



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 - 19

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

22 JULY - 2 AUGUST 1996

**"Protein structure prediction using
optimized energy functions"**

PART II

**Zaida LUTHEY-SCHULTEN
University of Illinois
School of Chemical Science
505 South Mathews Avenue
IL 61801 Urbana
U.S.A.**

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

How to "Learn" an Energy Function

I. Look for Features

The "Associative Memory Energy"

$$E = \sum_{\text{pairs } \langle i, j \rangle} \sum_{\text{examples}} \delta(\text{Target Sequence, Example}) \times \Theta(r_{ij} - r_{ij}^{(\alpha)})$$

How much do sequence in target and example agree?
 ≠ 0 if distance matches example distance

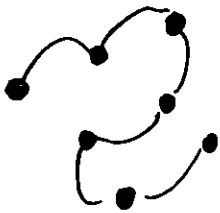
Examples



II. Assume simple physics

$$E = \sum_{\text{pairs}} \delta(A_i, A_j) \Theta(r_{ij} - d)$$

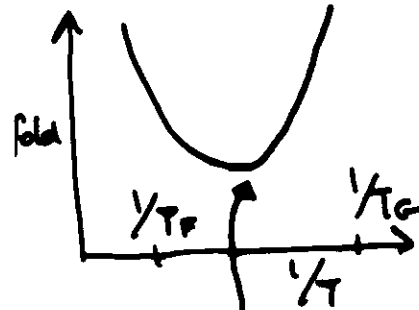
≠ 0 if pairs touch



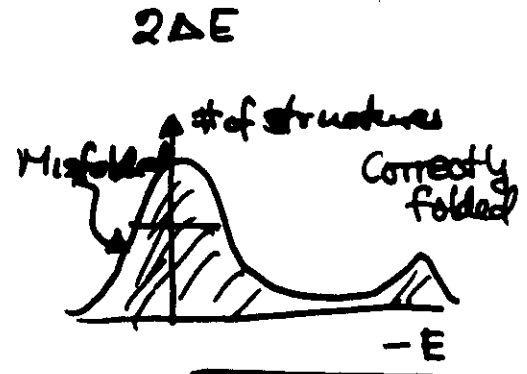
Use Minimum Frustration Principle to find δ 's from known "worked" examples!

Optimization of Energy Functions

R. Goldstein
 E. Luthey-Schulten
 P. Wolynes
 PNAS 92



Best simulated annealing if T_f/T_c is a maximum



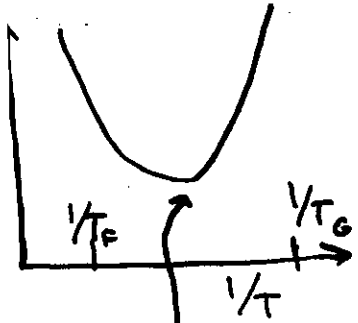
Best discrimination if $\delta E_s / \Delta E$ is a maximum

Same Variational Problem!

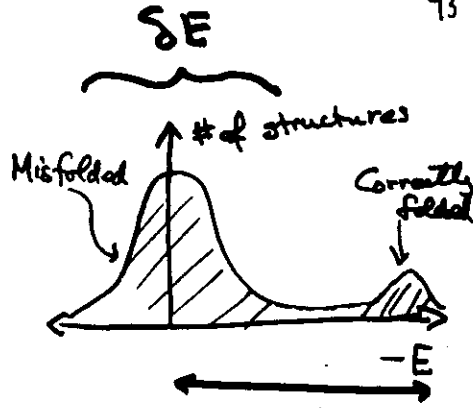
$$\text{Max } R = \frac{\delta E}{\Delta E} = \frac{(A \cdot \delta)}{\sqrt{\delta \cdot B \cdot \delta}} \Rightarrow \underline{\delta} = \underline{B}^{-1} \cdot A$$

Energy Landscaping and the Minimal Frustration Principle

R. Goldstein
Z. Lutnyk, Schulte
PGW
PNAS '92
'92
'93



Best simulated annealing if T_F/T_G is a maximum.



Best discrimination if $E_p/\Delta E$ is a maximum

Same Variational Problem!
within Random Energy Approximation

Minimize

$$\Gamma = \frac{E_p}{\Delta E} = \frac{(A \cdot \underline{x})}{\sqrt{\underline{x} \cdot B \cdot \underline{x}}}$$

$$\underline{x} = B^{-1} A$$

Heteropolymers and Protein Folding

What are the problems?

1. Chemistry should tell us

$$E(\{q_i\} | \{r_i\})$$

Biology give us $\{q_i\}$

Molecular Dynamics, simulated annealing $\{r_i\}$

2. Evolution's problem:

Given a $\{r_i\}$ find a $\{q_i\}$

3. Code breakers' problem:

Given $\{r_i\}$ and $\{q_i\}$ find

$$E'(\{q_i\} | \{r_i\})$$

such that

$$E'(\{q_i\}^{\text{new}} | \{r_i\}) \Rightarrow \{r_i\}^{\text{New}}$$

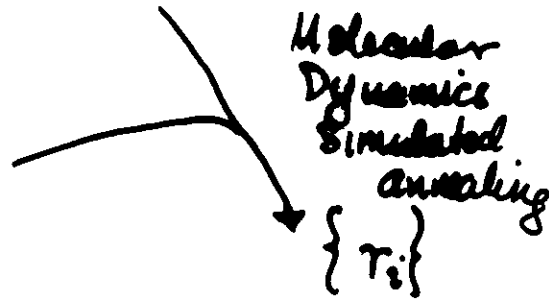
Code breakers need to understand both evolution and physics!

Protein Folding - What are the Problems?

1. Chemistry and physics should tell us

$$E(\{q_i\} | \{r_i\})$$

Biology gives us $\{q_i\}$



2. Evolution's problem: (Design)
Given a $\{r_i\}$ find a $\{q_i\}$

3. Code breaker's problem: μ ← database
Given $\{r_i\}^\mu$ and $\{q_i\}^\mu$ find

$$E'(\{q_i\} | \{r_i\})$$

Such that $E'(\{q_i\}^{new} | \{r_i\}) \Rightarrow \{r_i\}^{new}$

Optimized Energy Functions

I. Self-Consistently Optimized \mathcal{H}_{local} :

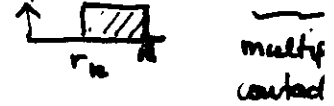
Threading Function

$$\mathcal{H}_{local} = E_{profile} + E_{contact} + E_{hydrogen\ bonds} + E_{gaps} + E_{exp\ constr}$$

(backbone α, β)

$$E_{profile} = \sum_{i=1}^N \delta^P(A_i, SS_i, SA_i)$$

$$E_{contact} = \sum_{i < j} \sum_{k=1}^2 \delta_k(A_i, A_j) \underbrace{U(r_k - r_{ij})}_{\substack{\uparrow \\ r_k \\ \downarrow}} + \frac{E}{3_{\text{loop}}}$$



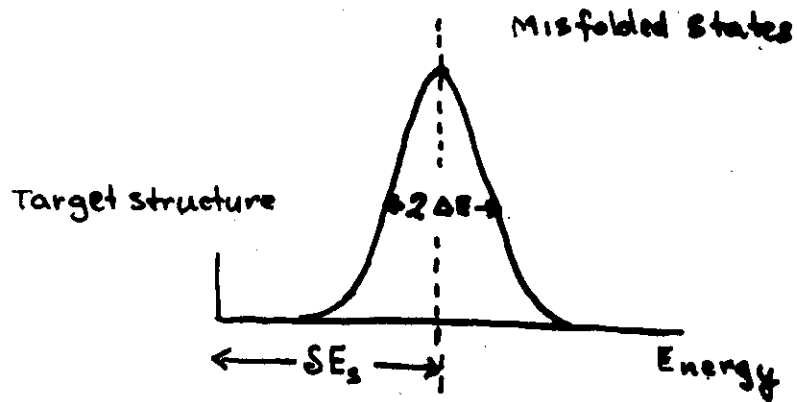
E_{gaps}

II Associative Memory Hamiltonian \mathcal{H}_{AM}
Molecular Dynamics

$$\mathcal{H}_{AM}^{(crs)} = - \underbrace{\sum_{\mu \text{ memories}} \sum_{i < j} \delta_{ij}^\mu \theta_{\text{gaussian}}(r_{ij} - r_{ij}^\mu)}_{\text{pattern recognition}} + \mathcal{H}_0$$

(backbone rigid)

Optimization of γ



Target Configuration: $E_T = \sum_i \lambda_i^T \gamma_i$

Misfolded State k: $E_k = \sum_i \lambda_{ki} \gamma_i$

$$\delta E_s = A\gamma$$

$$\Delta E^2 = \gamma B \gamma$$

where

$$A_i = \lambda_i^T - \langle \lambda_{ki} \rangle$$

$$B_{ij} = \langle \lambda_{ki} \lambda_{kj} \rangle - \langle \lambda_{ki} \rangle \langle \lambda_{kj} \rangle$$

Optimal γ : \Rightarrow maximise $R = \frac{E_T}{E_k}$ with respect to γ :

$$\gamma = B^{-1}A \quad (1)$$

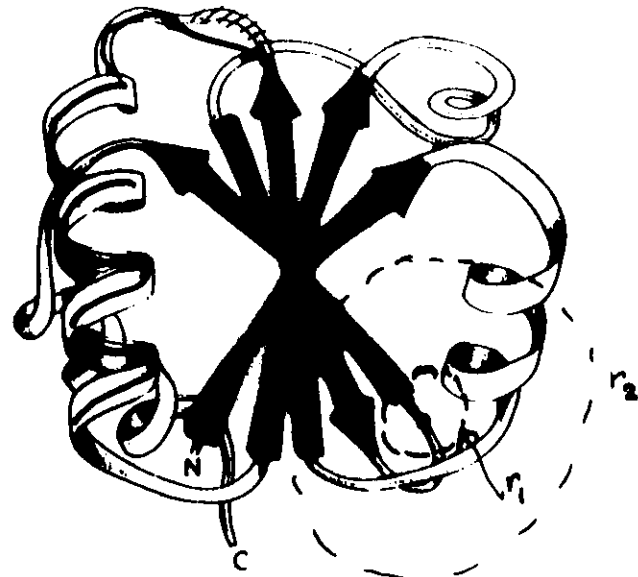
Mean-Field Sequence-structure Alignment

$$E_{\text{local}} = E_{\text{profile}} + E_{\text{contact}} + E_{\text{hydrogen bonds}} + E_{\text{gaps}} + E_{\text{exp. constraints}}$$

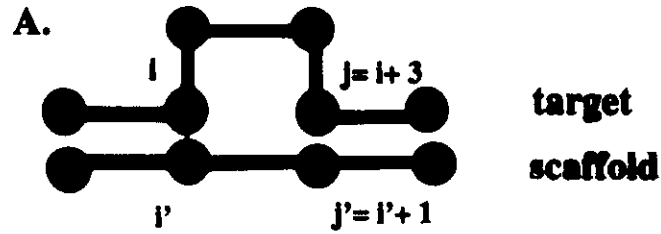
$$E_{\text{profile}} = \sum_{i \in I} \gamma_i^P [A_i, \text{secondary structure}(i), \text{surf. access}(i)]$$

$$E_{\text{contact}} = \sum_{i < j} \sum_{k=1}^2 \gamma_k(A_i, A_j) u(r_k - r_{ij}) + E_{\text{multiple bonds}}$$

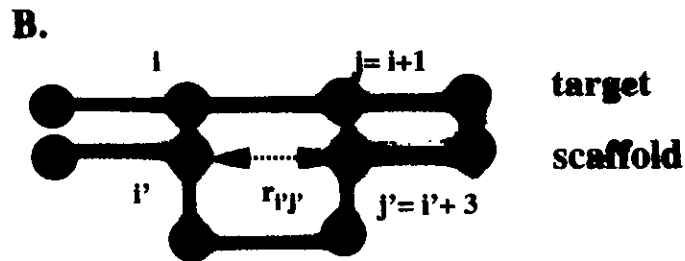
(Ref. Finkelstein & Deza, Skolnick et al, Thornton, et. al)



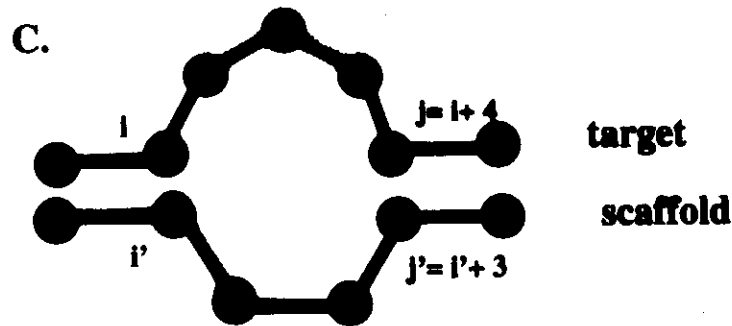
Statistical Mechanical Basis of Sequence-Structure Alignment Algorithms



$$E_g(A) = \gamma_1 + \gamma_2(j-1)$$



$$E_g(B) = \gamma_3 + \gamma_4 E_{1,j'} + \gamma_5 E_{1',j'}$$



$$E_g(B) = \gamma_6 + \gamma_7(j-1) + \gamma_8 E_{1,j'} / (j-1)$$

Conventional: (Smith-Waterman Algorithm)

$$A = (a_1, a_2, \dots, a_m) \quad 1-D \text{ information}$$

$$B = (b_1, b_2, \dots, b_n)$$

2 Metrics: ① $S(A, B) = \{s(a_i, b_j), \dots\}_{m \times n}$
 $s(a_i, b_j)$: Dayhoff Evol. scale
 (indels) ② $w = \alpha + \beta \log k$ \leftarrow length of indel
 penalty \uparrow non-uniform (Levitt, Chentis)

SMA: Sequence-Structure Alignment (Wolynes et al.)
 $F = E + \sum n_l \mu_l$ (free-energy)

where μ_l is a gap penalty $\mu_l(l=i-j, r_{i,j}^s)$

$$\mu_l = \frac{\Gamma_r^2}{\Delta E_{rc}} \log \left\{ \frac{P_{rc} P_c^{\lambda_2+1}}{P_r^{\lambda_2+1} P_{c2}} \right\}$$

\uparrow log length
 \uparrow distance

$P_{c2}, P_{r2} \dots$ = Prob. that gap of type l appears in correct / random alignments.

① Hawaii Conference, ② R. Elber, editor "New Developments..."

Smith-Waterman (1981)

Local & global alignments based on similarity score H where

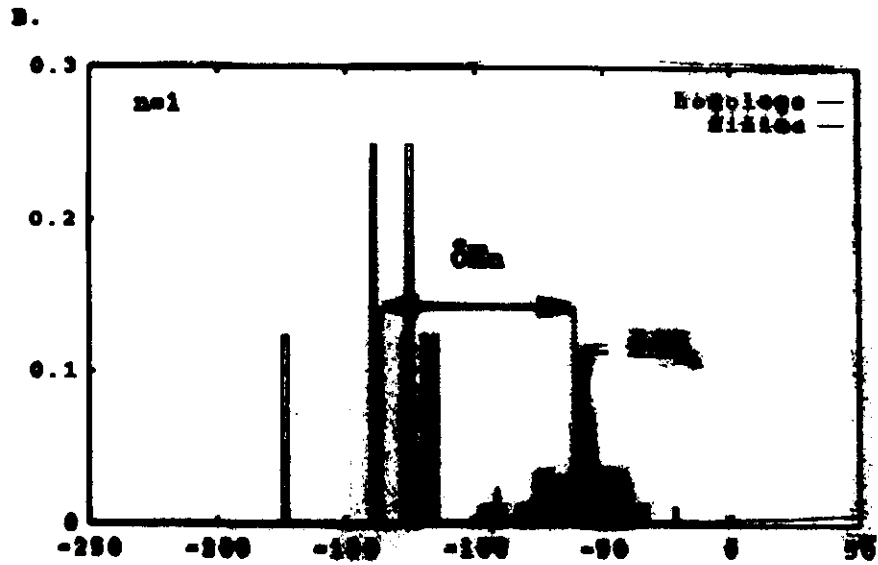
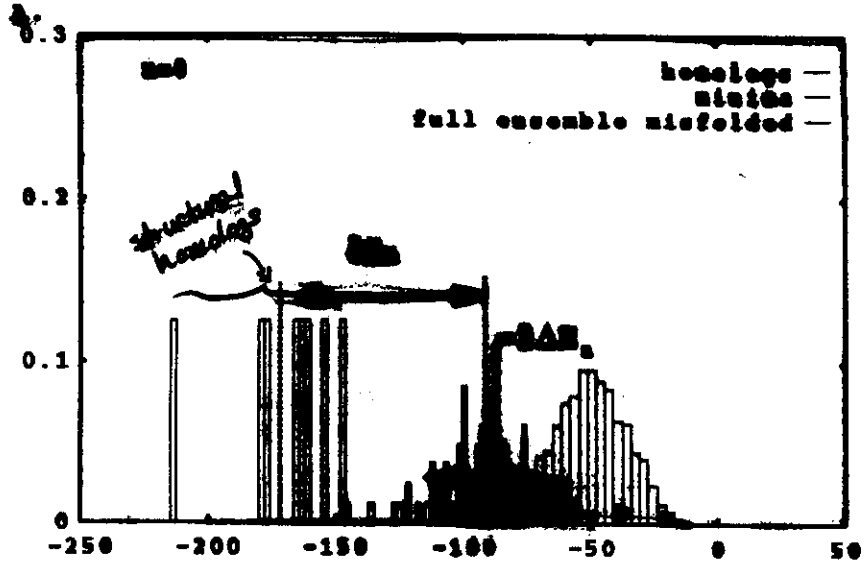
$$H_{ij} = \max \begin{cases} H_{i-1, j-1} + \delta(a(i), b(j)), \\ \max_{k \geq 1} H_{i, j-k} - W_k, \\ \max_{l \geq 1} H_{i-l, j} - W_l, \end{cases} 0$$

Discrimination Scores

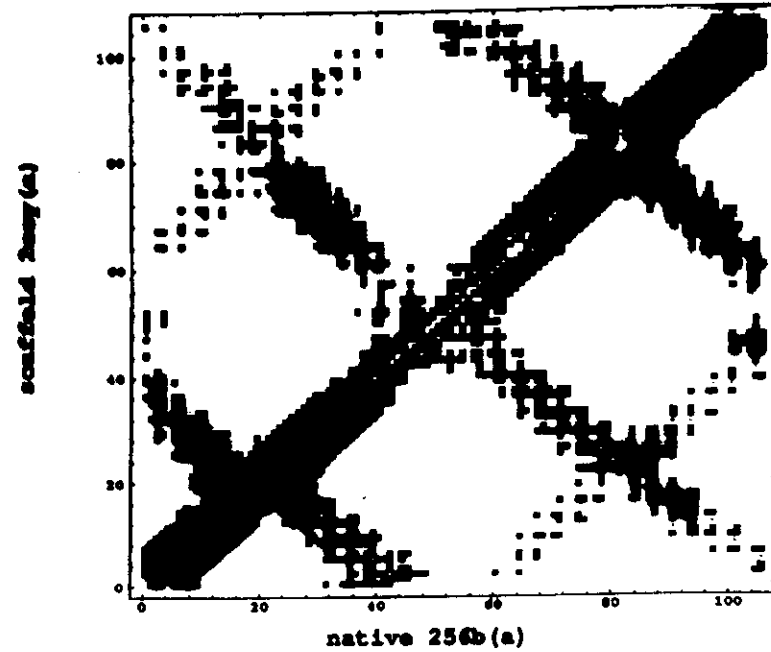
Training Proteins

PDB	NRRES	Name	$\frac{\delta E_n}{\Delta E_n}$			
			Zeroth	First	Second	Mean
α proteins						
2cro	63	434 Cro Protein	3.70	3.80	4.48	1.19
1wrp(R)	102	TRP Repressor (Trigonal Form)	2.41	3.80	4.44	1.84
1ccr	111	Cytochrome C	4.43	4.62	5.02	1.13
1hmq(A)	113	Hemerythrin (MET)	3.80	4.22	3.48	0.92
1bp2	123	Phospholipase	6.35	6.35	7.63	1.20
2hbb(A)	141	Hemoglobin (Deoxy)	2.37	2.90	3.36	1.42
3cln	143	Calmodulin	5.37	5.95	6.90	1.30
1fhh(G)	146	Hemoglobin (Deoxy, Human Fetal)	3.31	3.63	4.39	1.33
1mba	153	Sperm Whale Myoglobin (Deoxy)	3.79	4.62	5.49	1.45
β proteins						
6pcy	99	Plastocyanin	5.35	6.28	6.62	1.28
1rei(A)	107	Bence-Jones Immunoglobulin	3.93	4.15	3.50	0.91
2pas	123	Pseudoazurin (Cupredoxin)	6.55	8.69	8.85	1.35
2ilb	153	Interleukin-1 β	6.14	7.62	9.46	1.54
1f19(H)	215	R19.9 (IG*G2B=K=) Fab Fragment	2.14	2.79	2.55	1.20
3hfm(H)	215	IG*G1 Fab Fragment	4.71	5.40	6.62	1.41
2plv(B)	235	Poliovirus (TYPE 1, Mahoney Strain)	2.58	3.79	5.16	2.00
1ria(B)	238	Rhinovirus Serotype 1 (HRV1) Coat Protein	2.53	3.18	3.91	1.55
1cms	248	Chymosin B	3.36	4.19	5.56	1.65
$\alpha + \beta$ or α/β proteins						
1fdx	84	Ferredoxin	1.97	2.53	2.21	1.12
1alc	122	α -Lactalbumin	6.63	7.51	7.78	1.17
1rbb(A)	124	Ribonuclease B	5.68	6.49	6.91	1.22
1anc	126	Staphylococcal Nuclease	6.31	8.15	9.61	1.52
2dhf(A)	182	Dihydrofolate Reductase	3.36	4.28	5.01	1.49
2act	218	Actinidin (Sulfhydryl Proteinase)	5.06	6.59	7.80	1.50
2pk	279	Proteinase K	3.76	6.33	6.89	1.83
1pfr(A)	319	Phosphofructokinase (R-State)	3.85	5.56	7.10	1.84
2dx	331	Apo-Lactate Dehydrogenase	2.92	3.92	5.03	1.72
3gd(G)	334	D-Glyceraldehyde-3-Phosphate Dehydrogenase	3.44	4.72	7.24	2.10
2liv	344	Leucine/Isoleucine/Valine-Binding Protein	5.34	6.90	8.83	1.65
mean						
1.44						

Self-Consistent Optimization

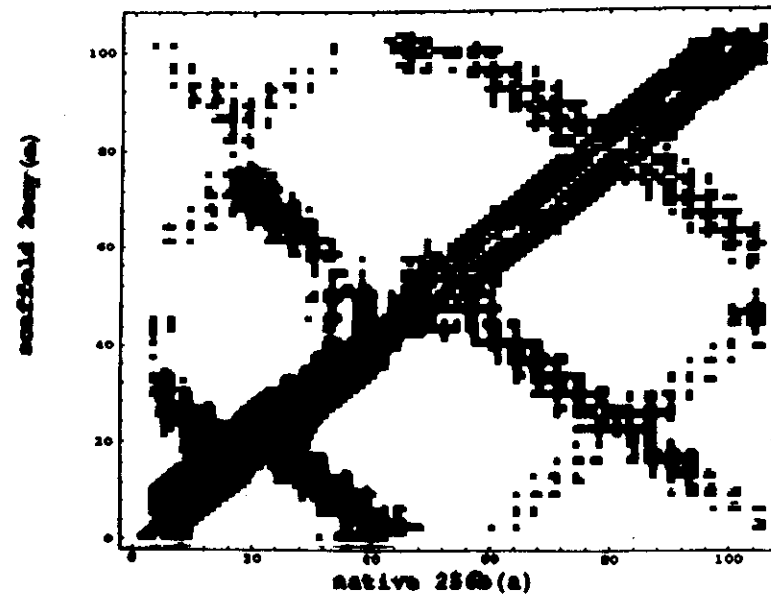


Energy Alignment



95%
Sec. Str.

Genetic Alignment



65%
Sec. Str.

!fold

Secondary structure assignment

Methanococcal adenylate kinase

131

```

W:  tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt
P:  ---tttttttt tttttttttt ---tttttt ---tttttt tttttttttt
E:  ---tttttttt tttttttttt ---tttttt ---tttttt tttttttttt
E:  ---tttttttt tttttttttt ---tttttt tttttttttt tttttttttt

```

181

```

tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt
tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt
tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt
tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt

```

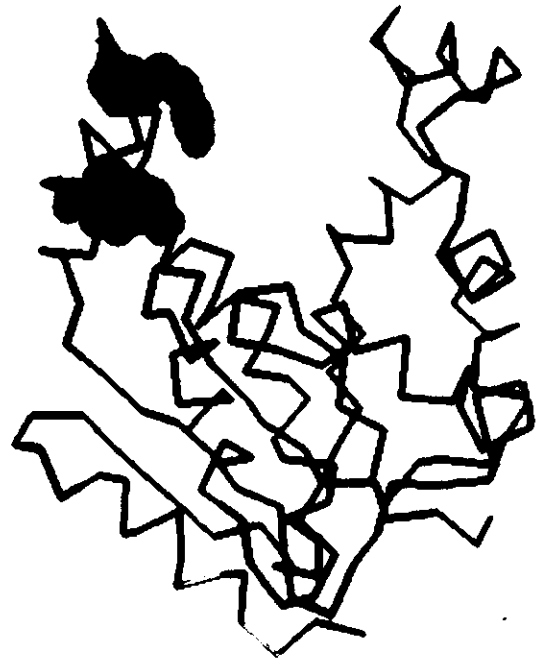
231

```

tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt
tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt
tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt
tttttttttt tttttttttt tttttttttt tttttttttt tttttttttt

```

(75% secondary structure)
(69% secondary structure)
(73% secondary structure)



Threading as good as the scaffold!

```

MKNKV VVIVGVPVG STVTNKAIE ELKKEGIEYK IVNFGTVMFE
MEEKLKSKI IFVVGPGSG KGTQCEKIVQ K-----YGYT HLSTGDLIRA

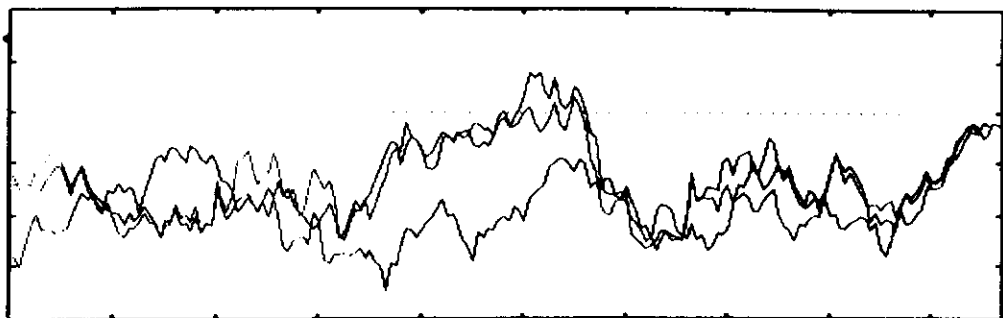
IAK-----EE GLVEHRDQLR KLPPEEQKRI OKLAGKKIAE MAKEFNIVVD
EVSSGSAROK MLSEIMEKGQ LVP--LETVL DMLRDAMVAK VDTSKGFLID

PHSTIKTPKG YLPGLPAWVL EELNPDIIVL VEAE---NDE ILMRRLKDET
---GYPREVK QGEEFERKIG QPT---LLLY VDAGPETMTK RLLKRGETSG

IQRDFESTED IGEHIFMNRG AAMTYAVLTG ATVKIKNRD F--LLKDAVQ
VDDNEETIK KRLETTYKAT EPTIAFYEKR GIVRVNAEG SVDDVFSQVC

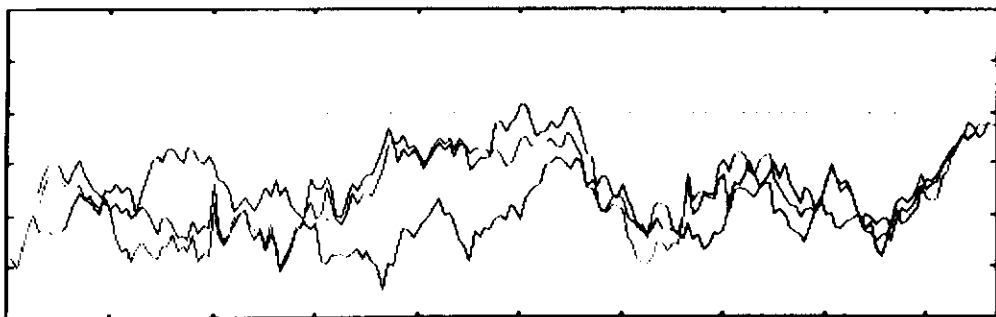
ELIEVLK
HLDTLK

```



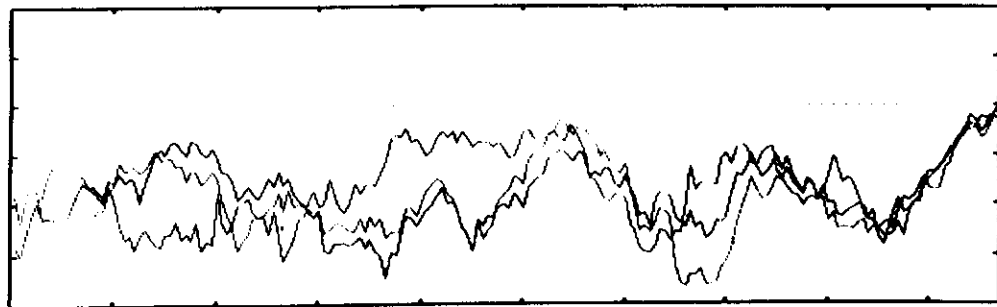
100

180



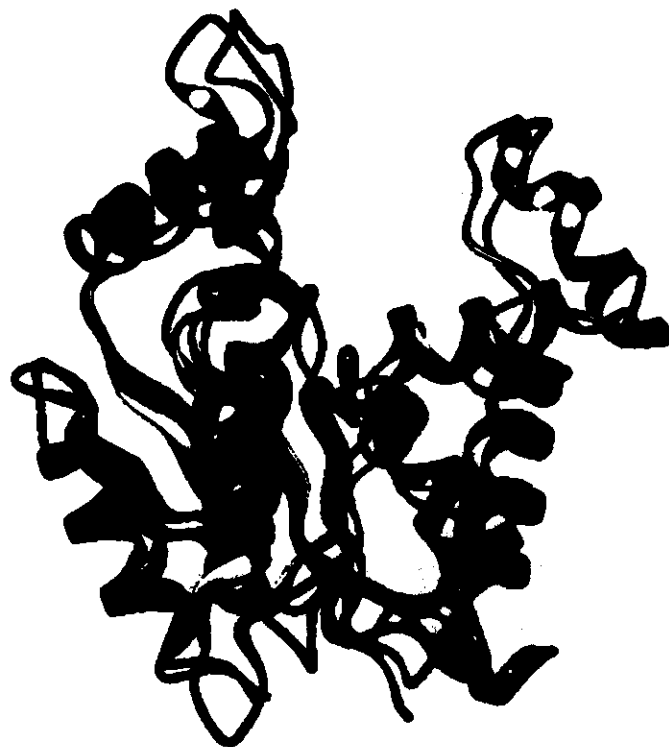
100

180

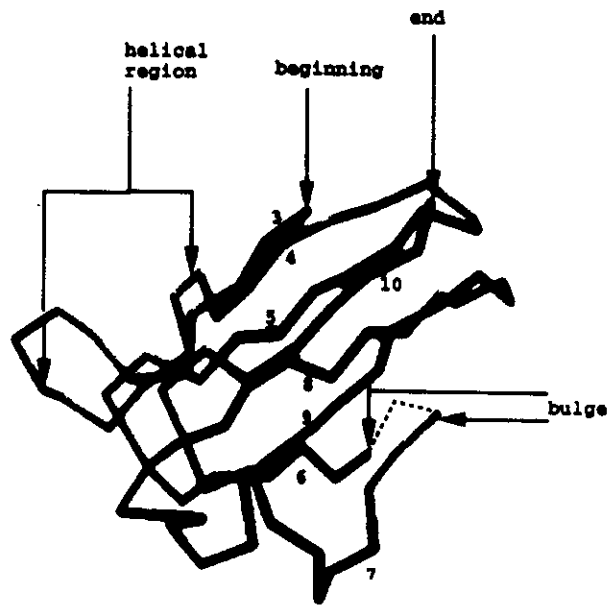


100

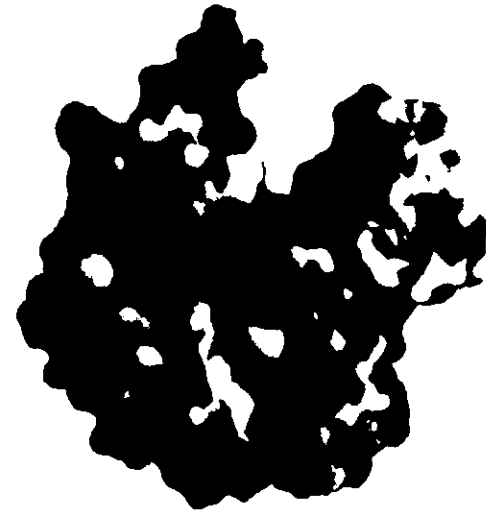
180



Prediction of Copper A Structure



A.



B.



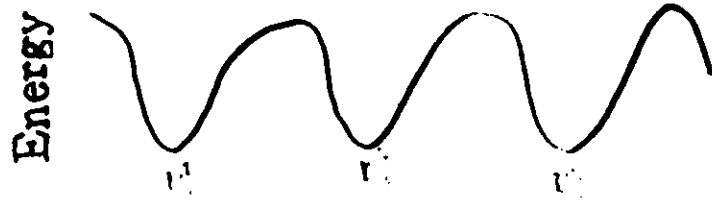
1D - Associative Memories Hamiltonian

Multiple Memories $\{\mu\}$

$$\mathcal{H}_{AM}(r_{ij}) = - \sum_{\mu} \sum_{\langle ij \rangle} \theta(r_{ij} - r_{ij}^{\mu}) + \mathcal{H}_0$$

↑ gaussian
↑ Backbone

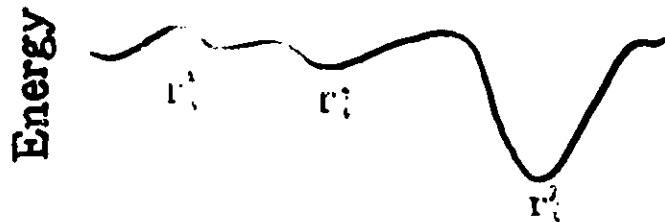
↑ pre-aligned



Including correlations between target and memory sequences

$$\mathcal{H}_{AM}(r_{ij}) = - \sum_{\mu} \sum_{\langle ij \rangle} \gamma_{ij}^{\mu} \theta(r_{ij} - r_{ij}^{\mu}) + \mathcal{H}_0$$

$$\gamma_{ij}^{\mu} = \gamma(S_i^T, S_j^T, S_i^{\mu}, S_j^{\mu})$$



Associative Memory Hamiltonian

Multiple Memories $\{\mu\}$

$$\mathcal{H}_{AM}(r_{ij}) = - \sum_{\mu} \sum_{\langle ij \rangle} \theta(r_{ij}^{c_{\alpha}, c_{\beta}} - r_{ij}^{\mu}) + \mathcal{H}_0$$

↑ gaussian
↑ memories pre-aligned

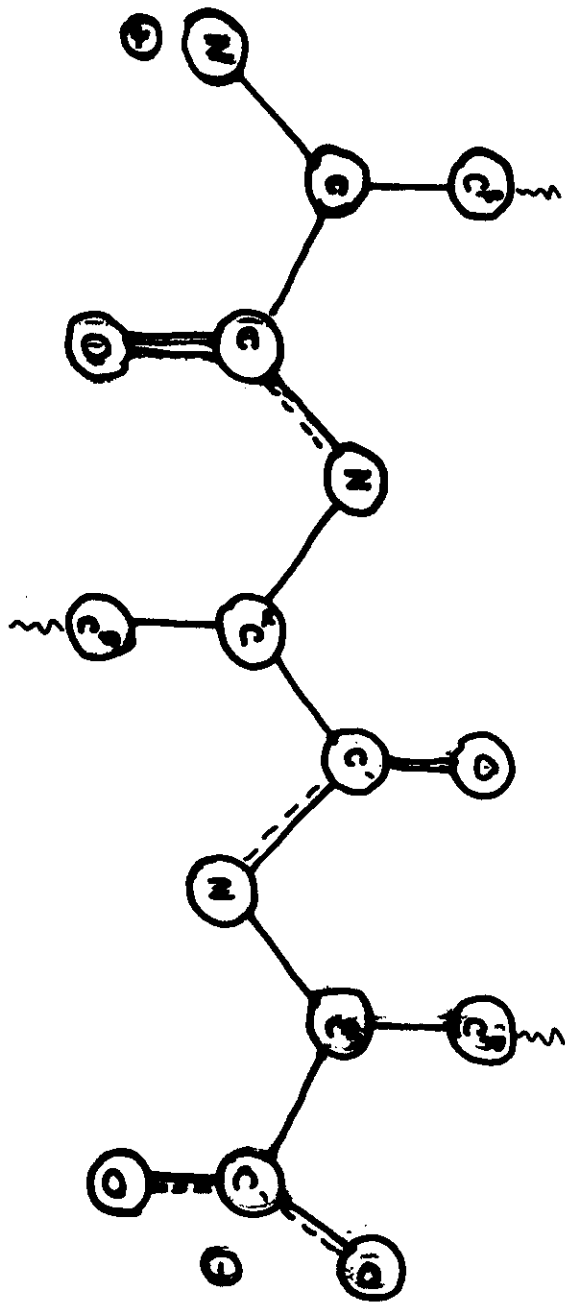
Including correlations between target and memory sequences

$$\mathcal{H}_{AM}(r_{ij}) = - \sum_{c_{\alpha}, c_{\beta}} \sum_{\mu} \sum_{\langle ij \rangle} \gamma_{ij}^{\mu} \theta(r_{ij}^{c_{\alpha}, c_{\beta}} - r_{ij}^{\mu}) + \mathcal{H}_0$$

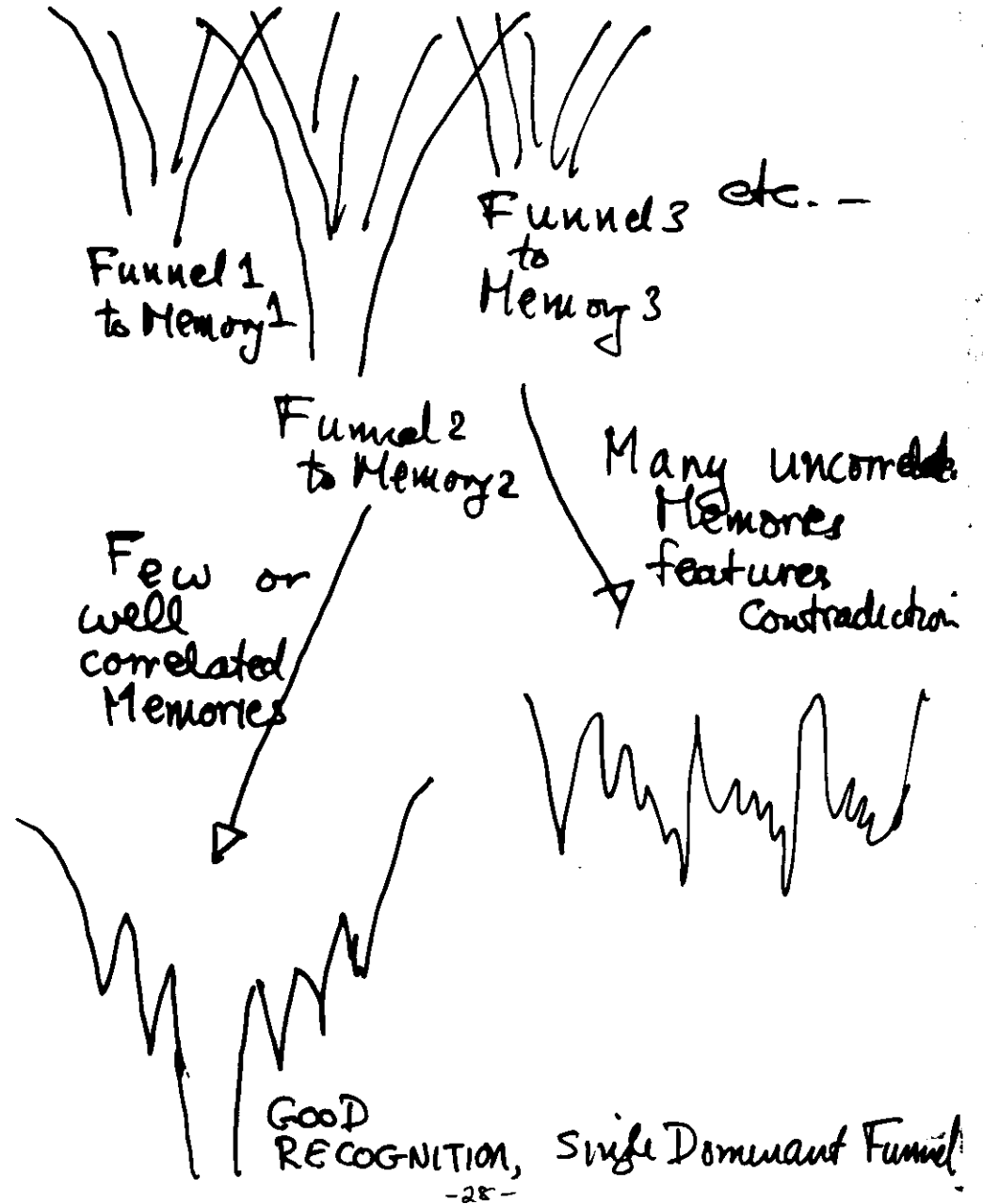
γ_{ij}^{μ} can encode:

- Hydrophobicity 2^4 (2 isoub ang consensus)
- Probability of mutation
- Predicted, Actual, or Instantaneous secondary structure
- Location of the residue in the protein proximity
- Properties of nearby residues

48 gammas: hydrophobicity x 3 categories $(j-i \in \mathbb{Z})$
 $(j-i) \in \mathbb{Z}, r_j \neq 1$



Energy Landscape of AM Energy Function

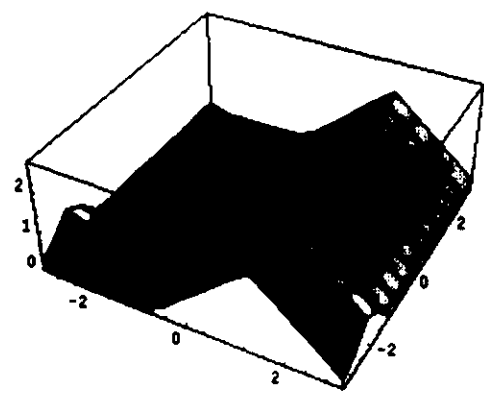


Full H_{AMH} and Protein Back bone Geometry

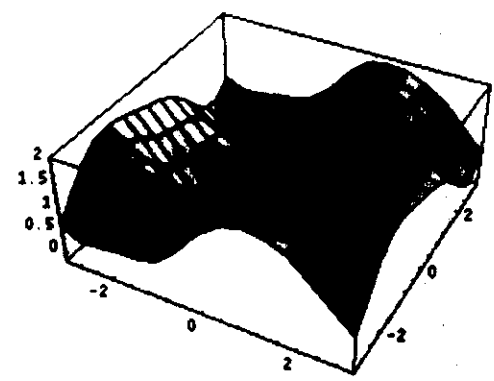
$$\begin{aligned}
 H(r_i^{C\alpha}, \phi | \{S_i^T\}) = & \sum_{\mu \text{ memories } i, j} \sum_{k, l = \alpha, \beta} \sum \delta_{ij}^{\mu} \theta(r_{kij}^{S\mu} - r_{kij}^{C\alpha, \mu}) \\
 & + H_{ev} : \Delta_{ev} \sum_{i, j} \begin{cases} (r_{ij} - 2A)^2 & r_{ij} < 3 \\ 0 & r_{ij} > 3 \end{cases} \\
 & + H_{chirality} : T [\vec{r}_{C_i' C_i} \times \vec{r}_{C_i' N_i}] \cdot \vec{r}_{C_i' C_i} \\
 & + T \Delta \phi + T H \phi \\
 & + T H_{\text{harmonic}} (r_{C_i' C_i}, r_{N_i C_i}, r_{N_i C_i})
 \end{aligned}$$

H_0
(reduced representation of backbone)

AMH Ramachandran Potentials

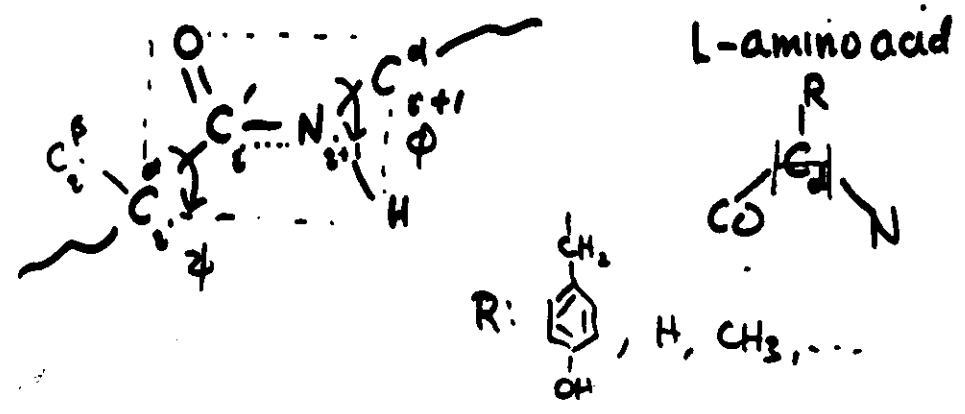


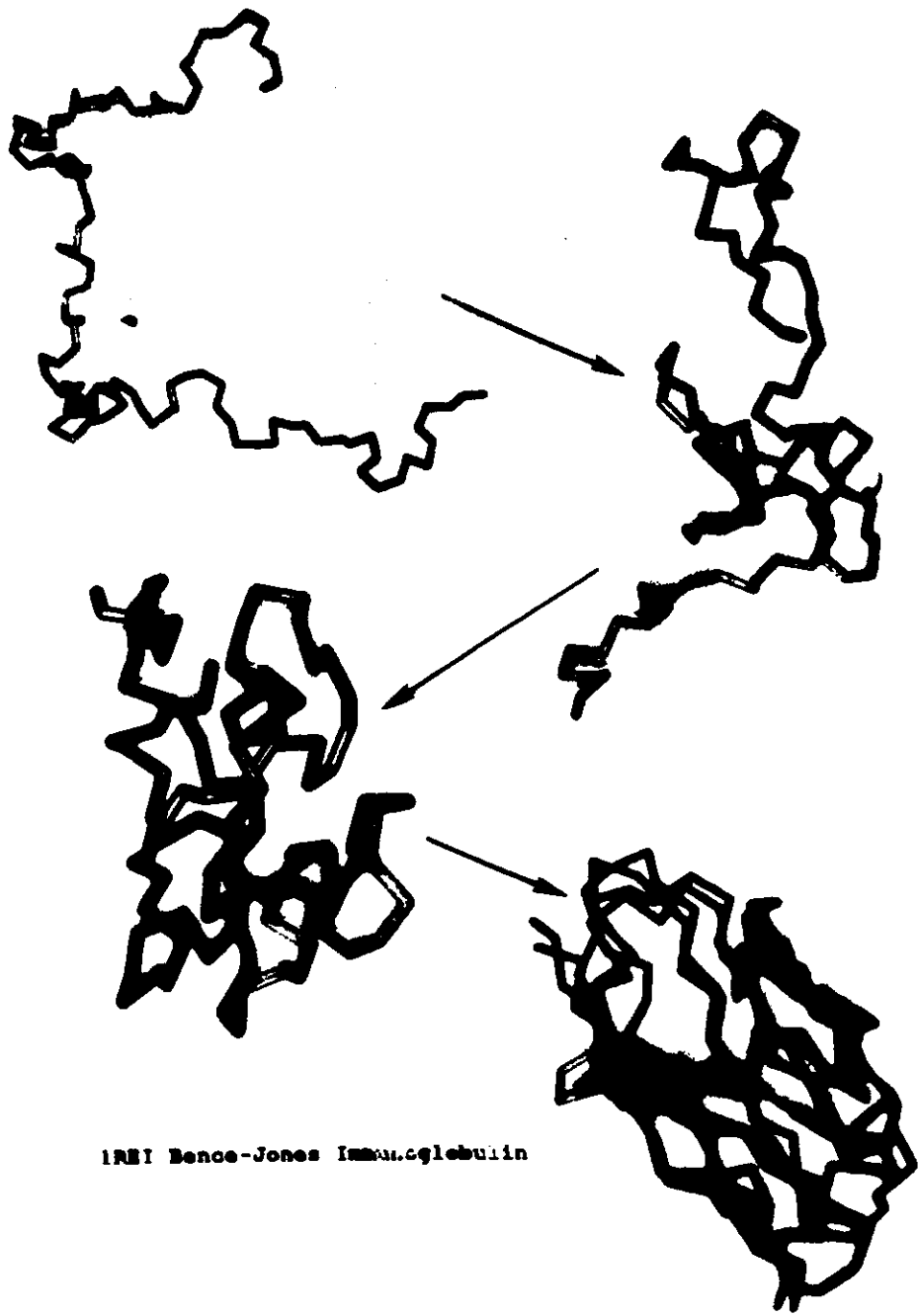
original potential



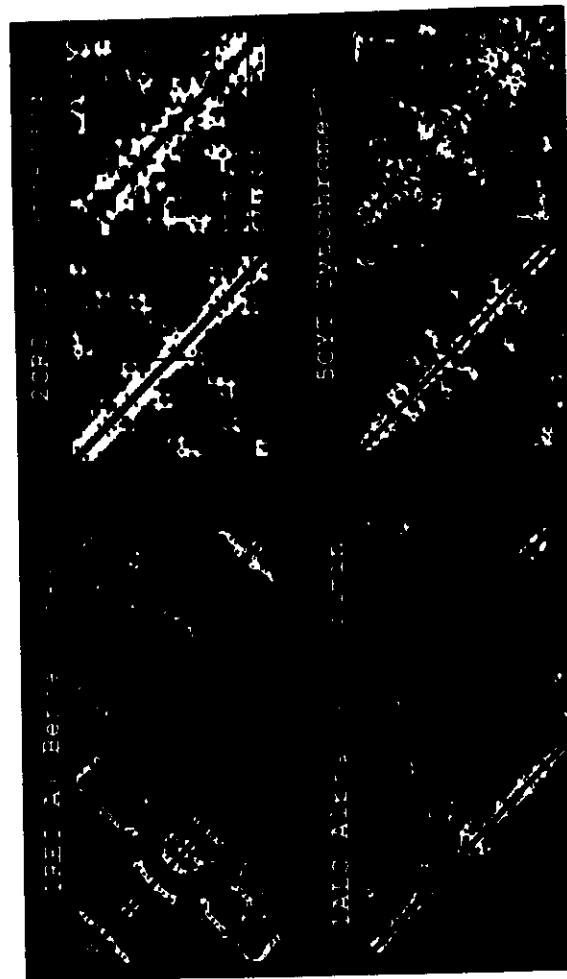
current potential

Planar Peptide Bond Geometry





1RBI Bence-Jones Immunoglobulin

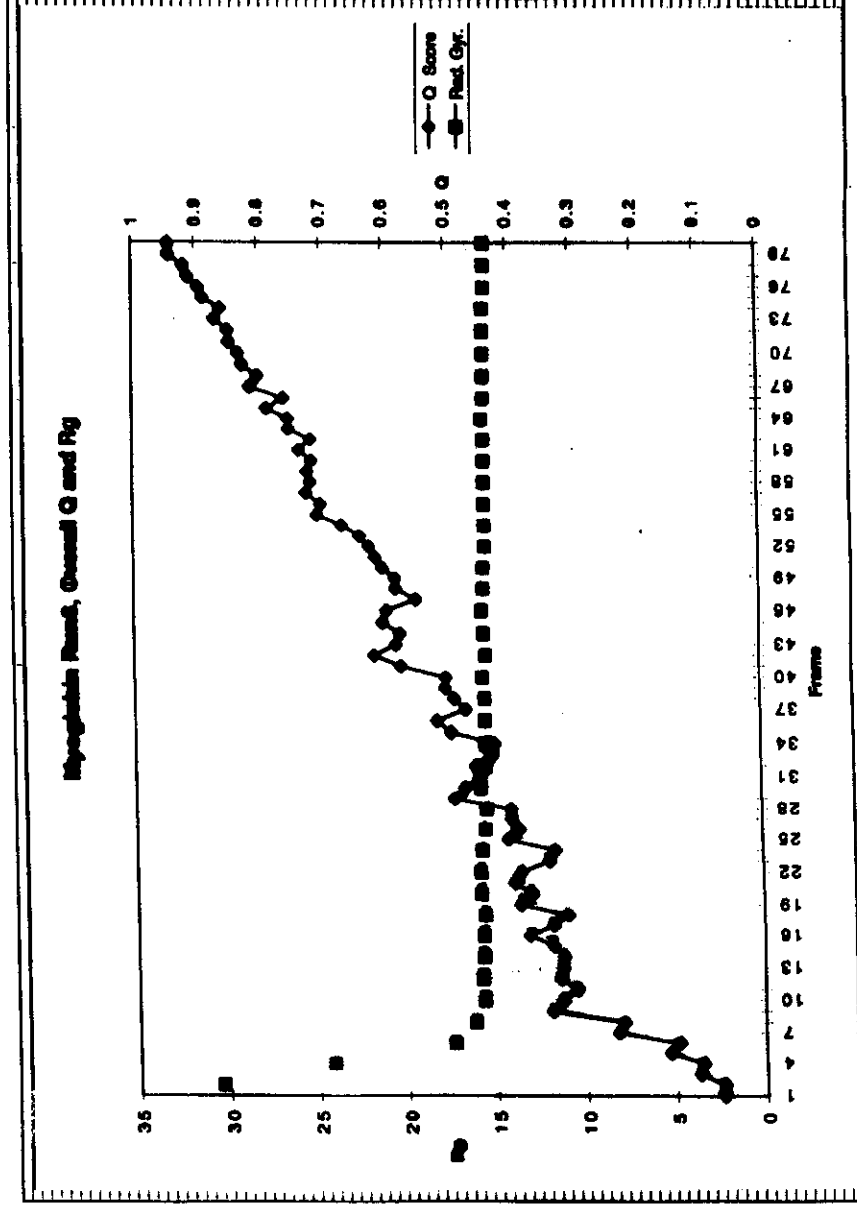


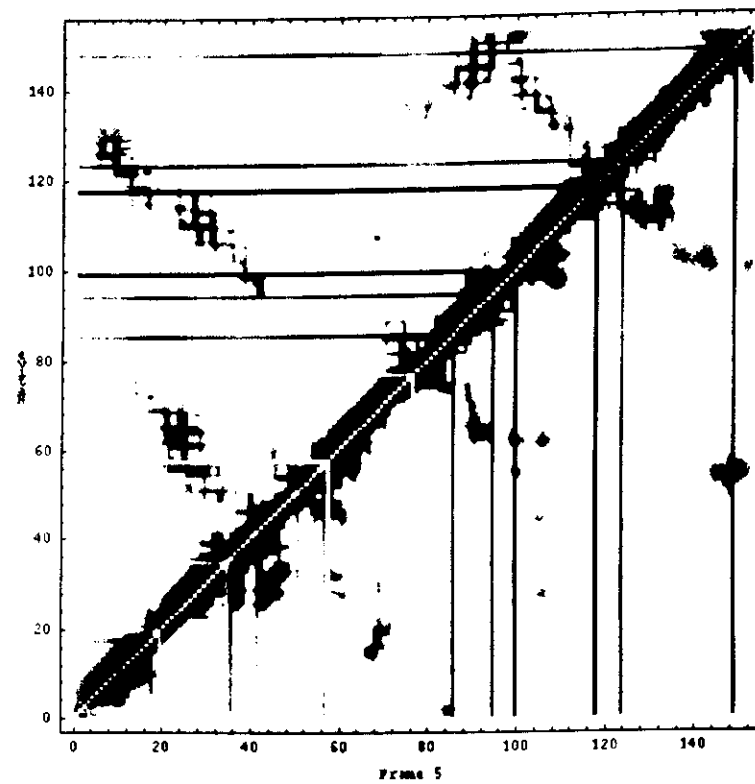
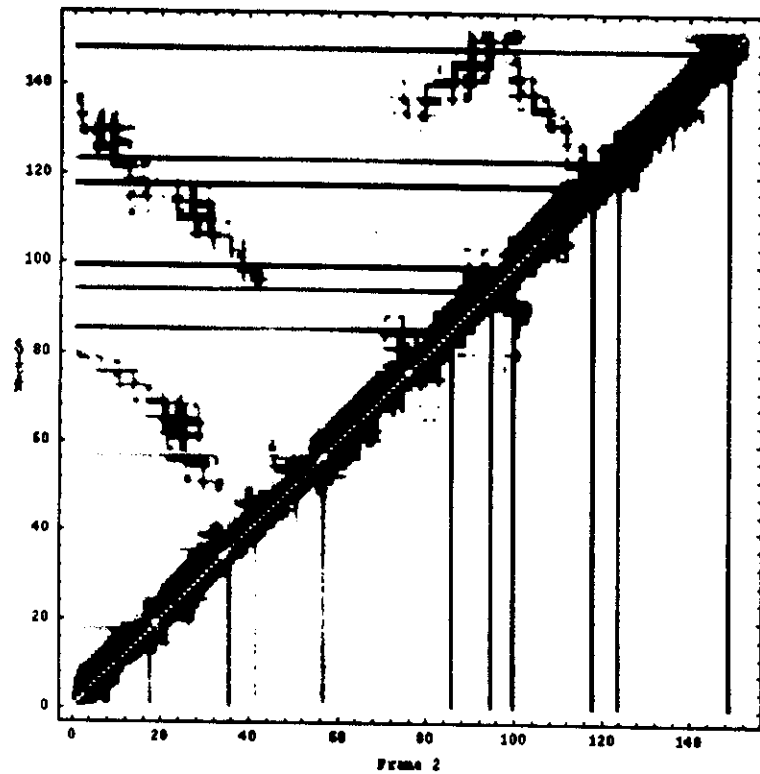
(1 homolog)

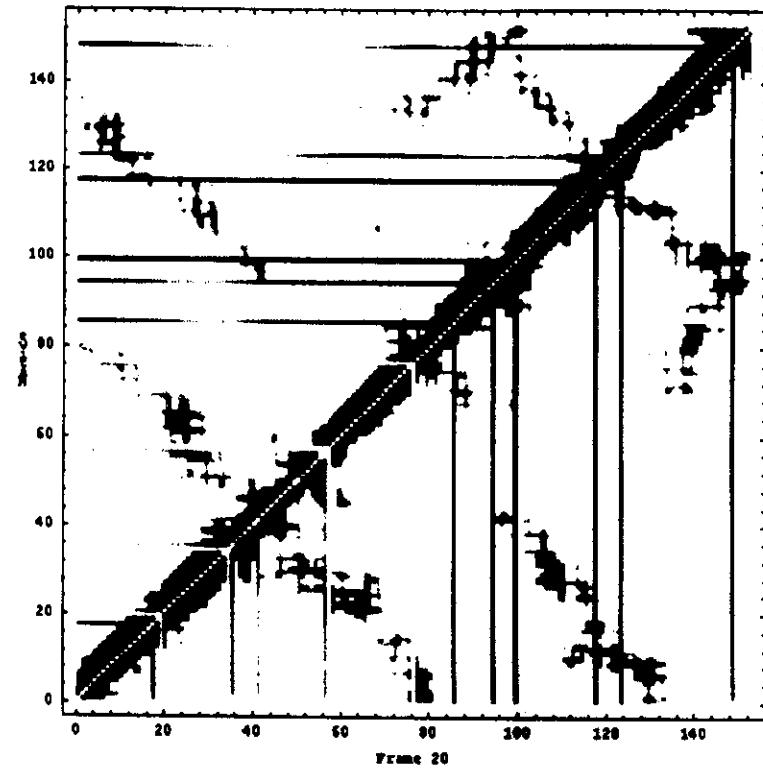
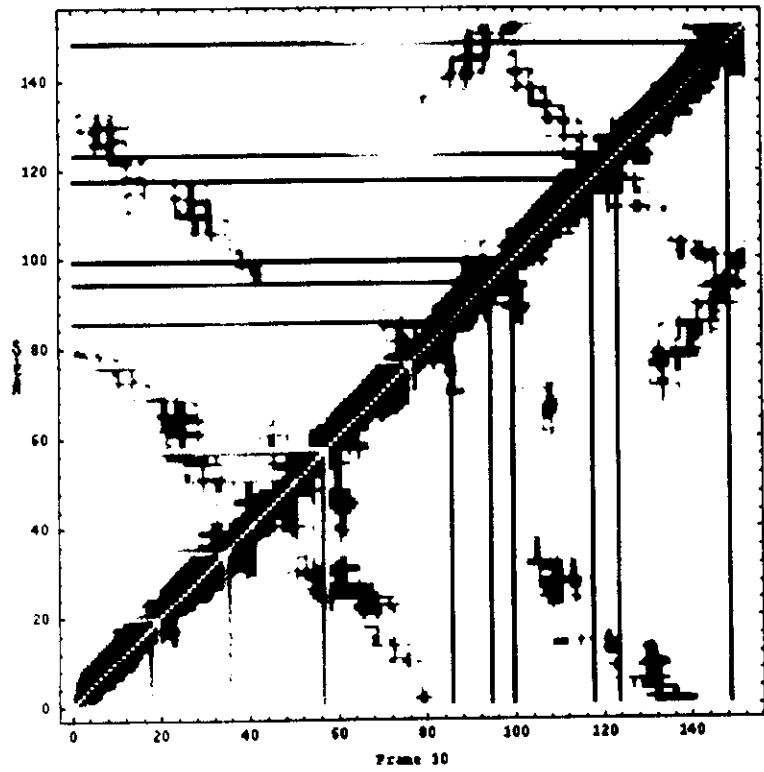
Molecular Dynamics Results

Target		Homolog				Predicted structure	
PDB	Name	PDB	Name	% I	q	q	RMS
Alpha							
2C80	204 Coo	1LBD	Lambda Repressor	16.4	0.46	4.21	6.06 0.22 6.41
2BUC	<i>P. carolinensis</i> Cytochrome C50A	3CXC	<i>E. rubrum</i> Cytochrome C5	20.2	0.47	3.40	6.15 0.25 6.19
3P28	Human Co-Binding	3CLN	Rat Calmodulin	21.3	0.48	3.35	9.86 0.23 10.14
2C7T	Furin Cytochrome C	3CXC	<i>E. rubrum</i> Cytochrome C5	20.2	0.46	4.21	
Beta							
1M8I	Human Insulin (Variable)	3HFM	Mouse F5b (H)	20.2	0.49	2.86	2.68 0.59 3.10
2BIA	Human Elastase (beta-2)	1MCP	Mouse F5b (L)	20.5	0.48	2.89	5.15 0.37 6.81
3PCY	Poplar Plastocyanin	1PAZ	<i>A. foveata</i> Plastocyanin	25.8	0.55	3.73	7.34 0.31
Other							
1ALC	Raboon Alpha-lactalbumin	2LYZ	Hen Egg White Lysozyme	20.7	0.53	1.96	7.58 0.47 7.81
209H	Barley Chymotrypsin Inhibitor II	2BEC	Leach Eglin-C	37.1	0.69	1.64	5.54 0.41 6.43 8.72

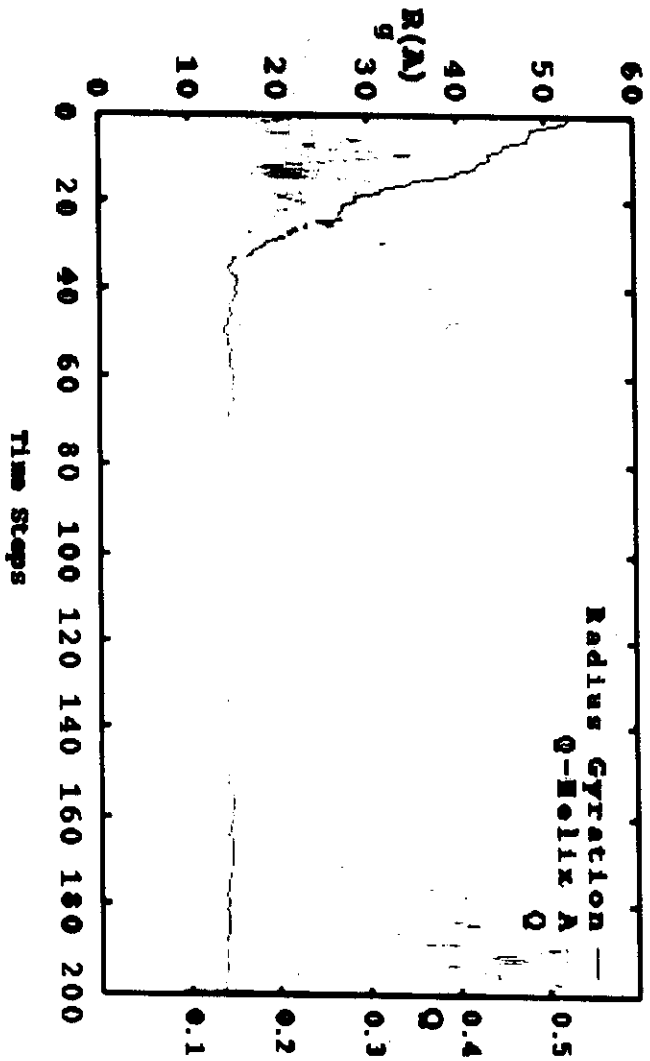
Q_table_my03



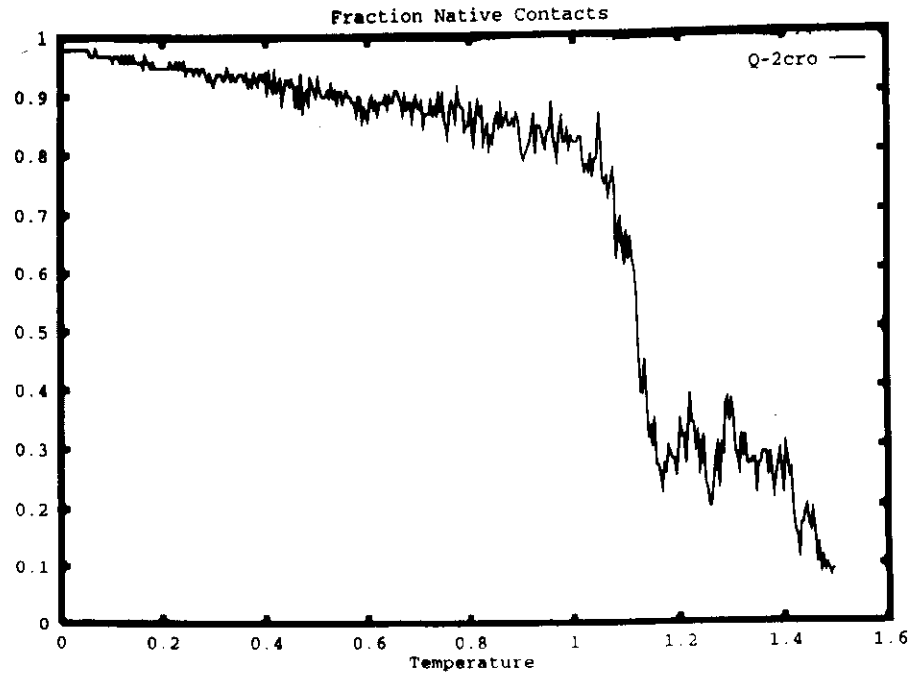




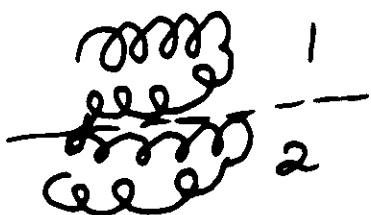
New ψ, ϕ Potential



Multiple Hemeris 2AM



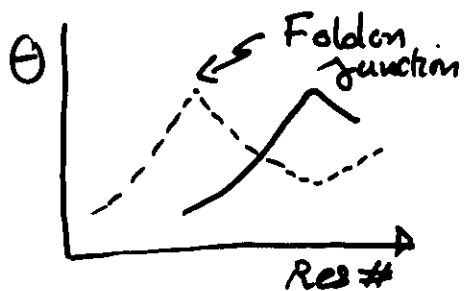
Foldons and Funnels



Is there a (sub)domain structure in Folding Kinetics?

Maximize over cuts

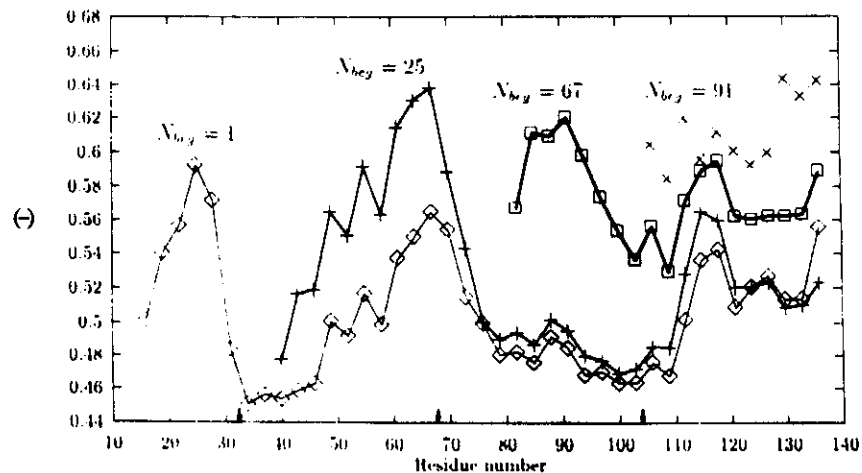
$$\Theta = \frac{\delta E_s^{(1)}}{\sqrt{N_1} \Delta E^{(1)}} + \frac{\delta E_s^{(2)}}{\sqrt{N_2} \Delta E^{(2)}}$$



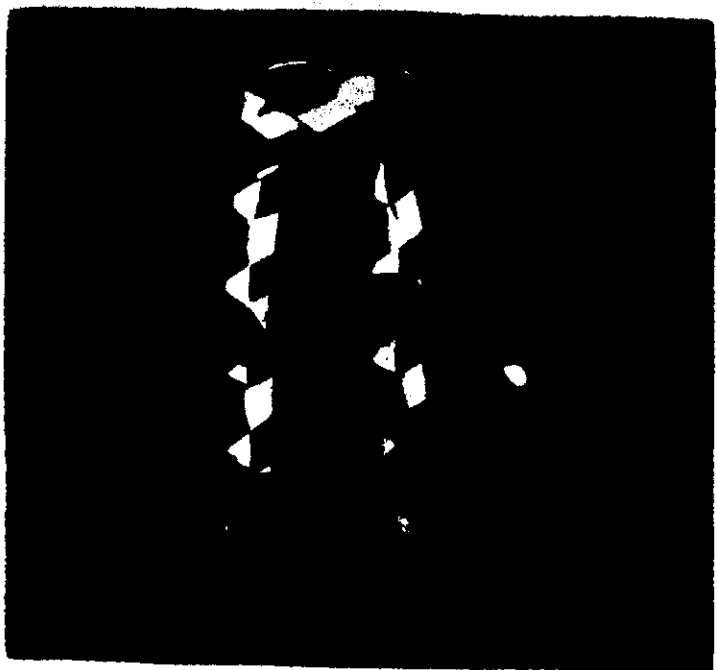
Rough Division into "Fast" Folding Contiguous Segments - "Foldons"

Biophysics: Panchenko et al.

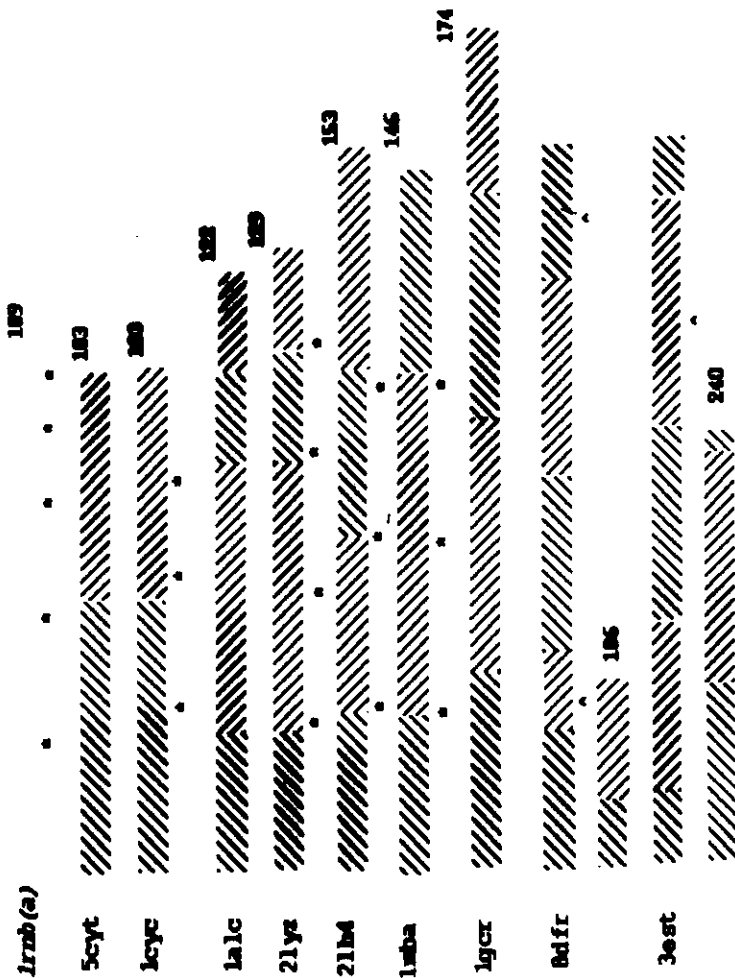
Proc. Natl. Acad. Sci. USA 93 (1996)



Comparing "Foldons" and Exons



Regions of maximum $\left(\frac{\delta E_s}{\Delta E}\right)_{local}$



Panchenko, Z. Schulten, P. Wolynes, PNAS (1995)

Table 1.

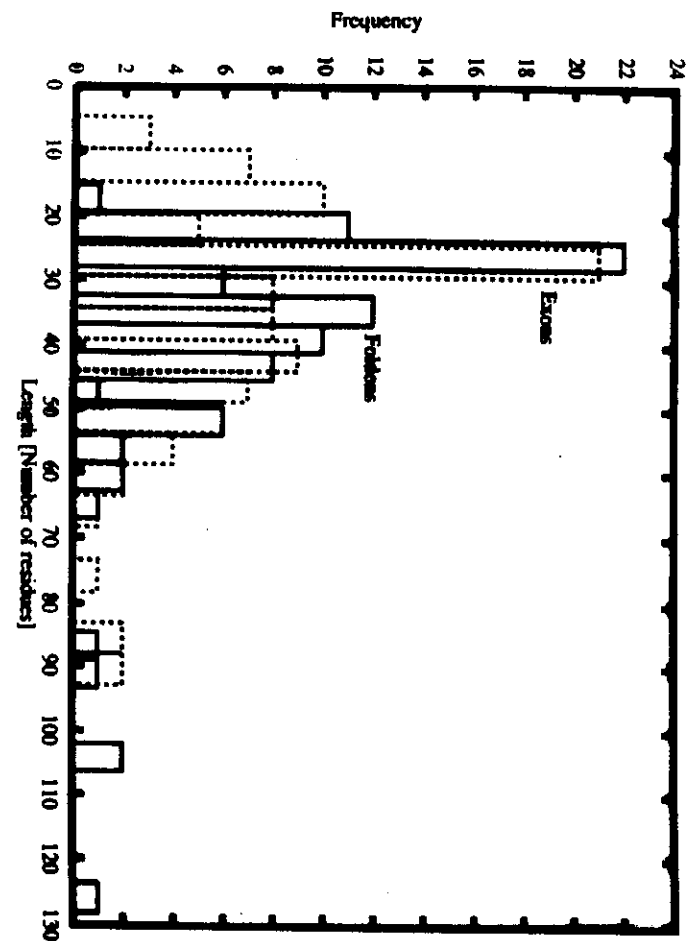
Foldons, exons and structural domain boundaries for some representative proteins: myoglobin (MBA), triosephosphate isomerase (TIM), lysozyme (LYZ), dihydrofolate reductase (DFR) and carboxypeptidase (CPA). Faldon positions are defined according to BPTT criterion, ones that satisfied the C1-2 criterion are indicated by asterisk.

Faldon junctions	Intron positions	Domain junctions	Faldon junctions	Intron positions	Domain junctions
Myoglobin			Lysozyme		
28	31	30 [†] , 34 [§]	28*	28	30 [†] , 38*, 41 [§]
64*		66 [†] , 67 [§]	55		55 [†] , 61 [§]
91*				82	84 [†] , 86 [§] , 88*
	105	102 [§] , 101 [†]	103*	108	109 [†]
Triose-phosphate isomerase			Dihydrofolate reductase		
22*	13	28 [†]	28	28	34*
	37	37 [†]		45	
64*	61	63 [†] , #	121	80	
	78	78 [†]	163	122	134*
88*		93 [†]		161	
	107	109 [†]	Carboxypeptidase		
121		125*, 127 [†]	22*	17	29 [§]
	132		55*	59	50*
	151	147 [†]		81	73 [§] , 96*
163	167	162 [†]	112*	126	122 [§] , 125*
	180/183	179 [†]	163*	156	141*, 144 [§]
202	209	201 [†] , 210*	193		189 [§] , 196*
	237/239	230 [†]	226*	215	211 [§] , 234 [§]
			271	260	271 [§]

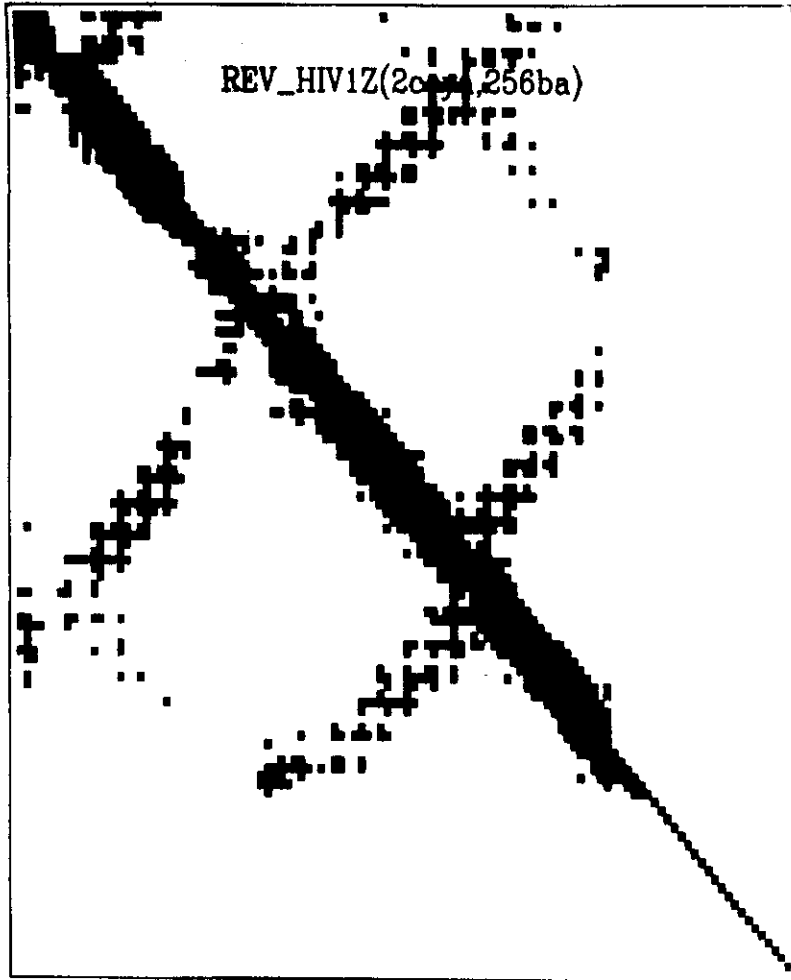
† - structural module boundaries calculated by C^α distance algorithm [8,10,33,31]

§ - domain boundaries defined by geometrical algorithm of Rose [6]

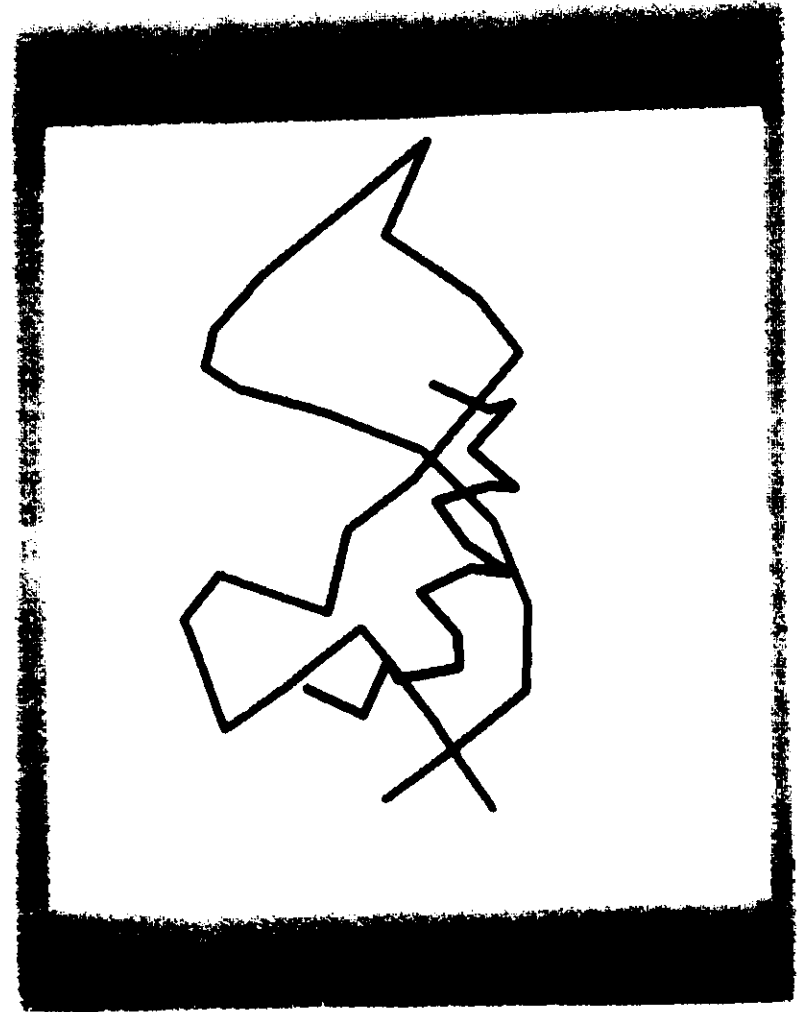
- domain boundaries found by surface accessibility algorithm [9]

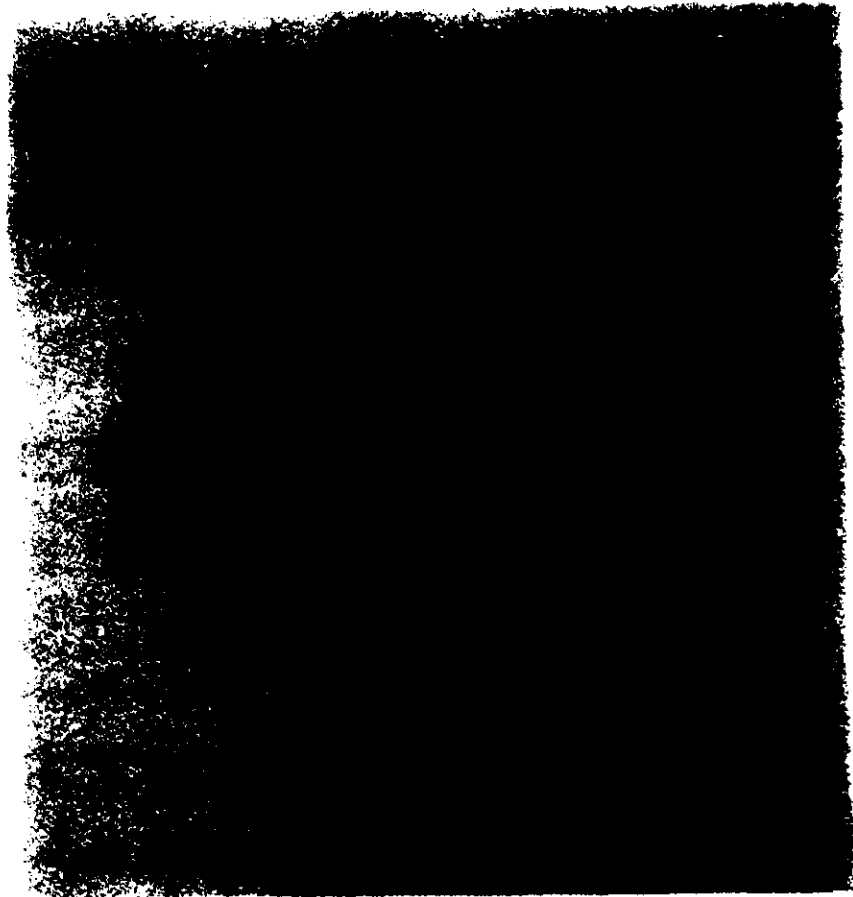


Helical Threading



REV Regulatory Protein in HIV
Prediction: 3-helix Bundle





C. Brooks, et. al
(pre-print) 1995

