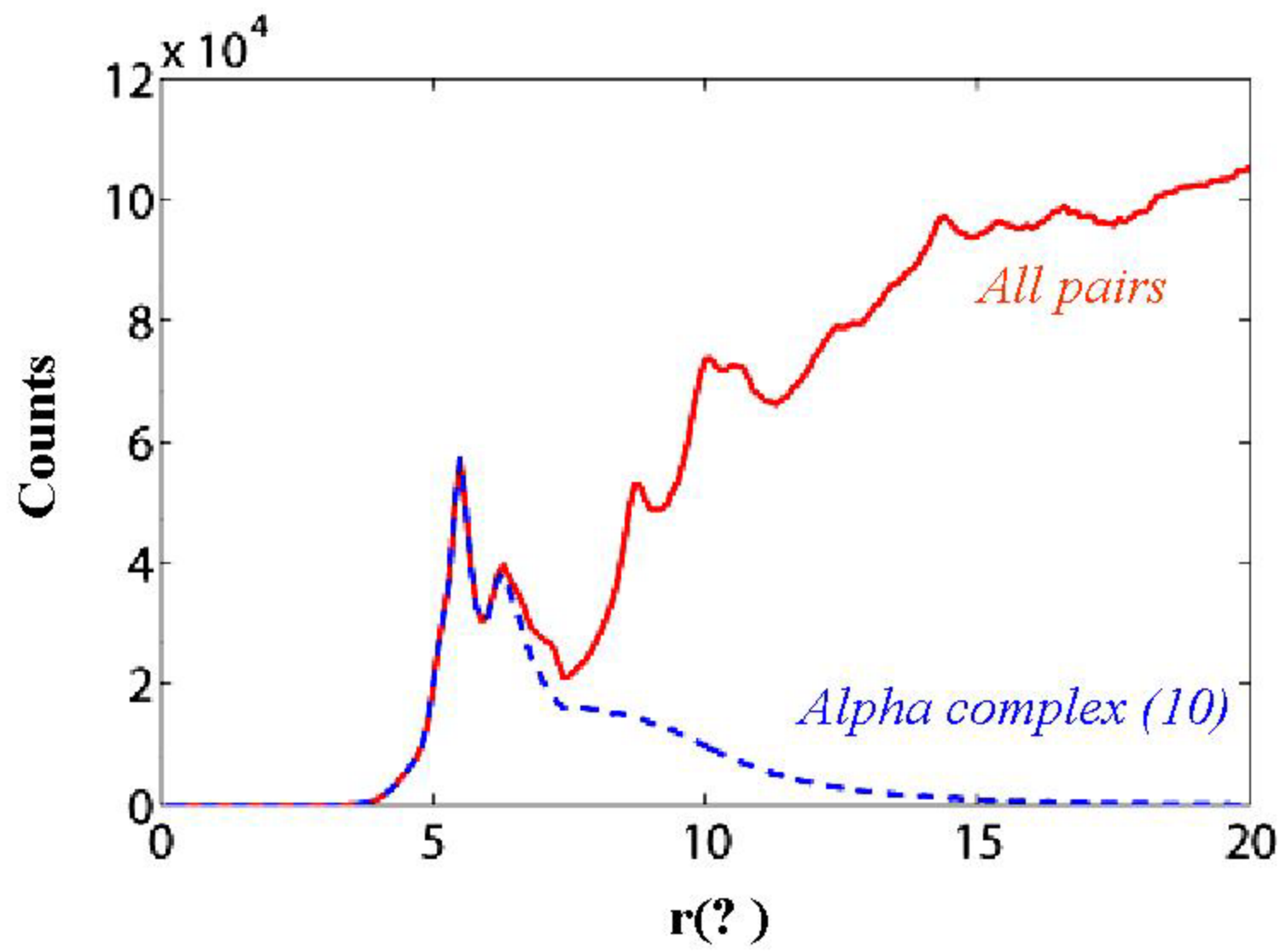
*Counts*

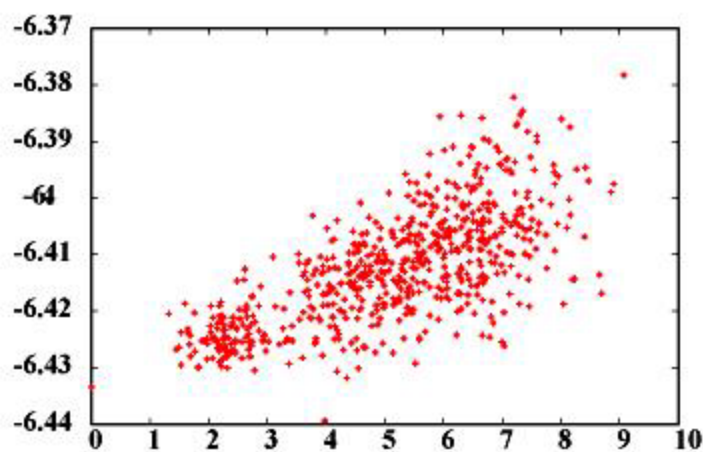*Energy*

# The Decoy Game
# Finding near native conformations
## 1CTF



$$E = \sum_{i<j} E(i,j) = -\sum_{i<j} \ln\left(\frac{P(a_i, a_j, r_{ij})}{P(r_{ij})}\right)$$

Geometric Filtering of the Residue Pairs

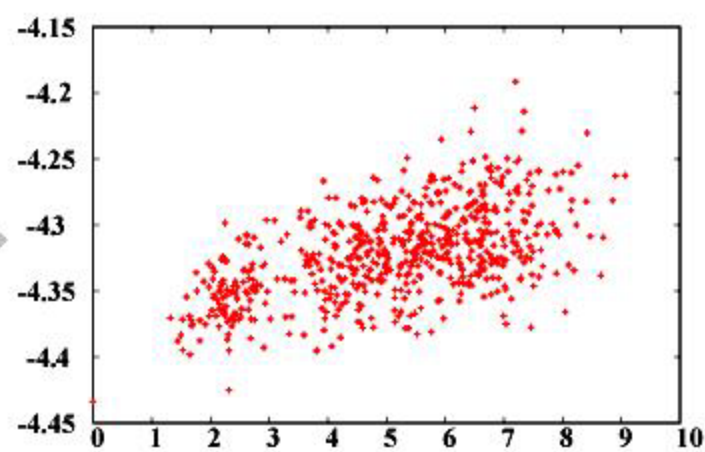All pairs

Alpha complex (10)

# Filtering Does Not Reduce Performance of PMF

# Outline

•Introduction

What is a Protein?

•Protein Energy Functions

Computational Geometry Tools

•Classifying Proteins

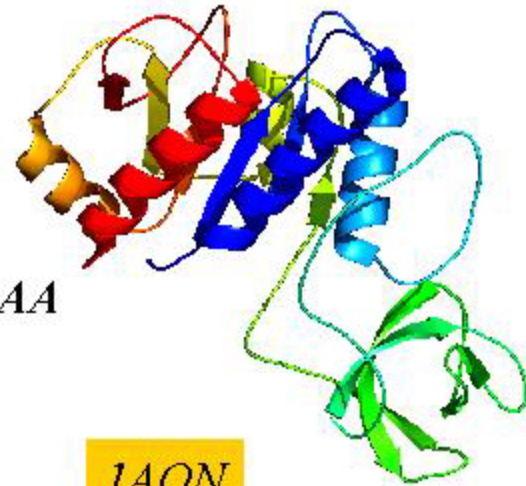The Shapes of Protein Structures

# Protein Structure Space

**1CTF**

68 AA

**1TIM**

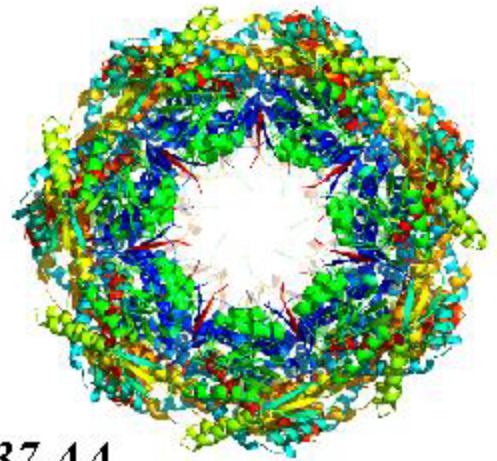247 AA

**1K3R**

268 AA

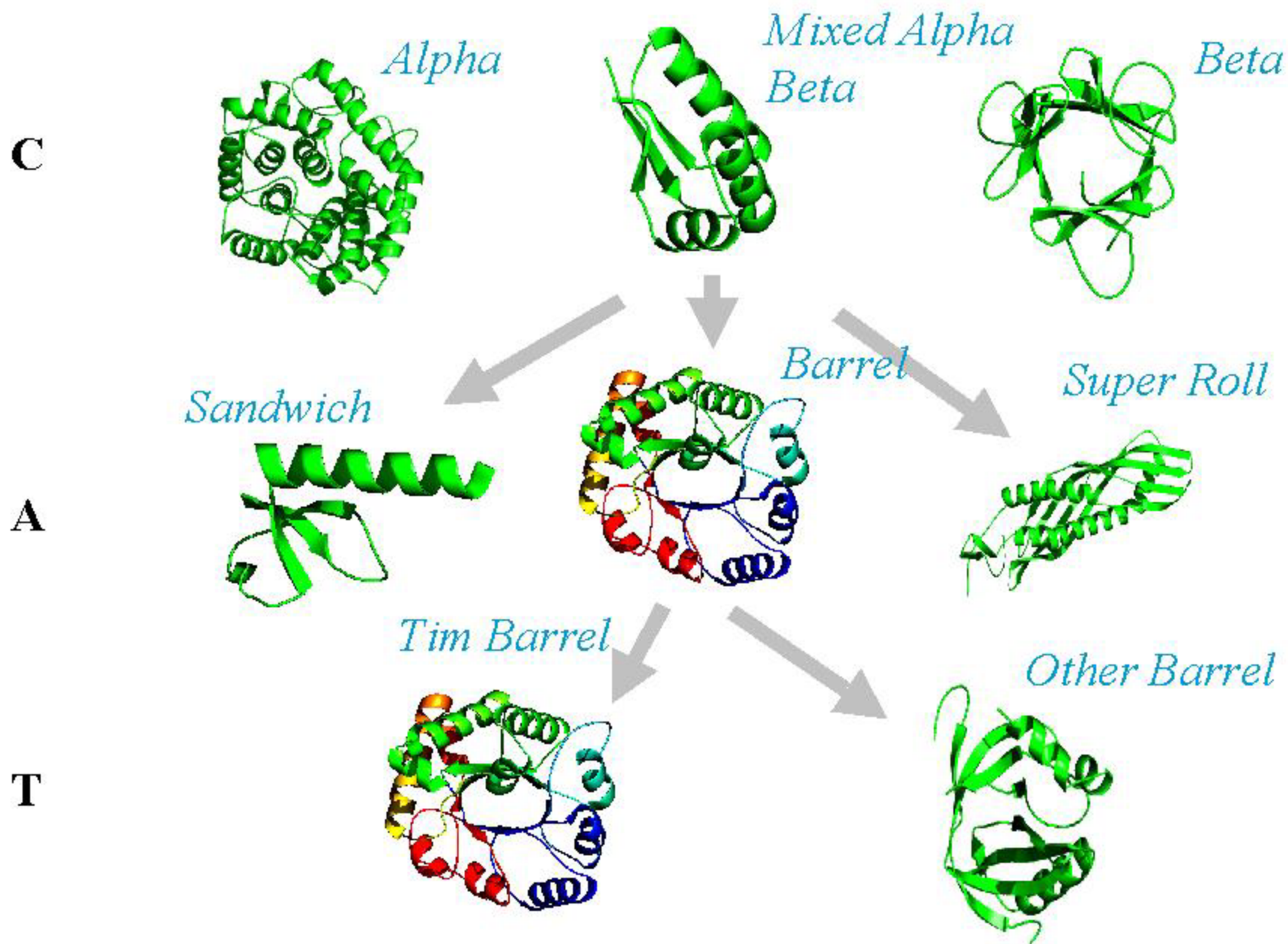**1A1O**

384 AA

**1NIK**

4504 AA

**1AON**

8337 AA

# Classification of Protein Structure: CATH

# Protein Structure Similarity

*Test set*

    2,930 proteins out of 23,000 proteins in PDB

    No sequence similarity (Fasta E-value < e-4)

*Reference structural similarity defined from CATH*

    769 folds

    104,000 pairs of similar structures out of 4,600,000 pairs

# Projecting Protein Structure Space

$$D = \begin{bmatrix} 0 & ... & d_{1N} \\ ... & 0 & ... \\ d_{N1} & ... & 0 \end{bmatrix} \quad \rightarrow \quad G = X^T X \quad \rightarrow \quad X$$

*Distance Matrix*       *Metric Matrix*       *Points in Space*

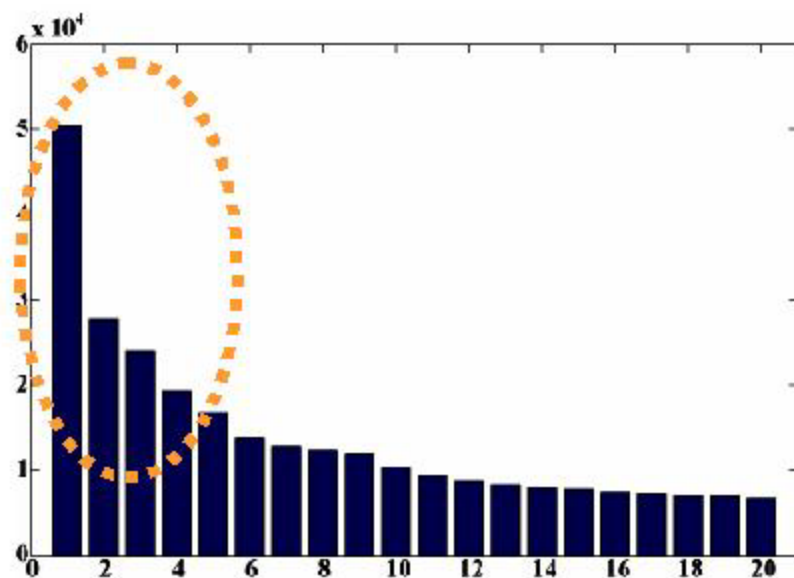# Projecting Protein Structure Space

# Protein Structure Classes
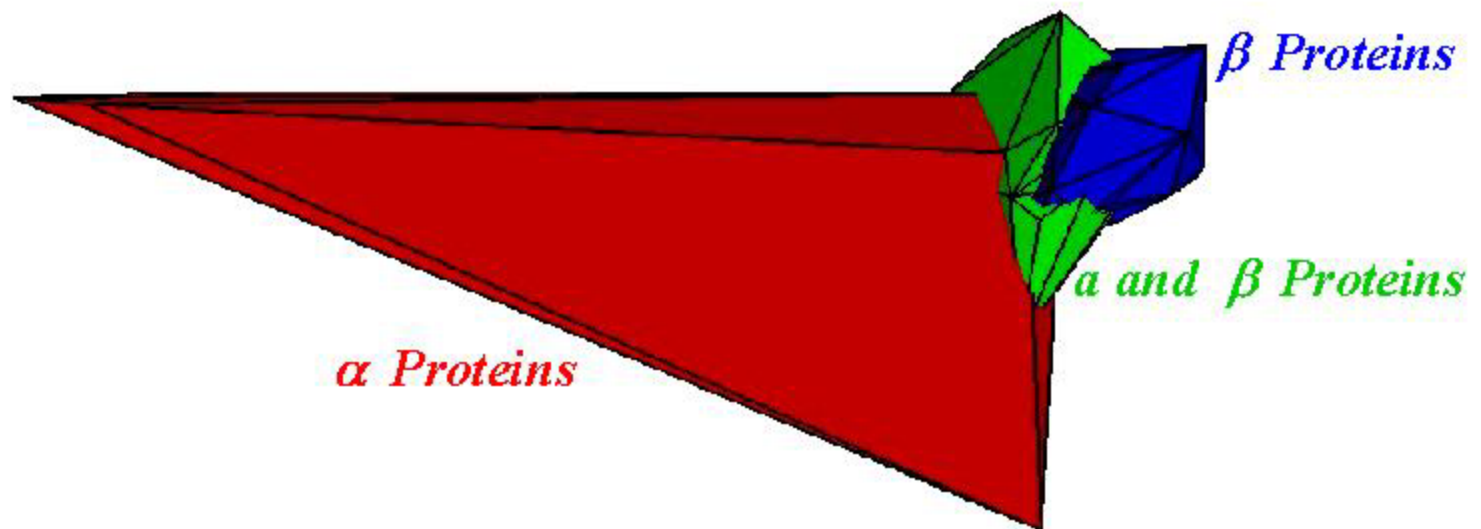
*Measure of Structure Similarity:*
*cRMS after Optimal Superposition*
*(Structal)*

*Eigenvalues of the Metric Matrix:*

# A Picture of the Protein Structure Space



*β Proteins*

*a and β Proteins*

*α Proteins*

# A Picture of the Protein Structure Space



1repC2

1bdo00

1a81G2

2bi6H0

β Proteins

α and β Proteins

α Proteins

1sfcK0

# cRMS is not a Metric



cRMS = 2.8 ?

cRMS = 2.85 ?

# Protein Fold Space
# ROC Analysis
## *(Receiver Operating Characteristic)*

# Protein Fold Space
# ROC Analysis
*(Receiver Operating Characteristic)*

*True positives*

pairs of proteins that belong to the same T class of CATH

*True negatives*

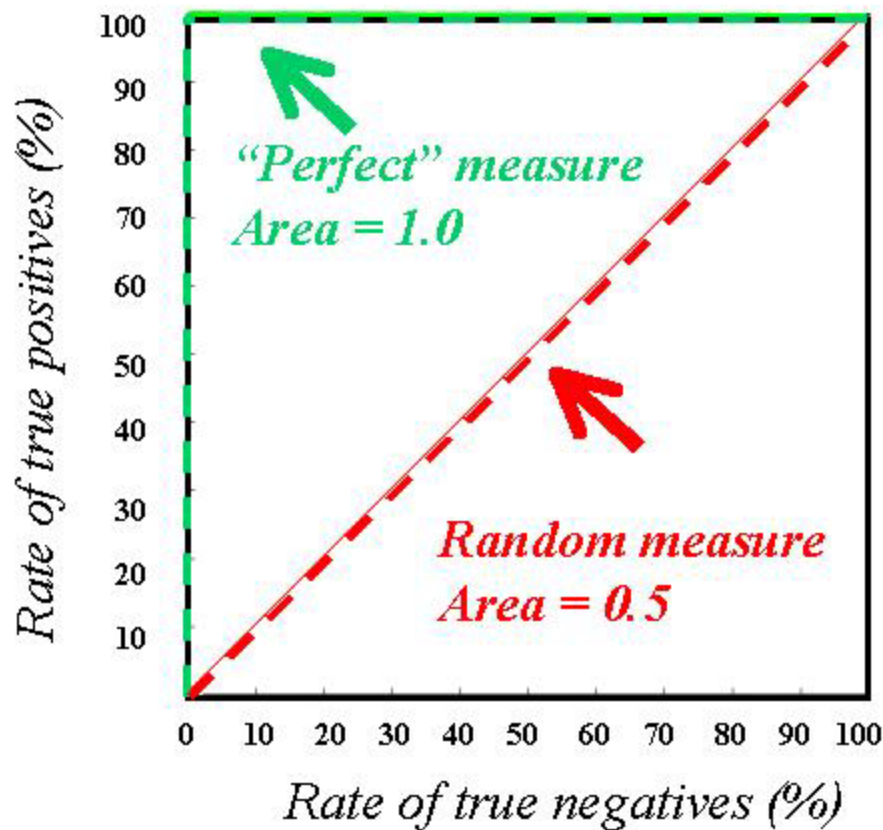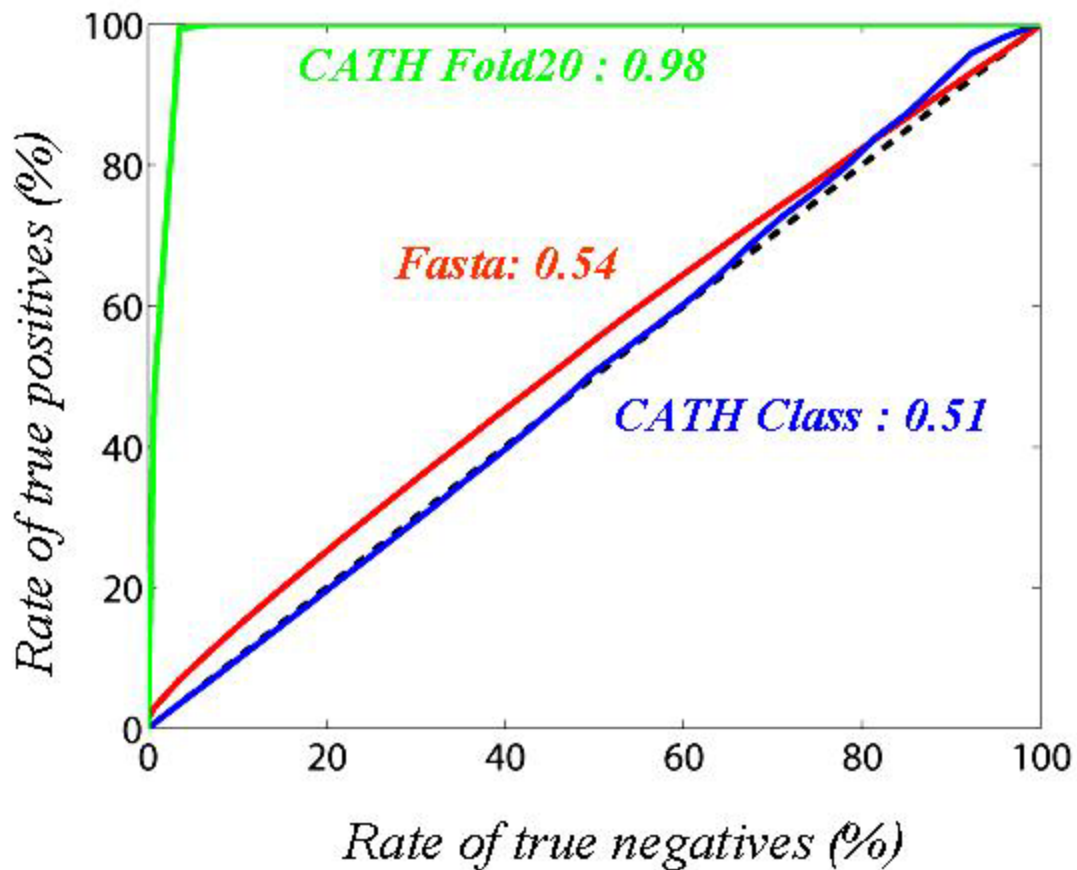pairs of proteins that belong to the same C class, but not the same T class.

# Protein Fold Space



CATH Fold20 : 0.98

Fasta: 0.54

CATH Class : 0.51

Rate of true positives (%)

Rate of true negatives (%)
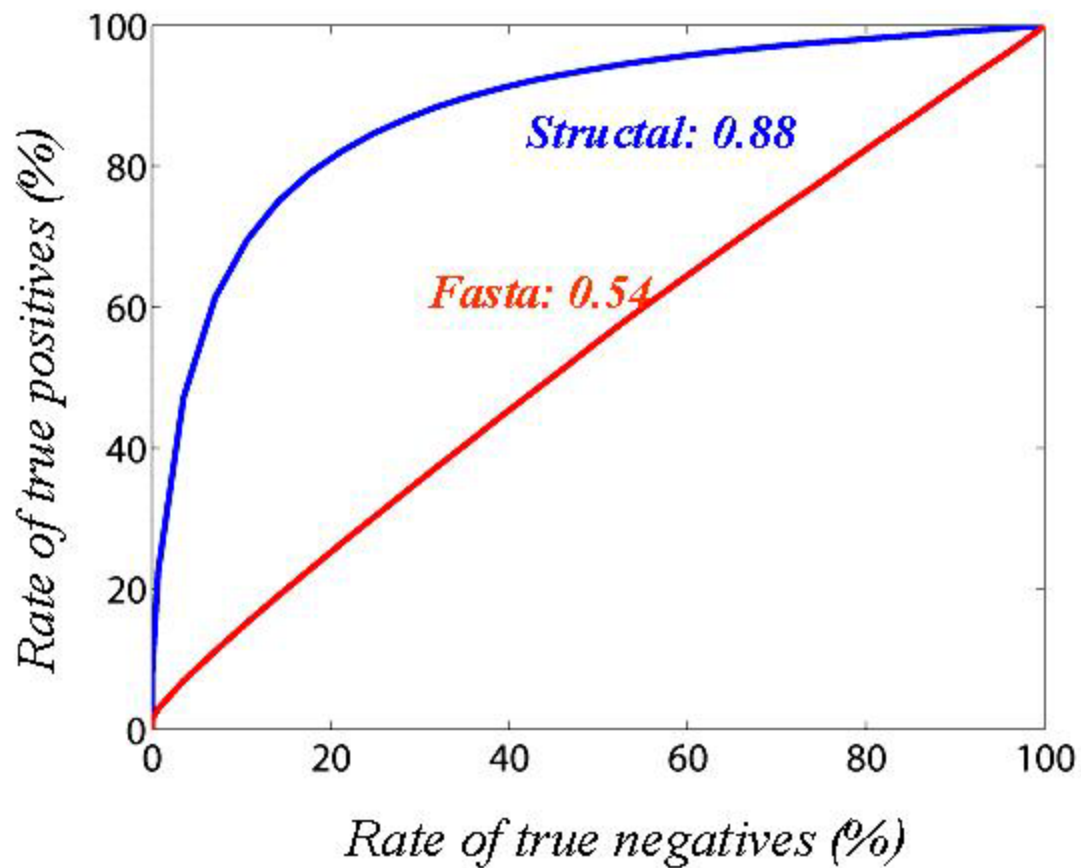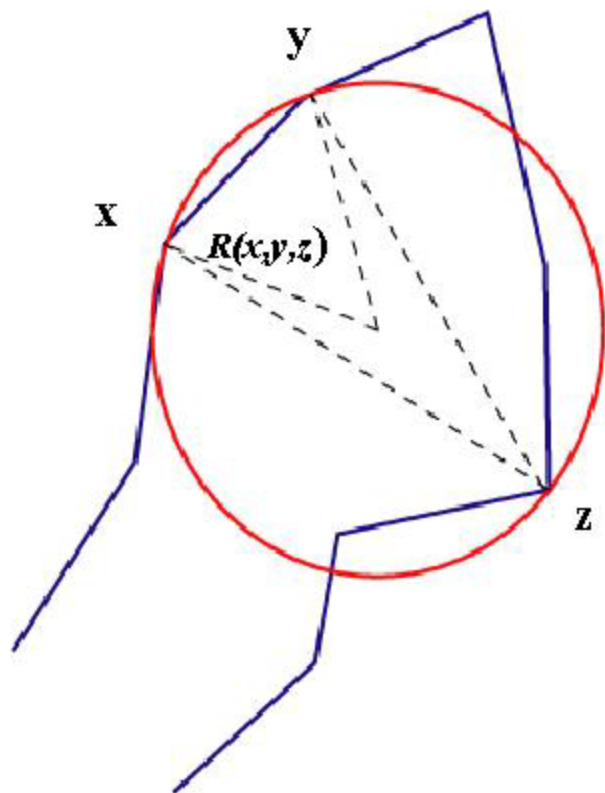
Fold20: first 20 coordinates derived from the CATH fold matrix

CATH class: first 3 coordinates derived from the CATH class matrix

# Protein Fold Space

# Protein Structure Features



$$R(\boldsymbol{x},\boldsymbol{y},\boldsymbol{z}) = \frac{d(\boldsymbol{x},\boldsymbol{y})}{2\left|\sin\left(\hat{\boldsymbol{z}}\right)\right|}$$

*Global radius of curvature:*

$$\rho(\boldsymbol{x}) = \min_{(\boldsymbol{y},\boldsymbol{z})}\left\{R(\boldsymbol{x},\boldsymbol{y},\boldsymbol{z})\right\}$$

*Thickness:*

$$\Delta = \min_{\boldsymbol{x}}\left\{\rho(\boldsymbol{x})\right\}$$

*(Gonzalez & Maddocks, PNAS, 1999, 96:4769)*

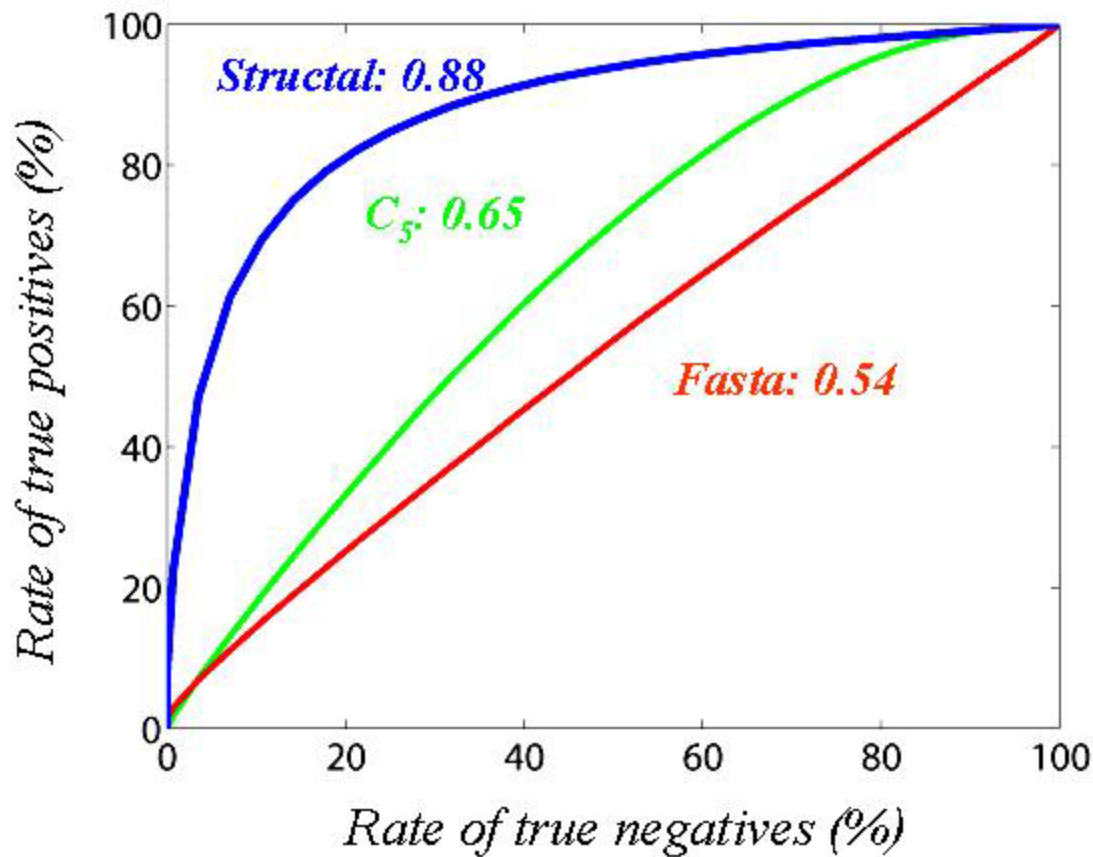# Thickness of a protein structure



$\Delta = 2.60$ ?

# Curvature Feature Vector

$$U_p = \left( \iiint \frac{1}{R(x,y,z)^p} dC_x dC_y dC_z \right)^{1/p}$$

$$C_5 = \begin{bmatrix} U_1 & U_2 & U_3 & U_4 & U_5 \end{bmatrix}$$
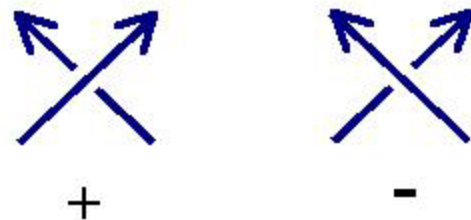
# Performance of the Curvature Feature Vector



*Curvature vector performs better than fasta.*

*Needs more features to match Structal.*
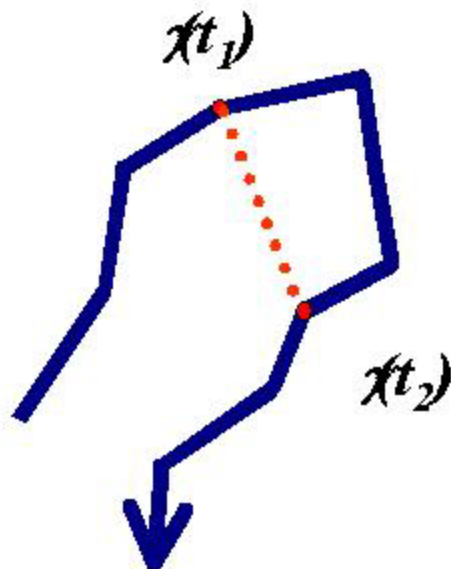
# Protein Structure Features: Writhing

**Sign of Crossing**

$+$      $-$

$\gamma(t_1)$

$\gamma(t_2)$

**Writhing Number**

$$Wr_1 = \frac{1}{4\pi} \iint_{\Delta^2} \omega(t_1, t_2)\, dt_1 dt_2$$

$$\omega(t_1, t_2) = \frac{\det\left(\gamma'(t_1), \gamma(t_1) - \gamma(t_2), \gamma'(t_1)\right)}{\left|\gamma(t_1) - \gamma(t_2)\right|^3}$$

**Writhe Feature Vector for Each Protein**

$$W_{10} = \begin{bmatrix} Wr_1 & |Wr_1| & Wr_{12} & |Wr_{12}| & \end{bmatrix}$$

# Protein Structure Features: Writhing

# Collaborators

- Leo Guibas
  Stanford University

- Herbert Edelsbrunner
  Duke University

- Michael Levitt
  Stanford University

- Afra Zomorodian
  Stanford University

- Rachel Kolodny
  Stanford University

# Thank You