



UNITED NATIONS EDUCATIONAL, SCIENTIFIC AND CULTURAL ORGANIZATION  
INTERNATIONAL ATOMIC ENERGY AGENCY  
**INTERNATIONAL CENTRE FOR THEORETICAL PHYSICS**  
I.C.T.P., P.O. BOX 586, 34100 TRIESTE, ITALY, CABLE: CENTRATOM TRIESTE



**SMR.961 - 14**

**WORKSHOP ON:  
PROTEINS, MEMBRANES and their INTERACTIONS**

**22 JULY - 2 AUGUST 1996**

---

***"Principles of protein structure"***

**Marc DELARUE  
Institut Pasteur  
25-28 Rue de Dr. Roux  
75015 Paris  
FRANCE**

---

*These are preliminary lecture notes, intended only for distribution to participants.*

# Principles of protein structure

## Overview.

1st course: Learn how to generate protein chains with correct

- stereochemistry
- shape
- distribution of  $(\phi, \psi)$
- self-avoiding.

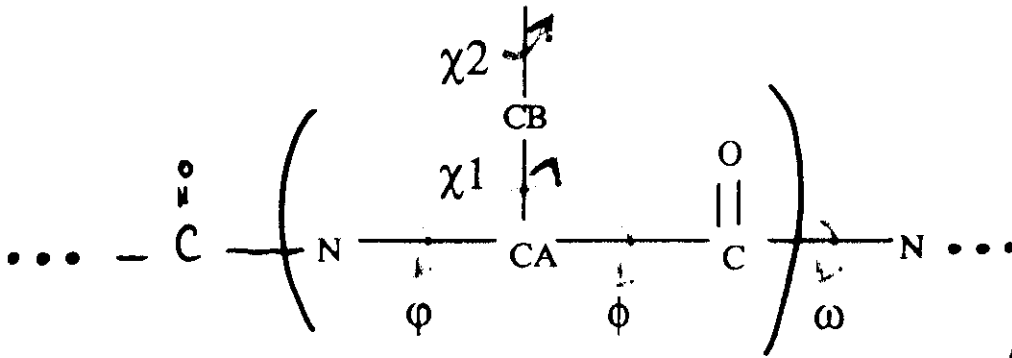
2nd course: Learn how to generate side chains with correct packing - Mean Field theory -

What is the potential energy to be minimized - Full atom models -

3rd course: The so called "inverse protein folding problem" - Sequence space available to a given fold.

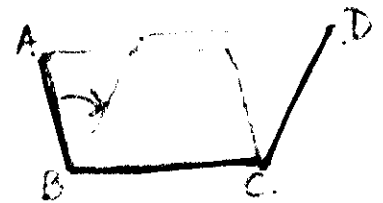
What is the amino acid - amino acid interaction table to be used - How to optimize it.

# PROTEIN CONFORMATION



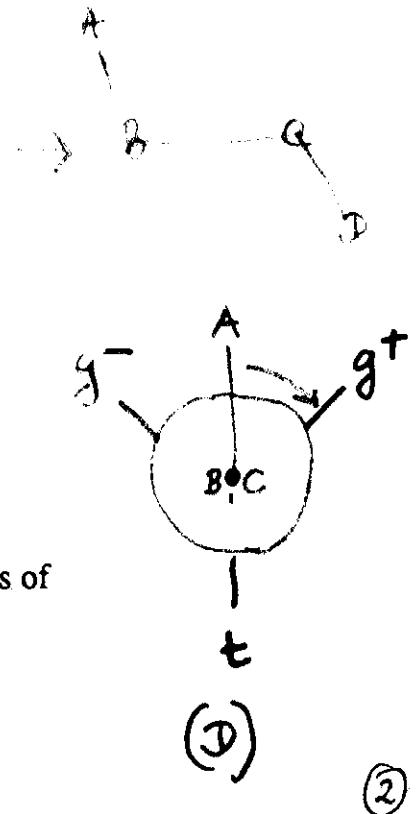
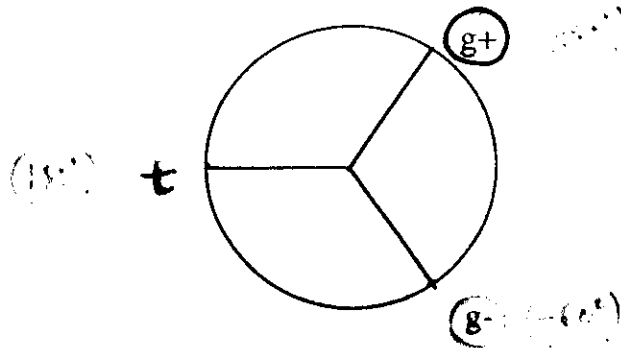
## A) BACKBONE

- 1) defined by 3 dihedral angles :  $\phi$ ,  $\phi$ , and  $\omega$
- 2) Planarity of the peptide bond :  $\omega = 180^\circ$  or  $\omega = 0^\circ$
- 3) For most residues, the  $(\phi, \phi)$  conformational space is limited (Ramachandran plot)



## B) SIDE-CHAINS

- 1) Defined by  $\chi$  dihedral angles
- 3) 3 possible conformations for each  $\chi$  angles

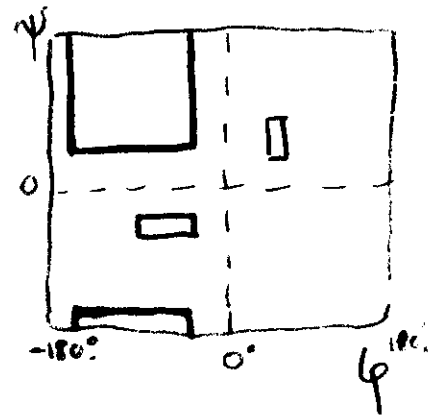


- 3) Side-chains are observed to cluster in limited sets of conformations, called rotamers



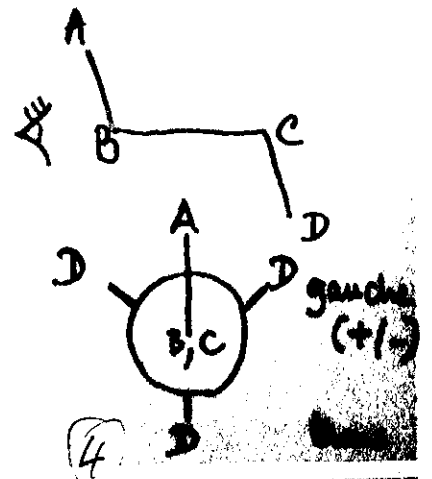
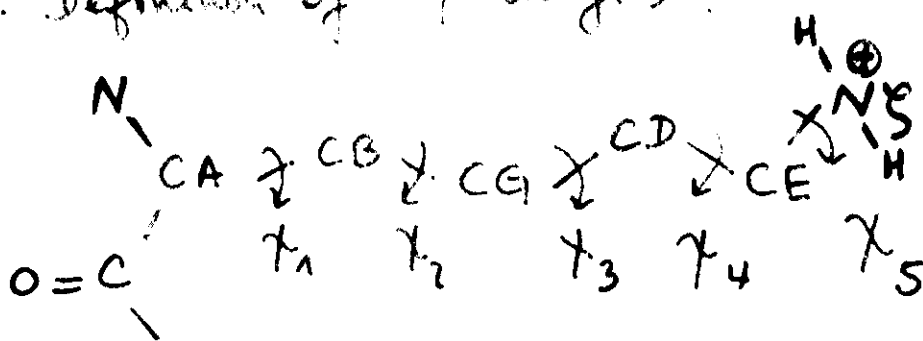
# The side chains

- small : Gly, Ala, Pro (cyclic)  $\Rightarrow \phi = -60^\circ$
- hydrophobic : Val, Ile, Leu, Met, Cys  
(+ sulfur atoms)
- Aromatic : Phe, Trp, Tyr
- Charged  $\ominus$  : Asp, Glu  
 $\oplus$  : Lys, Arg, His
- Amide groups : Asn, Gln
- Small & polar : Ser, Thr

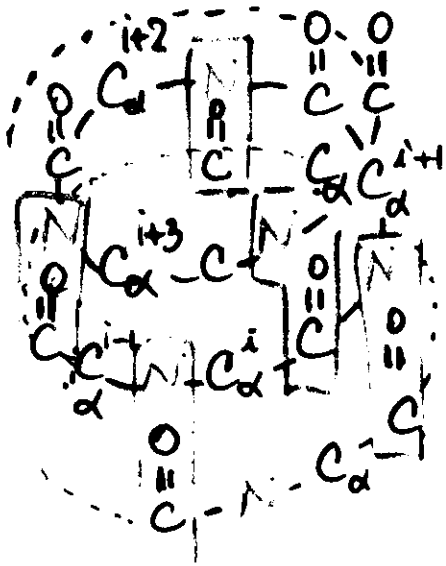


- All of them (except for Gly) have a  $C_\beta$  (at least one  $\chi_1$  angle). This means only some regions of the  $(\psi, \phi)$  diagram (the Ramachandran plot) are accessible because of steric clashes between  $C_\beta$  &  $C=O$ .

- Definition of  $\chi$  angles :



# Secondary structure elements



## • α-helix - Pauling (1953)

1.5 Å / residue

3.6 Å residues per turn

5.4 Å / turn

pitch = 5.4 Å w/o S.C.

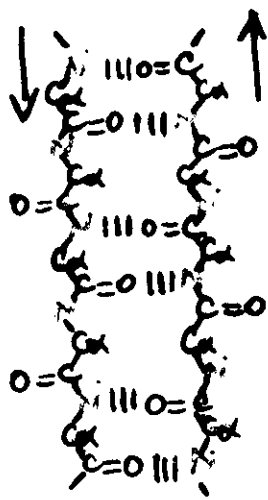
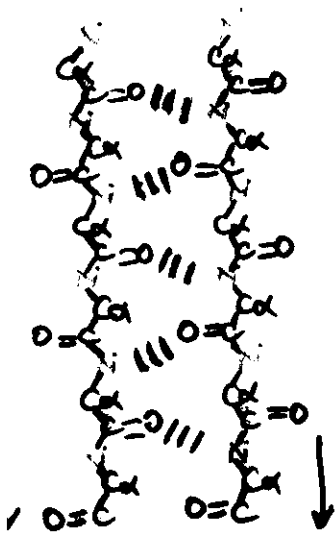
all peptide dipole moments are aligned.

## • β-sheets

3.5 Å / residue

Parallel

Antiparallel



side chains pointing alternately up and down (above & below the plane of the sheet).

twist of  $\beta$  sheets difficult to predict.

# Number of variables for a protein (Nseq)

- suppose  $N_{\text{seq}} = 250 \text{ aa} \Rightarrow \text{MW} \approx 28 \text{ kD}$   
(110 Da / a.a.)

- There are about 8 atoms / aa, on the average

• 4 (N, C $\alpha$ , C, O) for main chain

• 4 (C $\beta$ , O $\gamma$ , O $\delta_1$ , O $\delta_2$ ) for side chain.

Average weight of an atom is 14

- of conformational coords

$\Rightarrow 250 \times 8 \times 3 \approx \underline{6 \cdot 10^3}$  variables

- of dihedral (internal) coords

$\Rightarrow 250 \times 4 \approx \underline{10^3}$  variables.

The information missing is contained in the stereochemistry

- fixed bond lengths

- nature of peptide bond (planar,  $\omega = 180^\circ$ )

- fixed bond angles: tetrahedral for C  
trigonal for N

- if globular & spherical  $R \approx 24 \text{ \AA} - 25 \text{ \AA}$ .

$V \approx 42,000 \text{ \AA}^3$

## Other simplified representations of proteins

③ CA only (with the constraint  $CA-CA = 3.8 \text{ \AA}$ )

⇒ it is in most cases possible to reconstruct the whole backbone, using fragments of known proteins deposited in the PDB satisfying minimum bound with given CA positions - (bound  $d_{ij} < 1 \text{ \AA}$ ).

④ side chain only (with known CA's).

side chain =  $A_{ij} \times A_{ij}$ .

side chain = only 3-body information

$V_{ij}(s, i) = 1$  iff  $d_{ij} < \text{cutoff}$ .

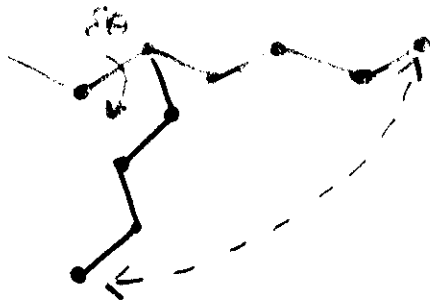
It is that kind of information that is obtained from NMR experiments.

# Simulate the dynamical behaviour of a protein

(A) cartesian coords  $m \frac{d^2 \vec{r}_i}{dt^2} = - \vec{\nabla}_{\vec{r}_i} V$

- Many variables - Add solvent molecules -
- integration time step limited by movement of highest frequency  $\rightarrow$  bond distance oscillations ...

(B) write down eqns of internal coords = pb of minimum consequences of one small perturbation of one angle



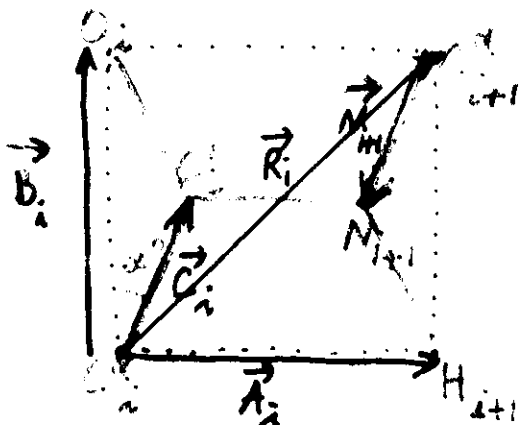
- Solutions
- R. Abagyan (NY) - Lagrangian etc ...
  - AT Brünger - Robotics & Kalman Filters..
  - The so-called "local moves" arrange to make moves that leave endpoints invariant.



# The Knapp algorithm for "local moves"

Knapp, J. Comput. Chem., 14:19-29 (1993).

- edge representation

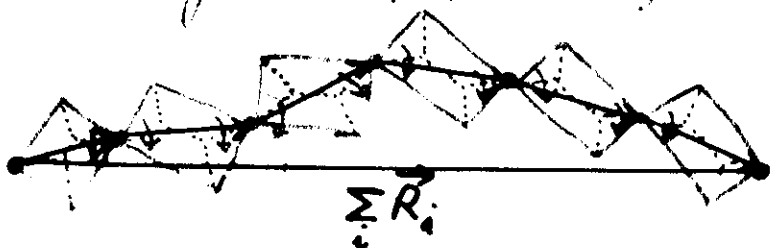


$$L = 1.436 \text{ \AA}$$

$$\text{tetrahedral} \Rightarrow \vec{R}_i \cdot \vec{C}_i = -\frac{1}{3}$$

$$\vec{C}_{i+1} = \vec{A}_i + \vec{B}_i$$

- window algorithm (nwindows)



$$\begin{cases} \vec{R}_{i+1} = \frac{1}{2} \vec{a}_i - \frac{\sqrt{3}}{2} \vec{b}_i \\ \vec{C}_i = \frac{1}{2} \vec{a}_i + \frac{\sqrt{3}}{2} \vec{b}_i \\ \vec{a}_{j+1} = \frac{1}{3} \vec{a}_j - \frac{2\sqrt{2}}{3} \vec{b}_j \\ \vec{b}_{j+1} = \frac{2\sqrt{2}}{3} \vec{a}_j + \frac{1}{3} \vec{b}_j \end{cases}$$

minimum  $\times 2$  - 4 variables can be picked at random  
the 4 remaining angles are given by:

$$\prod_{i=1}^3 \vec{C}_i \cdot \prod_{i=1}^3 \vec{N}_i \cdot \vec{a}_4 = -\frac{1}{3} \quad \text{if } n\text{window} = 4$$

$$\prod_{i=1}^3 \vec{R}_i = \prod_{j=1}^3 \vec{R}_j \quad \text{where } \vec{R}_j = \prod_{k=1}^j \vec{C}_k \vec{N}_k \vec{R}_j^{(0)}$$

- approximate (linear approx.)

$$\vec{C}_j = 1 + \psi_j \vec{C}_j \quad \text{and solve this linear set of eqns}$$

$$\vec{N}_j = 1 + \varphi_j \vec{n}_j$$

run an expansion of some more algorithms

→ Most of the moves will generate steric clashes

1 Expand the molecule  $d \rightarrow d(1+z)$

→ pick rotator at random, pick axis eqns,  
K randomly → rotant → rotant + translation

3 check energy (steric clashes)

20 2x3 many times ...

check, from time to time, if it is possible to

recompact the molecule, if we have reached

a # of iterations, we stop at local minimum.

Monte Carlo simulation of free stereochem.

19.8.1985

- no long-range multipolar instead of additional energy terms. Use interactively large matrix (sparse one) -
- excellent results in display of surfaces

Program available on request, commercially.

How do we approximate the main chain with one with idealized geometry?

Q: minimize some  $\text{error}^2 = \sum_i (x_i^{\text{real}} - x_i^{\text{idealized}})^2$

which is a global minimum.

However, we have to grow the chain, one residue at a time, one by one, i.e. make choices based on local strategy.

It may well be that a choice that is locally optimal at some point is not the best choice for best fitting of the rest of the chain.

Q: grow a population of chains of length  $n \rightarrow$

length  $n+1$  - estimate their Energy -

Replicate the best survivors and give

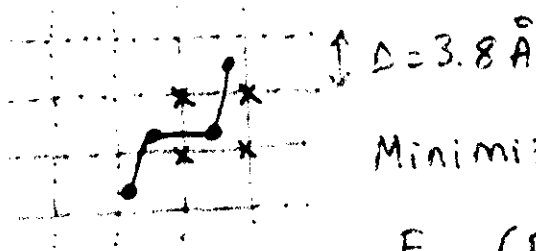
them weight  $\propto \exp(-E_{\text{chain } n+1} / k_B T)$ .

Normalize to keep population constant.

$\rightarrow$  Adjust temperature as  $n$  increases -

$\rightarrow$  See Garel, Orland, Baudin for more detail.

# Projections of protein mainchains on lattices.



Minimize

$$E_{\text{err}}(R_1 \dots R_N) = \sum_{i=1}^{N-1} U_i(R_i - R_{i+1}) \quad \text{Connectivity}$$
$$+ \sum_{i=1}^N f(x_i - R_i) \quad \text{error function}$$

... This can be solved by linear dynamic programming methods (Finkelstein & Reva).

... The final result strongly depends on the orientation of the protein with respect to the lattice.

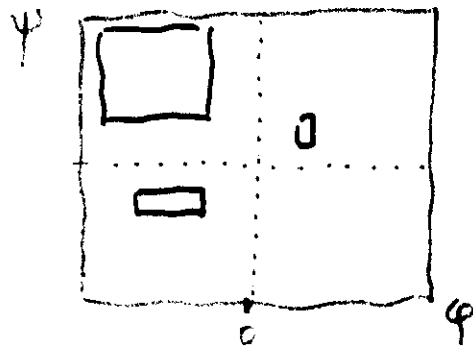
... Adding some other penalty for self avoidance naturally leads to search a solution in the framework of mean-field theory.

... see also lattices used by Shakhovitch & Coll, Dill & Coll, mainly  $3 \times 3 \times 3$  where all configurations are known & can be enumerated.

Ensemble of minimum criteria to be fulfilled by realistic protein main chains -

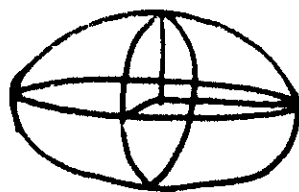
- ① Self avoiding
- ② Correct stereochemistry of peptide bonds: use  $(\phi, \psi)$  angles - Fixed bond lengths & bond angles.
- ③ Correct distribution of  $(\phi, \psi)$  angles in Ramachandran plot. Draw randomly  $(\phi, \psi)$  pair with a bias, given

by a priori information.



- ④ overall shape: Remain in a given ellipsoid  
of volume =  $168 \times N_{res} (\pm 10\%)$   
axial ratios  $e_1 = 0.61 \pm 0.11$   
 $e_2 = 0.78 \pm 0.12$

(Radkowsky, Hao and Schuraga, 1993, PNAS)

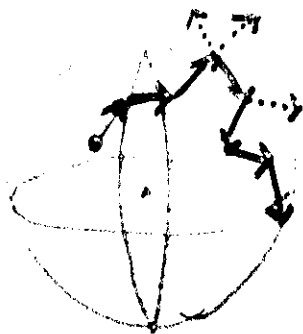


YAM: Use the recursive method of Cohen & Sternberg  
To build plausible protein models.

- ① Randomly place residue #1 inside allowed volume
  - ① Grow the main ~~oxide~~ <sup>chain as long as</sup> all desired criteria are fulfilled - Try again and again until it works -
  - ② Go back 3 or 5 residues upstream if no successful move has been accepted in the last 1000 trials
- Grow 3 residue **at** a time.
  - Abandon process if more than  $\times 10^6$  trials have already been made.
  - Introduce one more level of recursion ...

→ Method is slow but it works:

It can generate 200 structures that are acceptable overnight on a workstation.



# Mean Field Theory : applications to optimization lbs in Mol. Struct. Biology.

(M.D & I. KOEHL).  
I.L. Strasbourg.

What is Mean Field Theory and why do we need it.

Mean Field Theory helps to solve optimization lb. very rapidly. It is derived from Stat. Physics and is especially useful when the Energy contains 2-Body interaction terms. (Ising Model)

$$Z = \sum_{\{S_i\}} \exp -\beta H(\{S_i\}) \quad \text{where } S_i = \pm 1. \\ 2^N \text{ config if } i \in [1, N]$$

and use saddle-point method to evaluate

this  $\Rightarrow$  each spin sees an "average" (MF)

Energy that is determined by a self consistent Equation. (in a lot of them) -

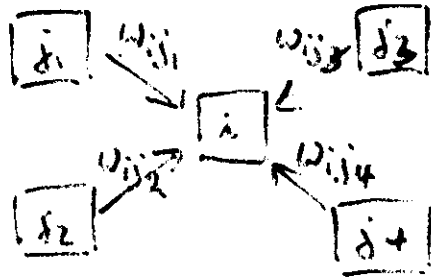
$$\langle S_i \rangle = \tanh(\beta E_{\text{local}}^i) \quad \text{where } E_{\text{local}}^i = \sum_j w_{ij} \langle S_j \rangle$$

if (2 body interactions)  $H = \sum_i \sum_j w_{ij} S_i S_j ; E_{\text{local}}^i = \frac{\partial H}{\partial S_i}$

# Application of Mean Field theory to predict protein side chain conformations.

- The optimization of side chain conf. cannot be treated by exploring exhaustively all possibilities.
- Side chain conformations are first discretized by building a library of rotamers, which are sufficient to describe most side chains of the PDB. (rmsd = 0.86 Å).
- Use mean field theory to explore all possibilities in only one passage.  
All possible rotamers are attached at each position, with a given weight that is updated from one cycle to the other.  
Each copy sees the average of the neighbouring environment (weighted by current weight).
- Only Van der Waals interactions are taken into account.
- update until CV is achieved.

- This method is also well known in the field of Neural Networks, where the response of each neuron is a function of the input of all the interacting neurons



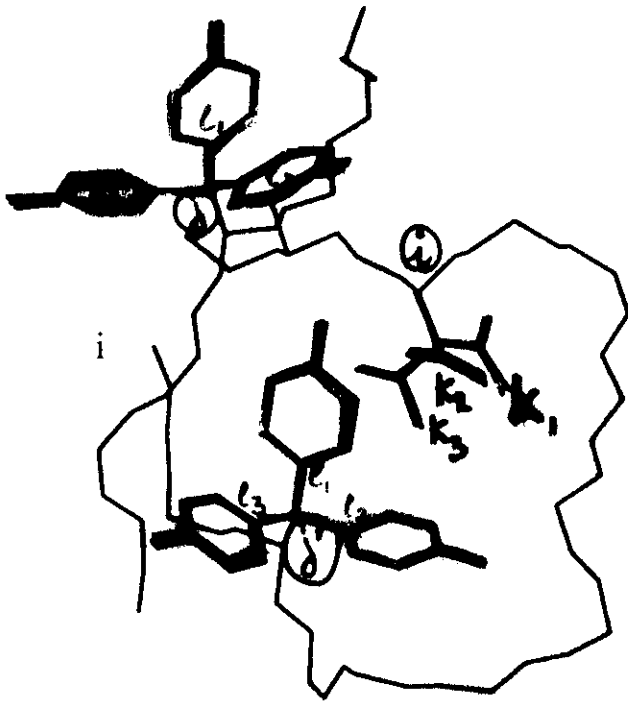
- In the core of a protein, each side chain "talks" to its neighbours and is influenced by them  $\rightarrow$  optimization pb of combinatorial nature

We wish to replace the difficult problem of finding the global minimum ~~logue~~  $E$  (enumerate and evaluate each configuration in turn and exhaustively because the energy landscape is rugged) by another problem, where the variables are the weights of the different possible copies at each position. (Make the pb discrete).

It turns out that the energy landscape is then much easier to explore.

Make the search LINEAR in term of the length of the seq. N seq. (Much more powerful than SA).  
 ⇒ Very rapid and efficient search -

### C) MEAN FIELD APPROXIMATION



Idea: generate multiple copies of the side-chains, which do not see each other, but see all possible copies of their neighbours, affected by their respective probability:  $CM(j,l)$  -

Mean energy of rotamer  $k$  for residue  $i$ :

$$E(i,k) = U(i,k) + U(i,k,B) + \sum_{j \in N_i} \sum_{l=1,2,0} CM(j,l) U(i,k,j,l)$$

position  
 ↓ possible rotamer

Intra	Side-chain - Backbone	Side-chain - Side-chain

- How to derive the local field in the general case:

Minimize  $F = U - TS$  (Finkelstein & Reva).

$$U = \sum_{\text{all possible interact. pairs}} \frac{1}{2} \sum_{\xi} \epsilon_{ij}^{\xi} f(r_i - r_j) p(\eta_i, \xi_j)$$

$$= \sum_i \sum_{j \in N_i} \sum_{\xi_j} \frac{1}{2} \sum_{\xi} \epsilon_{ij}^{\xi} f(r_i - r_j) \underbrace{p(\eta_i) p(\xi_j)}_{\text{MF hypothesis}}$$

$$-TS = +kT \sum_B \sum_{\eta} p(\eta_i) \log p(\eta_i)$$

$$\delta F = \delta U - \delta(TS)$$

$$= \sum_i \sum_{\eta} \delta p(\eta_i) \left[ -kT \log p(\eta_i) + \sum_{j \in N_i} \sum_{\xi} \epsilon_{ij}^{\xi} f(r_i - r_j) p(\xi_j) \right] + \text{const}$$

$$= 0$$

$$\Rightarrow p(\eta_i) = \exp \left[ -\frac{1}{kT} \sum_{j \in N_i} \sum_{\xi} \epsilon_{ij}^{\xi} f(r_i - r_j) p(\xi_j) \right]$$

- Description of the algorithm itself

Start from equiprobable  $p(\eta_i)$  and update until convergence is achieved. Some trick is needed.

$$\frac{\delta p}{\delta t} = k \left( p^{\text{calc}} - p^{\text{old}} \right) = \frac{p^{\text{new}} - p^{\text{old}}}{\delta t} \Rightarrow p = k \delta t p^{\text{calc}} + (1 - k \delta t) p^{\text{old}}$$

$$\Rightarrow p^{\text{new}} = \lambda p^{\text{calc}} + (1 - \lambda) p^{\text{old}}$$

where  $\lambda = k \delta t$   
 $\lambda = 0.5$

- In our implementation of the Method, we never physically construct a stereochemically and "side chain" packing correct chain. It is only at the end, when the algorithm has converged that we decide that, for each position, the rotamer to be constructed is the one that has the highest probability  $P(\eta_i) \rightarrow$  Energy minimization.

- In C. Lee implementation & H. Vasquez, actual physical realizations of a configuration  $\{S_i\}$  are constructed, using the so far derived probability distributions, to calculate  $E_{\text{local}}^i$

$$E_{\text{local}}^i = \left\langle E_{\text{physical}}^i \right\rangle \xrightarrow{\text{new}} P(\eta_i) \text{ by Boltzmann like formulae.}$$

$\alpha \approx 20 \neq$  configurations picked up using  $P^{\text{old}}(\eta_i)$

The system is always "physical".

Method called "Ensemble optimization".

It is always better to conduct the optimization process with a COULATION of CHAINS instead of just one chain.

## SIDE-CHAIN CONFORMATION

### 1) ROTAMER LIBRARY : Tuffery et al (1991)

– a total of 113 rotamers for 19 residues

### 2) THE CONFORMATIONAL MATRIX: CM

For a protein of N residues : CM of size NRES x K

CM(i,j) is the probability that residue i adopts the conformation described  
by rotamer j

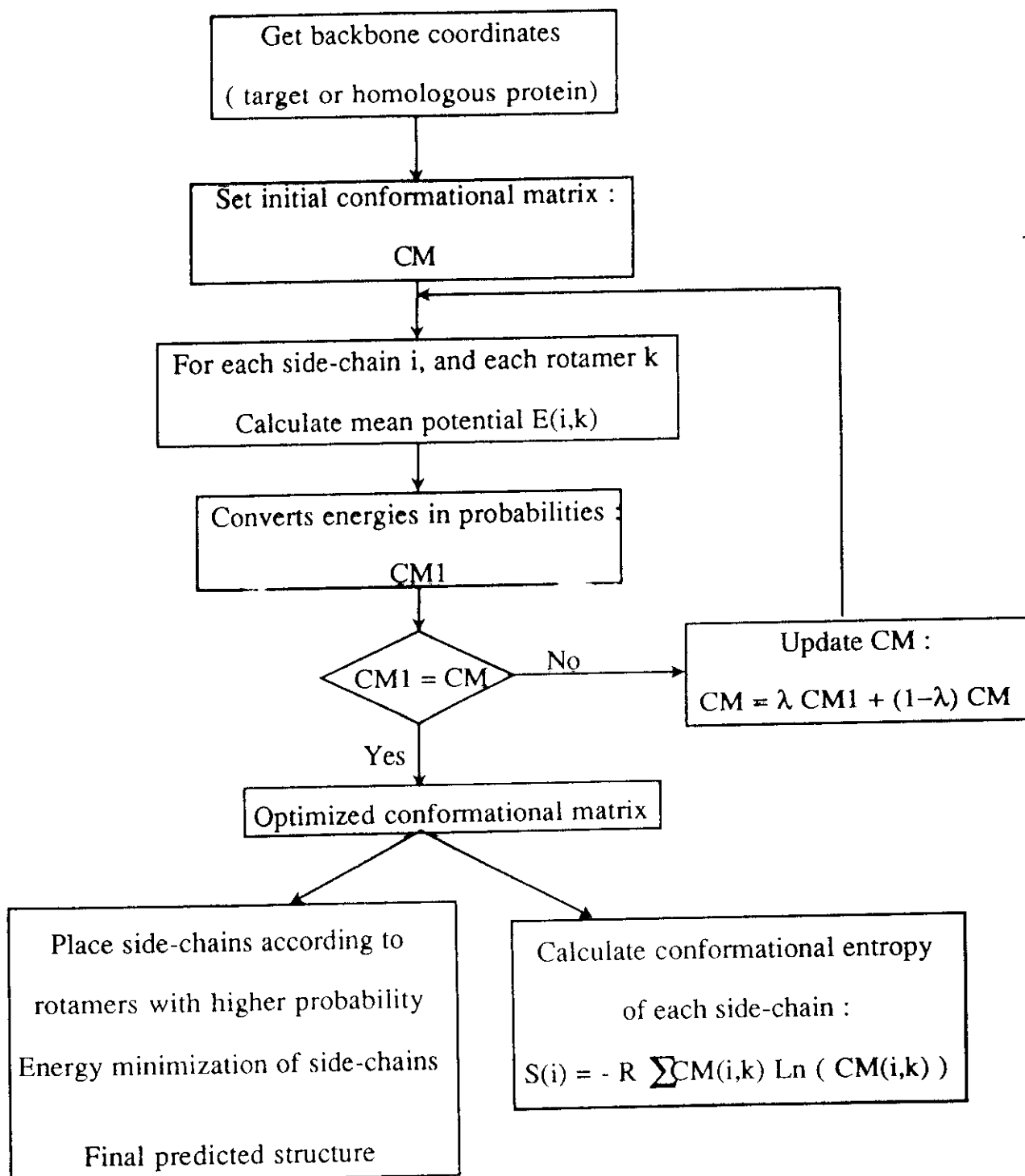
Two applications :

A) The rotamer with the highest probability describes the  
conformation of the side-chain

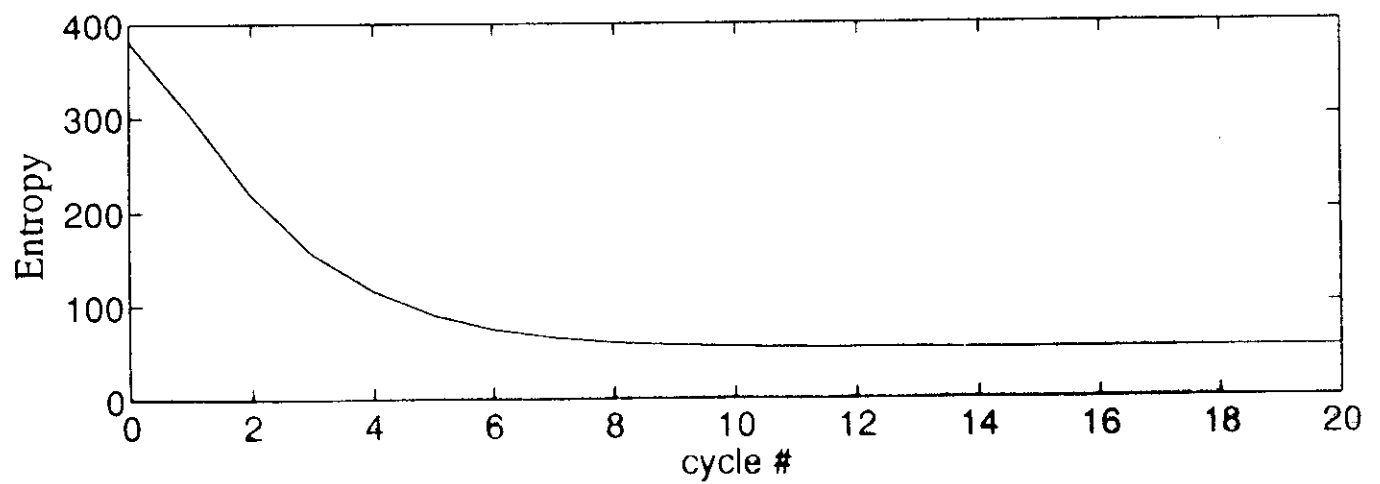
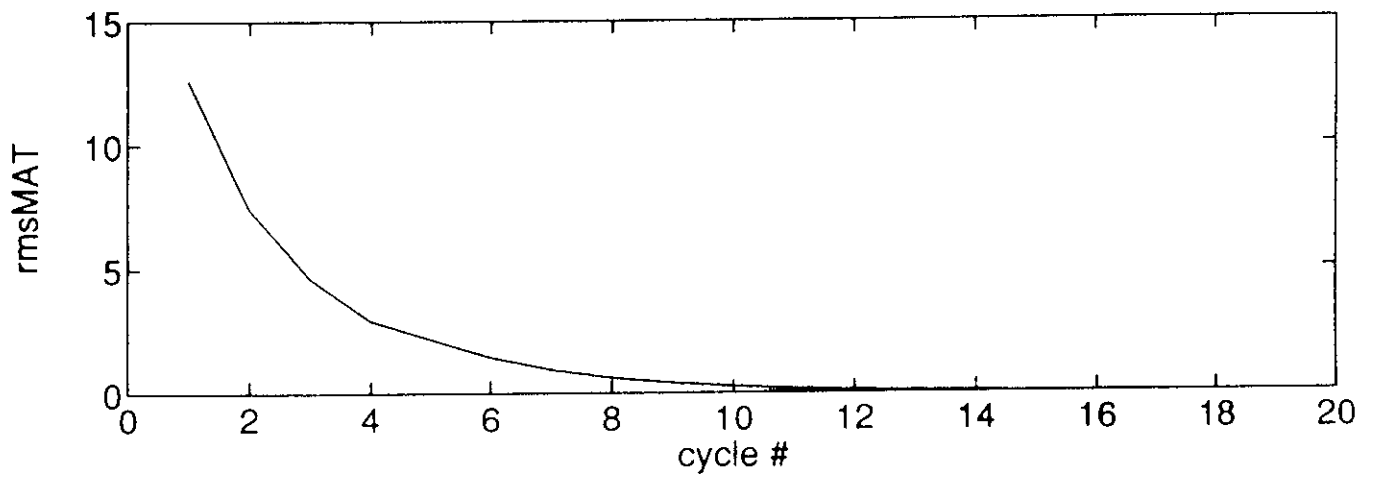
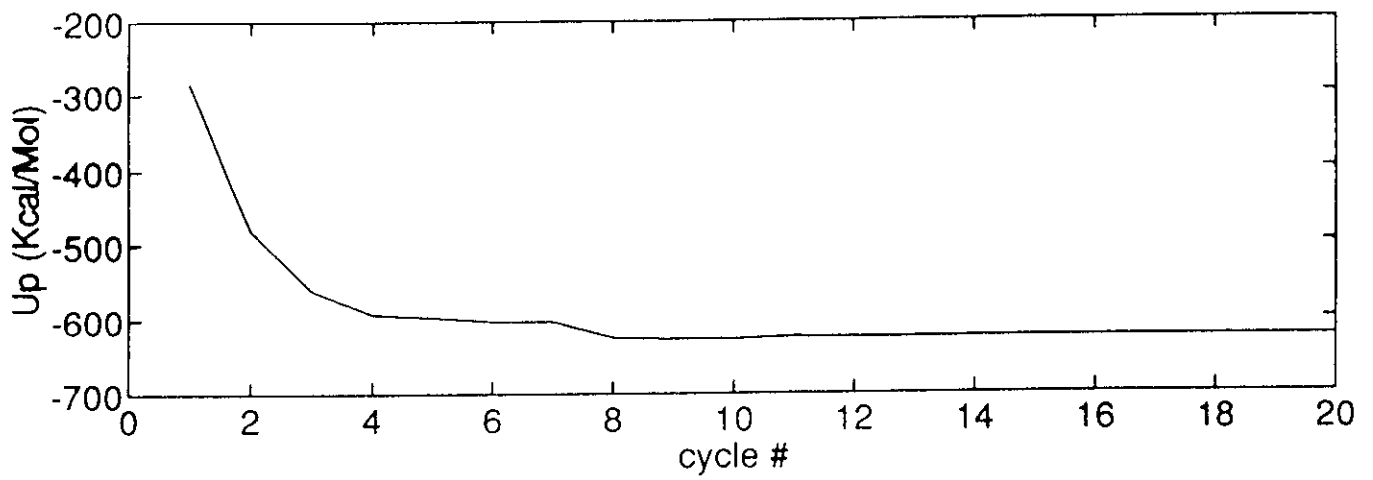
B) CM provides an estimate of the conformational entropy  
of the side-chains

For a residue i  $S(i) = -R \sum_k CM(i,k) \ln [ CM(i,k) ]$

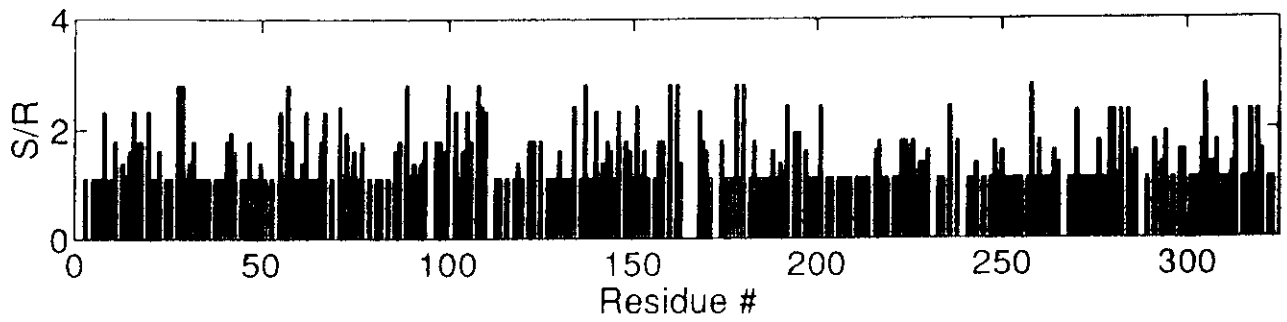
# CALCULATION OF THE SIDE-CHAIN CONFORMATION MATRIX



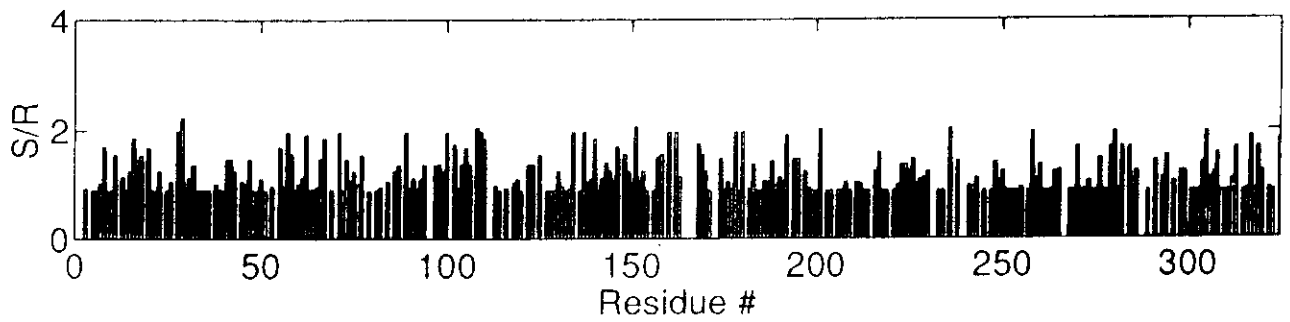
# SIDCHAIN POSITIONING FOR RHIZOPUSPEPSIN



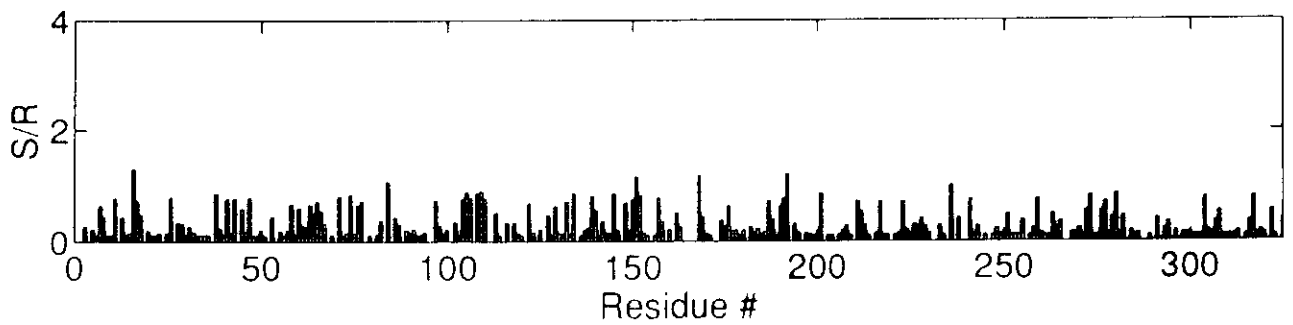
### Initial entropy



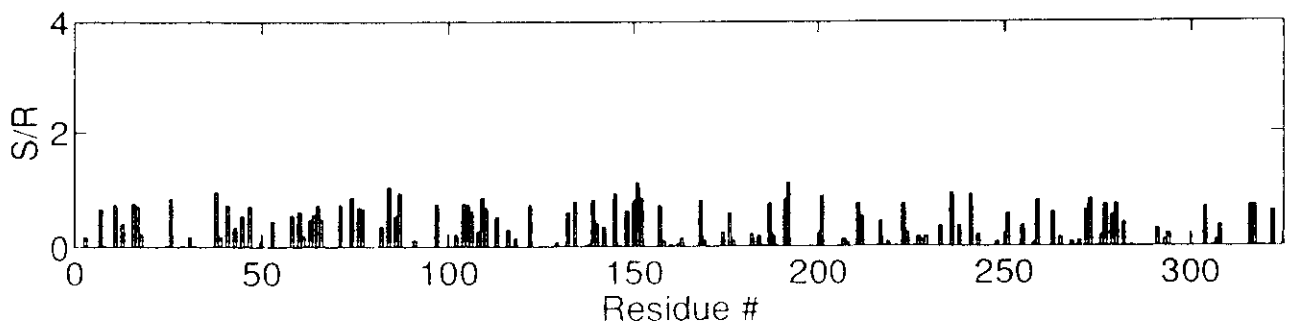
### After 1st cycle

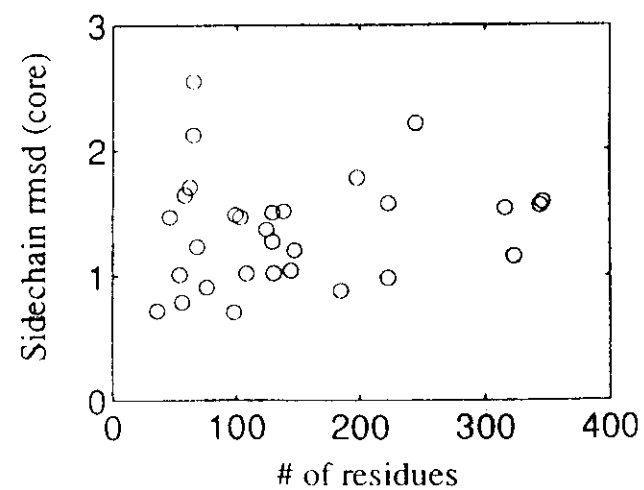
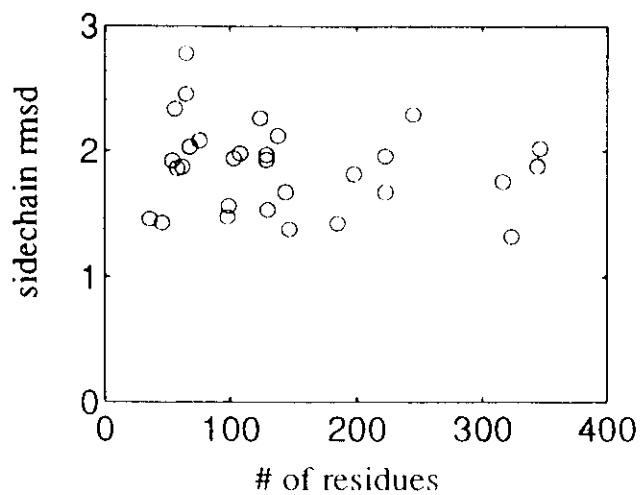
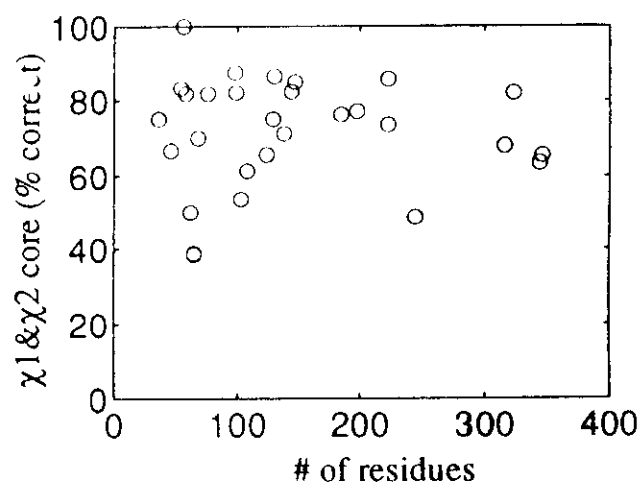
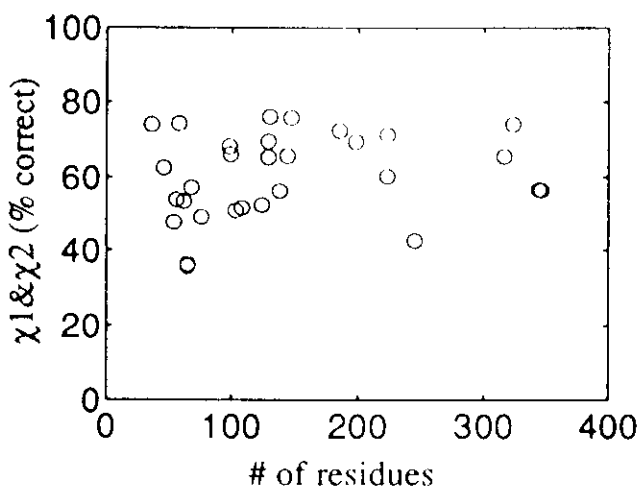
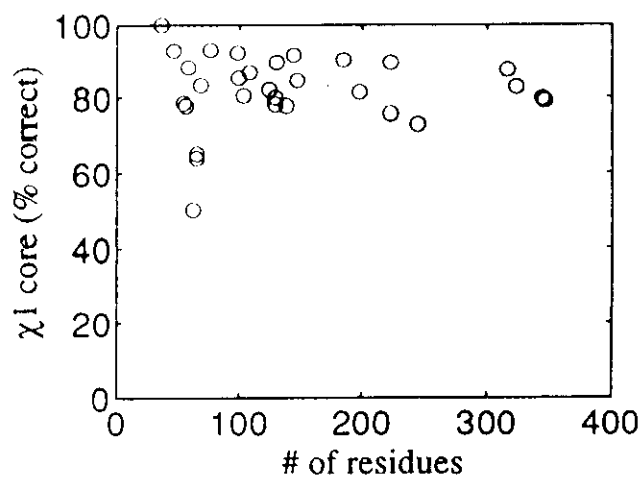
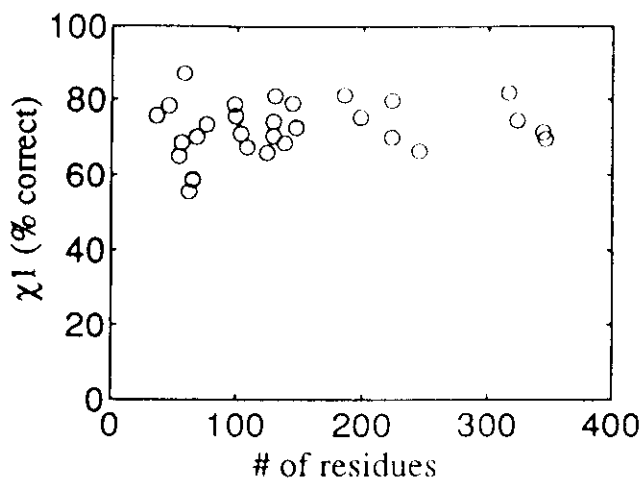


### After 5 cycles

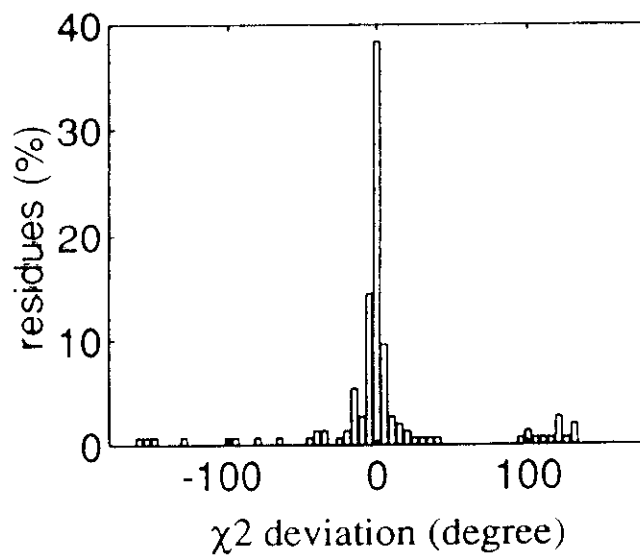
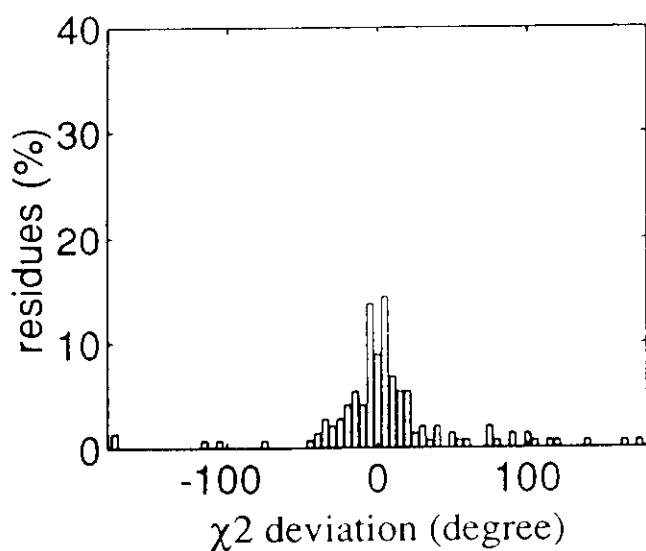
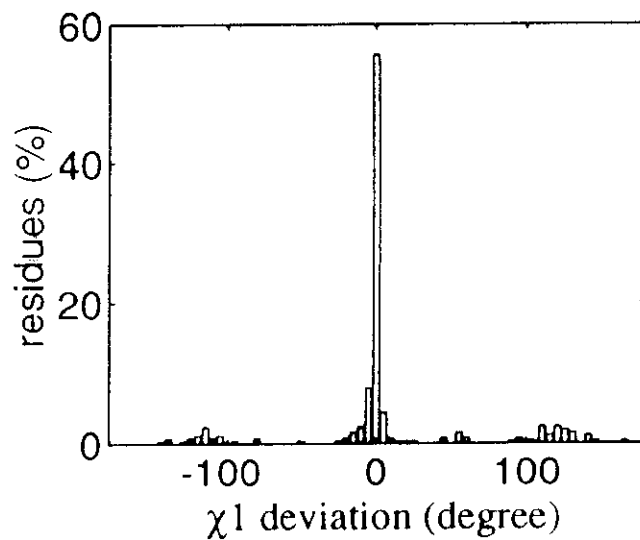
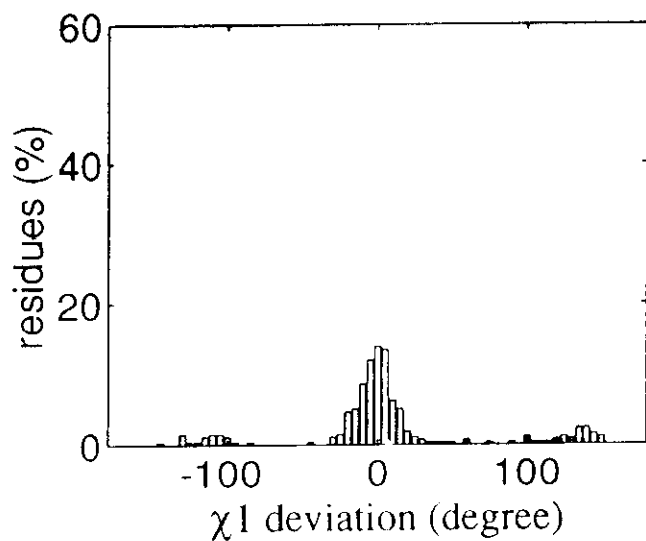


### Final





**Effects of the energy minimization on the accuracy of  
the dihedral angles prediction**



## Accuracy of side-chain prediction for specific side-chain types

Residue	$\chi_1$ (%)	$\chi_2$ (%)	$\chi_{1+2}$ (%)	Rmsd (Å)
<i>Large non polar</i>				
VAL	84			0.92
ILE	82	77	68	1.22
LEU	85	72	68	1.24
PHE	88	94	85	1.61
MET	84	73	62	1.81
<i>Aromatic side-chains</i>				
TRP	87	72	69	2.30
TYR	93	97	92	1.45
HIS	81	51	49	1.95
<i>Small side-chains</i>				
PRO	55			0.84
CYS	81			1.18
SER	42			1.69
THR	80			0.97
<i>Polar side-chains</i>				
ASN	73	57	46	1.84
GLN	73	74	60	2.02
ASP	64	81	56	1.83
GLU	66	66	45	2.15
LYS	68	72	50	2.44
ARG	66	81	54	3.05

## Possible causes for discrepancies

### A) Incomplete force-field

- 1) No specific terms for hydrogen bonds ( SER, THR ...)
- 2) Solvent not included ( surface residues )

### B) Inherent clashes due to the rotamer library

Test case : BPTI; energy : - 167 Kcal/Mol

"Ideal" structure : each sidechain replaced by the closest (r.m.s.d.) rotamer

Library	No. of rotamers	U (Kcal/Mol)		$\chi_1$ (%)	$\chi_{1+2}$ (%)
		Ideal structure	Final model		
Tuffery et al (1991)	303	81	-39.4	87	74
LIB1 (120°)	552	61	-56.7	85	67
LIB2 (60°)	3474	-75	-79.1	67	57
LIB3 (30°)	4512	-84	-110	65	54
BPTI			-167		

How do we know that a given sequence is compatible with a given 3D structure?

→ obvious answer if there is seq. homology.

→ Phenomenology of sequence evolution -

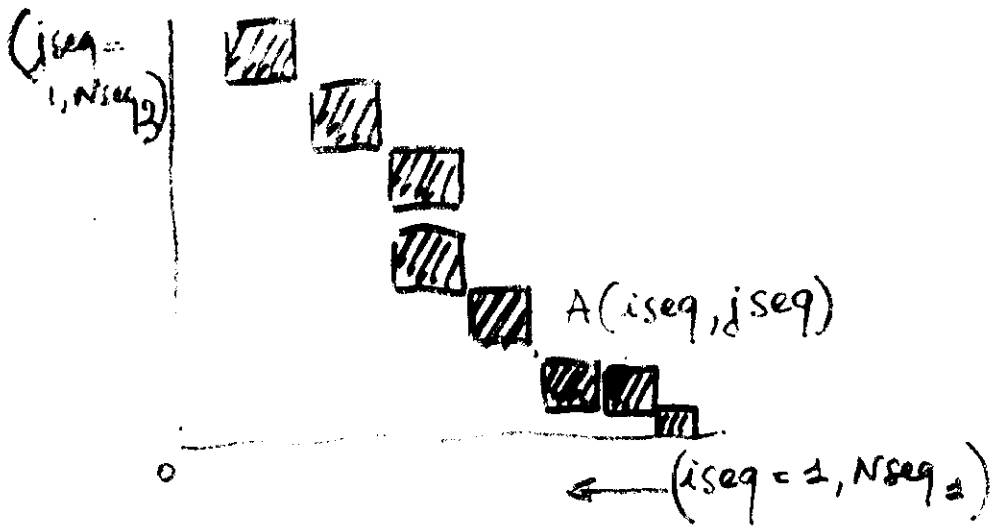
- Definition of conservative mutations
- The only strictly conserved residues for a family of proteins that has widely diverged are functional (only a handful of them ⇒ very difficult to detect) -
- Evolutionary information can and should be used to refine mutation matrix -

→ Molecular Mechanisms at stake

- point mutations - constant drift?
- Gene duplication & Gene fusion -

# Aligning two sequences ...

① The algorithm: Dynamic programming Method

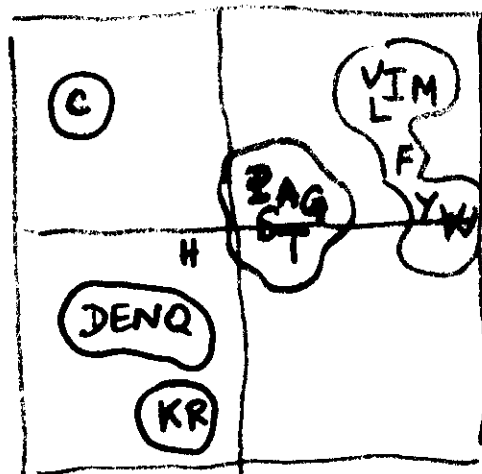


→ Table of amino acid substitution + gap penalties.

$$M(iaa, jaa) : \begin{matrix} iaa = 1, 20 \\ jaa = 1, 20 \end{matrix}$$

body empirical (M. Dayhoff, 70's)

Principal component analysis → projection  
on a plane



# Applications

## • Genome sequencing projects

- Mycoplasma ~ 480 genes

- Helicobacter pylori

- E. coli ~ 3000 genes (50%) - S. typhimurium -  
B. subtilis

- yeast 6000 genes. 12 Mbases.

- Rice, C. elegans, A. thaliana, D. melanogaster

- Human genome .  $3 \cdot 10^9$  bp ~  $10^5$  genes.

## • Automatic comparison of new seq. to existing seq. D (whole seq.)

→ 1/3 with clear homolog of known function

→ 1/3 with homolog of unknown function.  
(or unproved).

→ 1/3 with no homolog in Seq. Data Base.

# Going beyond obvious seq. homology (>30%)

## ① Evolutionary information - M. Gribskov

Define a position dependent Matrix  $M(i, j)$  from previous multialignment - Profile Method, confined to blocks of aligned sequences -

## ② Make use of structural information - J. Bowie

① define to which class each a.a. along the seq. belongs, according to criteria calc. from structure - solvent accessibility  
- secondary structure

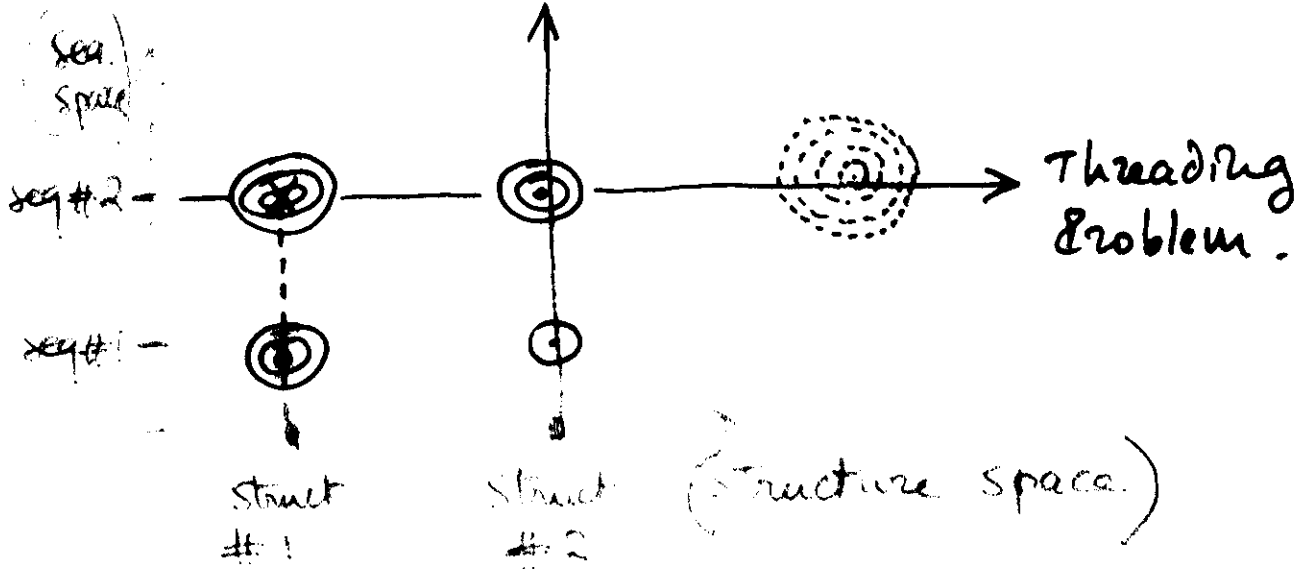
② convert belonging to one structural class into mutability to other a.a. types

③ Critique of this approach -

- solvent accessibility really conserved among structurally aligned positions - M. Sternberg
- Atomic properties - But Backbone can adjust

# Sequence Design

A. Godzik.



- The threading problem requires:

- good potentials

- a method to align the gaps, and to get rid of the "frozen approximation"

- The sequence design pb requires:

- good potentials

- a method to avoid "overdesigning" and to identify other folds.

Answer: Maximize

$$E = \langle E \rangle_{\text{all other folds}}$$

$$\sigma_{\text{all other folds}}$$

# The potential Energy Functions

- Residue-Residue contact potentials - (statistical pot.)

$$E = \sum_{i < j} \epsilon(r_i, r_j) f(\vec{r}_i - \vec{r}_j)$$

- Miyazawa & Jernigan, 1985 ; Sippl (dist. dep), 1990  
Bryant & Lawrence, 1991 ...

- All of them can distinguish the right structure among an ensemble of wrong structures -

- They do not have the same reference state.

from  $E_{ij}$ , define  $\rightarrow E_{ij}^{\text{ideal}} = \frac{E_{ii} + E_{jj}}{2}$  (Burial).

[Godzik, 1995].  $\rightarrow E_{ij}^{\text{excess}} = E_{ij} - E_{ij}^{\text{ideal}}$  (2-Body).

- They are not good for ~~primary~~ <sup>Seq. design</sup> experiments (MC).  
 $\rightarrow$  overdesigning -

- There are not truly speaking two body potentials that can be added directly - Need of 3 body-pot.

$\rightarrow$  Control experiment with  $3 \times 3 \times 3$  lattices, where all compact SAW are known: 103, 346 of them...

