### <u>Molecular Biophysics & Biochemistry</u> 400a/700a (Advanced Biochemistry)

#### Computational Aspects of: 1, Secondary Structure Prediction; 2, Protein Packing

Mark Gerstein

#### Classes on 9/24/98 & 9/29/98 Yale University

#### The Handouts I

#### • Notes

♦ this class (sec. str.) and next (packing)

#### • Presentation Paper

- Frishman D, and Argos P. (1997) The Future of Protein Secondary Structure Prediction Accuracy. *Folding & Design* 2:159-62.
- ♦ Controversial idea: secondary structure prediction to 80%?
- http://bioinfo.mbb.yale.edu/~mbg/clippings/frishman-fad-acc-secstr.pdf (guest:guest)

#### Secondary Structure Review

- ♦ Handout from D Frishman
- Garnier, J., Gibrat, J. F. & Robson, B. (1996b). GOR method for predicting protein secondary structure from amino acid sequence. *Methods Enzymol* 266, 540-53.
- Problem Set Paper (Sec. Struc.)
  - King, R. D. & Sternberg, M. J. E. (1996). Identification and application of the concepts important for accurate and reliable protein secondary structure prediction. *Prot. Sci.* 5, 2298-2310.

### The Handouts II

#### • Problem Set Paper (Packing)

- Joan Pontius, Jean Richelle, Shoshana J. Wodak (1996). Deviations from Standard Atomic Volumes as a Quality Measure for Protein Crystal Structures. Journal of Molecular Biology **264**: 121-136.
- http://bioinfo.mbb.yale.edu/~mbg/clippings/wodak-jmb-volume.pdf (guest:guest)
- Packing Review
  - M Gerstein & F M Richards. "Protein Geometry: Distances, Areas, and Volumes," (eventually) to appear in *International Tables for Crystallography*. (International Union of Crystallography, Chester, UK).
  - http://bioinfo.mbb.yale.edu/e-print/geom-inttab/
- Optional Fun Reading (Only via Web)
  - ◊ Barry Cipra (1998). "Packing Challenge Mastered At Last," Science 281: 1267
  - http://www.sciencemag.org/cgi/content/full/281/5381/1267
  - Simon Singh (1998). "Mathematics 'Proves' What the Grocer Always Knew," New York Times (August 25).
  - http://bioinfo.mbb.yale.edu/~mbg/clippings-u/nyt-sci-packproof.txt

### <u>The Raw Material</u> Molecular Biology Information - DNA

#### • Raw DNA Sequence

- ♦ Coding or Not?
- ♦ Parse into genes?
- ◊ 4 bases: AGCT
- ◇ ~1 K in a gene, ~2 M in genome

atggcaattaaaattggtatcaatggttttggtcgtatcggccgtatcgtattccgtgca gcacaacaccgtgatgacattgaagttgtaggtattaacgacttaatcgacgttgaatac atggcttatatgttgaaatatgattcaactcacggtcgtttcgacggcactgttgaagtg aaaqatqqtaacttaqtqqttaatqqtaaaactatccqtqtaactqcaqaacqtqatcca gcaaacttaaactggggtgcaatcggtgttgatatcgctgttgaagcgactggtttattc ttaactgatgaaactgctcgtaaacatatcactgcaggcgcaaaaaaagttgtattaact qqcccatctaaaqatqcaacccctatqttcqttcqtqqtqtaaacttcaacqcatacqca ggtcaagatatcgtttctaacgcatcttgtacaacaaactgtttagctcctttagcacgt gttgttcatgaaactttcggtatcaaagatggtttaatgaccactgttcacgcaacgact qcaactcaaaaaactqtqqatqqtccatcaqctaaaqactqqcqcqqcqqcqqcqqtqca tcacaaaacatcattccatcttcaacaqqtqcaqcqaaaqcaqtaqqtaaaqtattacct gcattaaacggtaaattaactggtatggctttccgtgttccaacgccaaacgtatctgtt gttgatttaacagttaatcttgaaaaaccagcttcttatgatgcaatcaaacaagcaatc aaaqatqcaqcqqaaqqtaaaacqttcaatqqcqaattaaaaqqcqtattaqqttacact gaagatgctgttgtttctactgacttcaacggttgtgctttaacttctgtatttgatgca gacgctggtatcgcattaactgattctttcgttaaattggtatc . . .

### <u>Molecular Biology Information:</u> <u>Protein Sequence</u>

#### • 20 letter alphabet

- ♦ ACDEFGHIKLMNPQRSTVWY but not BJOUXZ
- Strings of ~300 aa in an average protein (in bacteria), ~200 aa in a domain

#### ~200 K known protein sequences

dldhfa_	LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMTTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr	LNSIVAVCQNMGIGKDGNLPWPPLRNEYKYFQRMTSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_	ISLIAALAVDRVIGMENAMPWN-LPADLAWFKRNTLNKPVIMGRHTWESI
d3dfr	TAFLWAQDRDGLIGKDGHLPWH-LPDDLHYFRAQTVGKIMVVGRRTYESF
d1dhfa_	LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMTTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr	LNSIVAVCQNMGIGKDGNLPWPPLRNEYKYFQRMTSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_	ISLIAALAVDRVIGMENAMPW-NLPADLAWFKRNTLDKPVIMGRHTWESI
d3dfr	TAFLWAQDRNGLIGKDGHLPW-HLPDDLHYFRAQTVGKIMVVGRRTYESF
dldhfa_	VPEKNRP LKGRINLVLS RELKEP PQGA HFLSRSLDDALKLTE QPELANKVD MVWIVGGSSVYKEAM NHP
d8dfr	VPEKNRPLKDRINIVLSRELKEAPKGAHYLSKSLDDALALLDSPELKSKVDMVWIVGGTAVYKAAMEKP
d4dfra_	G-RPLPGRKNIILS-SQPGTDDRV-TWVKSVDEAIAACGDVPEIMVIGGGRVYEQFLPKA
d3dfr	PKRP <b>LPERTNVVLT</b> HQEDYQ <b>AQGA-VVVHDVAAVFAYAK</b> QHLDQ <b>ELVIAGGAQIFTAFK</b> DDV
d1dhfa_	-PEKNRP LKGRINLVLS RELKEP PQGA HFLSRSLDDALKLTE QPELANKVD MVWIVGGSSVYKEAM NHP
d8dfr	-PEKNRP <b>LKDRINIVLS</b> RELKEA <b>PKGAHYLSKSLDDALALLD</b> SPELKSKVD <b>MVWIVGGTAVYKAAM</b> EKP
d4dfra_	-GRPLPGRKNIILSSSQPGTDDRV-TWVKSVDEAIAACGDVPE IMVIGGGRVYEQFLPKA
d3dfr	-PKRP <b>LPERTNVVLT</b> HQEDYQ <b>AQGA-VVVHDVAAVFAYAK</b> QHLDQ <b>ELVIAGGAQIFTAFK</b> DDV

#### <u>Molecular Biology Information:</u> <u>Macromolecular Structure</u>

- DNA/RNA/Protein
  - ♦ Almost all protein

(RNA Adapted From D Soll Web Page, Right Hand Top Protein from M Levitt web page)



"Identity elements' in Eschetichia coli glutamine tRNA.





### Molecular Biology Information: **Protein Structure Details**

#### Statistics on Number of XYZ triplets

- 200 residues/domain -> 200 CA atoms, separated by 3.8 A
- Avg. Residue is Leu: 4 backbone atoms + 4 sidechain atoms, 150 cubic A

- => ~1500 xyz triplets (=8x200) per protein domain
- ♦ 10 K known domain, ~300 folds

ATOM	1	С	ACE	0	9.401	30.166	60.595	1.00	49.88	1GKY	67
ATOM	2	0	ACE	0	10.432	30.832	60.722	1.00	50.35	1GKY	68
ATOM	3	CH3	ACE	0	8.876	29.767	59.226	1.00	50.04	1GKY	69
ATOM	4	Ν	SER	1	8.753	29.755	61.685	1.00	49.13	1GKY	70
ATOM	5	CA	SER	1	9.242	30.200	62.974	1.00	46.62	1GKY	71
ATOM	б	С	SER	1	10.453	29.500	63.579	1.00	41.99	1GKY	72
ATOM	7	0	SER	1	10.593	29.607	64.814	1.00	43.24	1GKY	73
ATOM	8	CB	SER	1	8.052	30.189	63.974	1.00	53.00	1GKY	74
ATOM	9	OG	SER	1	7.294	31.409	63.930	1.00	57.79	1GKY	75
ATOM	10	Ν	ARG	2	11.360	28.819	62.827	1.00	36.48	1GKY	76
ATOM	11	CA	ARG	2	12.548	28.316	63.532	1.00	30.20	1GKY	77
ATOM	12	С	ARG	2	13.502	29.501	63.500	1.00	25.54	1GKY	78
• • •											
ATOM	1444	CB	LYS	186	13.836	22.263	57.567	1.00	55.06	1GKY1	510
ATOM	1445	CG	LYS	186	12.422	22.452	58.180	1.00	53.45	1GKY1	511
ATOM	1446	CD	LYS	186	11.531	21.198	58.185	1.00	49.88	1GKY1	512
ATOM	1447	CE	LYS	186	11.452	20.402	56.860	1.00	48.15	1GKY1	513
ATOM	1448	NZ	LYS	186	10.735	21.104	55.811	1.00	48.41	1GKY1	514
ATOM	1449	OXT	LYS	186	16.887	23.841	56.647	1.00	62.94	1GKY1	515
TER	1450		LYS	186						1GKY1	516



# Explonential Growth of Data Matched by Development of Computer

- CPU vs Disk & Net
  - $\diamond$  As important as the increase in computer speed has been, the ability to store large amounts of information on computers is even more crucial

Num.

Protein

Domain

Structures

**Driving Force in** • **Bioinformatics** 

> (Internet picture adapted from D Brutlag, Stanford)



8

#### **Overview**

#### • Why interesting?

- ◊ Not tremendous success, but many methods brought to bear.
- What does difficulty tell about protein structure?
- Start with TM Prediction (Simpler)
- Basic GOR Sec. Struc. Prediction
- Better GOR
  - ♦ GOR III, IV, semi-parametric improvements, DSC
- Other Methods
  - ◊ NN, nearest nbr.



Credits: Rost et al. 1993; Fasman & Gilbert, 1990

- Not Same as Tertiary Structure Prediction -- no coordinates
- Need torsion angles of terms + slight diff. in torsions of sec. str.

	TRUE	154152	COLUMN	TIERRO	Dis Devices	COLUMN TO C	State of the state
CONTRACTOR OF	100.0	1122	1000	COPIE ST	OTHER DESIGNATION.	office where	ALC: N
1.1.1	CERET	<b>EARENE</b>	LINESE	T	NUMBER OF TAXABLE PARTY	COLUMN TO A	darm -
THEY HERE	EVER IS	CONTRACTOR OF	DEPER	Tion of	Distance in the local	COLUMN TO A	
					4		
AA IPPLING	TAXABLE IN CONTRACTOR	APPROX NO.	SCHOPPER L	ABIIGNES	EPHRAPYMACC	0.779	VLIBER I
-861 E	TELLE 11	100.00	100004	a setu	IN NOTION OF COLUMN	COORE IN	COLUMN 1
ALC: NO.		DOD BOD B	10000	1.000	THE R. LEWIS CO.	110.00	COLUMN 1
THE COLOR	SALES I	TRADE OF STREET	10.80305	A MER IN	A NOT STREET, SALES	CO BOOK OF	COLUMN 1
NOT STREET	ADDRESS R	COLUMN STATES	10.0010-0	1.00	NO YOU WANTED	a sink of	10000
	*********				l		
A HOLEND	DO-BHEFTIDO	Real Provides	CONTRACTOR OF THE	DATAMELO	OTPETLAPETS	LINGY	NAUD I
2041 68	LEE	6.835	120		AND INCOME.		100.0
CHIRAL BREAK	1 1000	100	No. of Concession, Name	DITT	100		INVESTIGATION OF
THINEER	CREE	PERCENT.	TC 8		THE REPORT OF THE		100
NOI DONE	LOCAL DOCK OF	10100	0822	100	10110		100
1000						-	
U. I MARLING	<b>LITERAACYP</b>	PPT-MODPIG	<b>LTER</b> LYES	10/10/10	NERGON AND 1 M		
OC LEVENSES	OR MONTH AND I	HINK	NAME AND ADDRESS OF		100 Million Colored	and the second second	
THE ROOMAGE	ADD NOT TO A	14	KENERGING		A TRANSPORT		
TH DEEPER	TT DE TE DE	1.00.0	OTHER DOCTOR	a	IN MARKING STATE	and the second second	
10 HOURSEN	OR ADDRESS OF		the set of some		in many shares	Concernance of the local division of the loc	
						-	-
	WHECKMAN	THEFT					B
A DESCRIPTION							
A IFORM	10.00						

(a) Residue-by-residue comparison of experimentally observed (005) and predicted (COM<sup>32</sup>, ETH<sup>36</sup>, PHD (Ref. 35 and B. Rost and C. Sonder, submitted) structures of the catalytic subunt of the coMP-dependent protein kinase (3cpk), Wi is the armon acid sequence taken from Protein Data Bank entry 1cpk inesidues 27-2871. Secondary structure: H = whelk, E = (5-sheet leatenake), trank = loop. Predicted or-helices and (5-strands that have insufficient overlap with an observed segment of the same type are underlined. Note the relatively good prediction of the location of segments for the ETH and PHD methods and overprediction of ehelices for the COM method.





(b) Ribbon view of the domain used in this blind test. The X-ray is structure of catalytic suburit of the cAMP-dependent protein kinase. Drawn using Molecopt<sup>24</sup>.



Figure 1 Column model for the core of the reaction center from Rep. windls. Reproduced, with permission, from Ref. 18.

		Some TM s	scales:		
F	-3.7	CEC		I	4.5
М	-3.4			V	4.2
I	-3.1			L	3.8
L	-2.8			F	2.8
V	-2.6			С	2.5
С	-2.0			М	1.9
W	-1.9			А	1.8
A	-1.6			G	-0.4
Т	-1.2			Т	-0.7
G	-1.0			W	-0.9
S	-0.6			S	-0.8
P	+0.2			Y	-1.3
Ŷ	+0.7			Ρ	-1.6
- Н	+3.0			Η	-3.2
0	+4.1			Ε	-3.5
∑ N	+4.8			Q	-3.5
E	+8.2			D	-3.5
ĸ	+8.8			Ν	-3.5
D	+9.2			Κ	-3.9
R	+12.3			R	-4.5

### How to use GES to predict proteins

- Transmembrane segments can be identified by using the GES hydrophobicity scale (Engelman et al., 1986). The values from the scale for amino acids in a window of size 20 (the typical size of a transmembrane helix) were averaged and then compared against a cutoff of -1 kcal/mole. A value under this cutoff was taken to indicate the existence of a transmembrane helix.
- H-19(i) = [ H(i-9)+H(i-8)+...+H(i) + H(i+1) + H(i+2) + . . . . + H(i+9) ] / 19

#### Graph showing Peaks in scales



Figure 3.12. Representative profiles of three membrane proteins used to predict membrane-spanning helices. The amino acid scales of Kyte-Doolittle (804), Goldman-Ergelman-Steitz (GES) (389), and Rao-Angos (1194) were used. A computer solware package (SEQANAL) provided by Dr. A. Crofts (Uaiv, of Illinois) was used to pretente duste profiles. For comparative purposes, the Kyte-Duolittle and GES plots were obtained using a window of 19 residues and then smoothed using a second pass with a window of 7. The average value at each residue position is plotted as a function of residue nearbor starting with the amino terminus on the left in each case. The values plotted for the Kyte-Doolittle and GES scales represent average hydropethy and transfer free energy per residue (local/mol). The Rao-Argos plot used a span of 7 residues and was structured used two additional passes with the same span of 7, as recommended by the authors. The scale values reflect the relative preference for being in a membrane-spanning helin. Note that the version of the GES algorithm which was used does not take into account possible ion pair formation. See text for details. Illustrations Adapted From: von Heijne, 1992; Smith notes, 1997



(a)



#### **Removing Signal sequences**

 Initial hydrophobic stretches corresponding to signal sequences for membrane insertion were excluded. (These have the pattern of a charged residue within the first 7, followed by a stretch of 14 with an average hydrophobicity under the cutoff).



### Relation between Energy Scales and Statistics

- Statistical Mechanics: Energies give Probabilities of Observing an Object
- Inverse Statistical Mechanics: Derive Energies from observing probabilities in the database
- E = = kT In P
- dE = E1 E2 = -kT ln P1/P2
- dE = E(state-1) E(~state-1)
- $dE = -kT \ln P(state-1)/P(~state-1)$
- P(state-1)/P(~state-1) = odds ratio => dE = lod score

### Ex. Pr(S) probability that residue j has secondary structure i

- Problem of DB Bias
- f(A) = frequency of residue
  A to have a helical conf. in
  db
- f(A,i) = f(A) at position i in a particular sequence
- E(α)=statistical energy of helix over a window
- p(i, α) = probability that residue i is in a helix

$$E_{\alpha} = \sum_{i}^{N} \ln f_{\alpha}^{i}$$

$$p_{\alpha}^{i} = \frac{e^{-Ea/RT}}{\sum_{j}^{N} e^{-E_{j}/RT}}$$

### Statistics Based Methods: Persson & Argos

 Propensity P(A) for amino acid A to be in the middle of a TM helix or near the edge of a TM helix

 $P(A) = \frac{\frac{n(A, \text{TM})}{\sum_{A} n(A, \text{TM})}}{\frac{n(A, \text{everywhere})}{\sum_{A} n(A, \text{everywhere})}}$ 



Figure 1. Positional amino acid compositions of transmembrane segments. But that showing the amine acid compositions for 15.N and Ceterminal positions relative to the centre of partitive transmembrane segments listed in feature tables of the Swim-Prot database. For each position, the percentage contribution of each amine acid shown according to the hydrophilic (top) to hydropholic (bottom) order, given in the ruler but at the left. The hydropholic residue contributions are illustrated in white, the hydrophilic is dark-gray, and intermediate in light gray. The compositions of positions 11 to 15 at the N-terminal side and 12 to 15 at the C-terminal side differ significantly from the others, especially for the most hydropholic and charged, hydropholic is residues. These results suggest that in general transmembrane spans consist of a hydropholoke parties 21 residues in length.

Illustration Credits: Persson & Argos, 1994

# Refinements: Charge on the Outside, Positive Inside Rule

 for marginal helices, decide on basis of R+K inside (cytoplasmic)



Credits: von Heijne, 1992



Figure 4. (a) Hydrophobicity plot for the SecY protein. The upper and lower cutoffs are marked. A tentative transmembrane segment with a mean hydrophobicity falling between the 2 cutoffs is marked by an arrow. (b) Two possible topologies for the SecY protein based on the hydrophobicity plot. The putative transmembrane segment is shown in black. The number of Arg + Lys vesidues is shown next to each polar segment. Note that the correct alternative (bottom, including the putative ransmembrane segment) has a much higher charge-bias han the incorrect one.



- How to train to find right threshold? Not that many TM helices
- Marginal TM helices are not that hydrophobic but 1/3 of TM's are very hydrophobic, so focus on these.
- Sosui, Klein & Delisi, Boyd
- Discriminant analysis: set threshold to be best partition of dataset

## <u>Secondary Structure Prediction:</u> <u>GOR is based on Info. Theory</u>

- Analogous to P&A, but more data and more involved stats
- Based on information theory, want to maximize the "Information"
  - I(x;y) = information that event y carries about the occurrence of x
  - $\Diamond I(x;y) = \log (P(x|y)/P(x))$
  - $\Diamond$  I = 0 for no information
  - $\Diamond$  I > 0 if A favors helix
  - $\Diamond$  I < 0 if disfavors
- Information sort of like entropy I increases, entropy increases

### **GOR: Simplifications**

- For independent events just add up the information
- I(S<sub>j</sub>; R<sub>1</sub>, R<sub>2</sub>, R<sub>3</sub>,...R<sub>last</sub>) = Information that first through last residue of protein has on the conformation of residue j (S<sub>j</sub>)
  - Could get this just from sequence sim. or if same struc. in DB (homology best way to predict sec. struc.!)
- Simplify using a 17 residue window: I(S<sub>j</sub>=H ; R[j-8], R[j-7], ...., R[j], .... R[j+8])
- Difference of information for residue to be in helix relative to not: I(dSj;y) = I(Sj=H;y)-I(Sj=~H;y)

 $\diamond$  odds ratio: I(dSj;y)= In P(Sj;y)/P(~Sj;y)

I determined by observing counts in the DB, essentially a lod value 21

<u>Basic</u> <u>GOR</u>

- Pain & Robson, 1971; Garnier, Osguthorpe, Robson, 1978
- I ~ sum of I(Sj,R[j+m]) over 17 residue window centered on j and indexed by m
  - ◊ I(Sj,R[j+m]) = information that residue at position m in window has about conformation of protein at position j
  - ◊ 1020 bins=17\*20\*3
- In Words
  - Secondary structure prediction can be done using the GOR program (Garnier et al., 1996; Garnier et al., 1978; Gibrat et al., 1987). This is a well-established and commonly used method. It is statistically based so that the prediction for a particular residue (say Ala) to be in a given state (i.e. helix) is directly based on the frequency that this residue (and taking into account neighbors at +/- 1, +/- 2, and so forth) occurs in this state in a database of solved structures. Specifically, for version II of the GOR program (Garnier et al., 1978), the prediction for residue i is based on a window from i-8 to i+8 around i, and within this window, the 17 individual residue frequencies (singlets)<sub>22</sub>

### Directional Information

helix

#### strand

**Table 3.** Directional informational parameters:  $h(S) = x_i x^i$ ; Rj + m) for residue position versus residue type for  $\alpha$ -helices<sup>4</sup>

	1-8	i-7	1-6	1-5	1-4	1-3	1-2	i-1	i	i+1	1+2	i+3	1+4	1+5	1+6	1+7	i+8
8	19	21	22	24	34	36	44	47	60	60	53	50	44	40	31	23	24
c	-47	-45	-44	-47	-44	-36	-44	-55	-56	-58	-54	-55	-58	-58	-59	-53	-66
đ	14	15	14	15	17	21	15	17	-7	-11	-31	-42	-28	-12	-8	1	-5
	14	16	15	20	26	27	34	52	62	57	32	15	19	12	6	7	9
£	-19	-14	-10	-4	-2	-1	6	-1	10	10	12	12	-4	-5	2	0	2
g	5	2	1	-5	-22	-30	-50	-70	-92	-52	-28	-21	-13	-17	-8	-6	-6
h	-22	-20	-9	-10	-19	-10	-14	-7	-11	-4	0	-3	-2	2	6	11	12
i	7	7	0	0	1	1	2	-5	1	2	1	7	-6	-3	10	8	6
x	-2	-1	-1	-1	-6	-9	-6	5	17	17	21	27	35	33	21	22	23
1	0	-1	0	6	9	16	30	33	45	47	51	53	37	32	30	25	18
ㅋ	4	3	15	23	30	30	39	36	45	54	57	53	44	29	30	14	1
n	2	3	2	-5	-9	-10	-16	-17	-31	-16	-17	-16	-9	-8	-9	-10	-5
P	-12	-15	-14	-19	-23	-25	-30	-48	-82	-195-	-145	-104	-67	-49	-43	-33	-17
q	-4	3	7	4	13	8	10	24	35	32	31	21	18	18	9	8	6
r	5	3	6	13	3 7	13	19	27	34	32	36	41	33	29	23	21	18
	-10	-7	-10	-10	-16	-17	-25	-21	-39	-35	-39	-41	-32	-35	-34	-35	-33
t	1	-1	-6	-8	-6	-11	-16	-25	-48	-47	-48	-46	-34	-31	-34	-26	-24
v	-5	-12	-13	-14	-13	-19	-17	-20	-15	-22	-22	-20	-26	-19	-15	-10	-5
w	0	-4	-12	-19	-7	14	16	12	18	17	12	8	1	-6	1	3	-13
Y	-22	-19	-17	-20	-16	-21	-30	-32	-8	-10	-4	-12	-17	-9	-10	-14	-15

"Note that the convention used is the reverse of that adopted by (Gamier et al., 1978), for example the first entry for alanine at position j-8 is the amount of information that an alanine residue eight positions toward the N terminus has for predicting at α-helix

Table 4.	Directional	informational	parameters.	for	residue position	versus	residue	type	for	B-strands	ļ
----------	-------------	---------------	-------------	-----	------------------	--------	---------	------	-----	-----------	---

	1-8	1-7	1-6	1-5	1-4	i-3	1-2	i-1	i	i+1	1+2	1+3	1+4	1+5	1+6	1+7	1+8	_
a	-8	-7	-13	-17	-23	-33	-26	-32	-43	-37	-30	-30	-26	-27	-26	-25	-25	
	1 3	13	-9	-20	-15	-3	9	33	47	51	21	19	9	-5	7	-5	-14	
ā	-7	-5	0	-9	-4	-14	-42	-73	-83	-59	-21	10	22	24	16	11	13	
	-14	-5	-5	-11	-21	-27	-45	-44	-57	-54	-46	-29	-25	-12	-12	-2	D	
ž	-9	-20	-32	-34	-30	-12	24	44	49	39	24	2	-9	-23	-24	-29	-23	
a	-3	9	24	29	34	30	18	-23	-48	-27	6	27	39	38	33	23	23	
ĥ	6	11	17	22	12	16	0	-2	3	-2	5	3	8	4	-1	1	-3	
1	-21	-30	-31	-21	-12	-3	26	58	76	64	33	11	-14	-24	-20	-14	-11	
k	20	12	15	14	8	4	-8	-14	-25	-40	-39	-27	-20	-24	-20	-15	-15	
1	-2	-10	-18	-27	-30	-27	-6	15	27	21	2	-19	-31	-29	-28	-26	-25	
-	-22	-26	-29	-40	-31	-17	-7	23	24	28	17	2	-15	-31	-53	-36	-16	
n	1	8	14	5	0	-6	-30	-65	-62	-28	-6	11	18	21	16	10	3	
	9	7	12	24	20	8	-22	-65	-108	-64	-8	17	25	30	32	31	21	
a .	6	12	8	16	8	-5	-22	-27	-30	-52	-49	-34	-22	-17	-9	2	20	
-	0	8	3	-3	5	2	1	-14	-26	-32	-30	-35	-27	-26	-25	-25	-21	
	16	14	17	19	14	5	-3	-13	-15	-4	15	27	32	32	31	28	21	
ŧ.	6	8	14	15	16	21	19	25	31	22	13	9	12	25	34	34	34	
÷	1	-11	-15	-11	4	25	51	75	91	81	49	19	-6	-12	-16	-11	-11	
÷.	-8	-8	-28	-19	-9	5	23	44	45	30	13	-18	-22	-40	-15	-7	-9	
Y	13	13	4	14	12	20	24	37	48	31	20	-1	2	11	7	٥	-4	
																		_

#### coil

	1-8	1-2	1-6	1-5	1-4	1-3	1-2	1-1	1	1+1	1+2	1+3	1+4	1+5	1+6	1+7	148
	-12	-15	-12	-12	-17	-13	-15	-24	-32	-35	-33	-29	-24	-20	-13	+5	-6
•	36	26	41	50	4.5	31	25	1.9	7	5	27	25	38	4.8	42	45	59
4	-8	+10	-13	-8	-13	-14	12	25	50	43	35	27	7	-7	-4	-9	-5
	-3	-11	-19	-11	-10	-7	-5	-23	-26	-23	-3		-1	-3		-5	-9
£	32	35	28	25	21	. 2	-23	-34	-43	-40	-25	-12		20	3.3	1.6	13
	-3	-8	-18	-17	-7	- 2	26	68	97		19	-3	-18	-14	-18	-11	-11
h.	15	,	-4	-7		+3	12			. 8	-4	1	-3	-5	-5	+10	-9
4	. 7	13	1.9	1.4		1	-31	-43	-66	-55	-26	-14	2.4	1.8	4	- 2	1
x	-12	-7	-10	-9	-1	5	11	5	. 6	. 9	5	-8	-20	-15	-7	-10	-12
1	- 2		31	11	11	- 2	-23	-42	-65	-63	-52	-39	-15	-11	-10	-6	0
-	11	14	4	- 3	-3	-16	-33	-52	-62	-77	-71	-14	-32	-7	3	. 9	
-	-3	-8	-11	1		13	32	51	61	31	18	. 6	-6	-0	-6	- 2	2
P	4		4	-1	5	15	39	. 26	128	159	.98	59	32	17	11	3	. 0
9	-1	-11	-12	-15	-17	-4	5	-5	-13	1	- 1	1	-2	-5	-1	- 9	-30
=	-4	-9.	-8	-10	-10	-13	-18	+16	-14	-3	-14	-15	-14	-11	-5	-3	-2
	-3	-4	-4	-4	4	11	22	26	41	31	20	13	3	5	4		11
E	-5	-5	-5	-4	-7	-5		- 2	15	21	29	30	1.9	7	3	-4	-5
w.	3	17	20	28		-2	-26	-66	-50	-51	-20	3	25	24	23	15	11
÷.	5		28	28	12	-16	-32	-65	-53	-30	-20	5	13	38	. 9	2	16
7	10		12	7	6	1	7	-1	-31	-14	-11	11	13	1	3	32	35

Credits: King & Sternberg, 1996

**Table 3.** Directional informational parameters: I(Sj = x:x': Rj + m) for residue position versus residue type for  $\alpha$ -helices<sup>\*</sup>

Types of		i-8	i-7	i-6	i-5	i-4	<b>i-</b> 3	i-2	i-1	i	. +1	i+2	i+3	i+4	i+5	i+6	i+7	i+8
<u>1 ypc3 01</u>	a	19	21	22	24	34	36	44	47	60	10	53	50	44	40	31	23	24
<b>—</b> • • •	с	-47	-45	-44	-47	-44	-36	-44	- 5	-56	- 5 B	-54	-55	-58	-58	-59	-53	-66
Raciduac	d	14	15	14	15	17	21	15	17	-7	-11	-31	-42	-28	-12	-8	1	-5
	e	14	16	15	20	26	27	34	5.7	62	57	32	15	19	12	6	7	9
	Ľ	-19	-14	-10	-4	-2	-1	- 6	-1	10	10	12	12	-4	-5	2	0	2
	g	5	2	1	-5	-22	-30	-50	-70	-92	-52	-28	-21	-13	-17	-8	-0	-0
	п	- 44	- 70	-9	-1	1	-10	-14	-/	-11	-4	1	- 3	- 6	4	10		12
		, '	_1	_1	_1	- 6		- 5	-5	17	17	21	27	25	33	21	22	23
		-2	-1	-1	-1	-0	16	30	22	45	47	51	53	37	12	30	25	10
		4		15	23	30	30	30	36	45	54	57	53	44	29	30	14	1
	1 7	2	3	2	-5	-9	-10	-16	-17	-31	-16	7	-16	-9	-8	-9	-10	-5
		-12	-15	-14	-19	-23	-25	-30	-48	-82	-195	-14	-104	-67	-49	-43	-33	-17
	a	-4	-3	7	4	13	8	10	24	35	32	- 51	21	18	18	9	8	6
	r	5	3	6	1:	3 7	13	19	27	34	32	36	41	33	29	23	21	18
	s	-10	-7	-10	-10	-16	-17	-25	-21	-39	-35	-39	-41	-32	-35	-34	-35	-33
	t	1	-1	-6	-8	-6	-11	-16	-25	-48	-47	-48	-46	-34	-31	-34	-26	-24
	v	-5	-12	-13	-14	-13	-19	-17	-20	-15	-22	-22	-20	-26	-19	-15	-10	-5
	w	0	-4	-12	-19	-7	14	16	12	18	17	12	8	1	-б	1	3	-13
	У	-22	-19	-17	-20	-16	-21	-30	-32	-8	-10	-4	-12	-17	-9	-10	-14	-15

Credits: King & Sternberg, 1996

<sup>a</sup>Note that the convention used is the reverse of that adopted by (Garnier et al., 1978), for example the first entry for alanine at position j-8 is the amount of information that an alanine residue eight positions toward the N terminus has for predicting an  $\alpha$ -helix at position j.

- Group I favorable residues and Group II unfavorable one:
- A, E, L -> H; V, I, Y, W, C -> E; G, N, D, S -> C
- P complex; largest effect on proceeding residue
- Some residues favorable at only one terminus (K)

# <u>GOR III</u>

- Improvements in GOR -- Full GOR decomposition
- I(S<sub>j</sub>; R<sub>1</sub>, R<sub>2</sub>, R<sub>3</sub>,...R<sub>last</sub>) = Information that first through last residue of protein has on the conformation of residue j (S<sub>j</sub>)

 $\diamond\,$  looked at singlets, now pairs, eventually triplets....

- GOR III
- I(Sj; R[j+m], R[j]) = information that pair of residues at postions 0 & m in window has about conformation of protein at position j
  - ◊ 16 pairs =>16\*20\*20\*3=19200 bins

## <u>GOR IV</u>

 I(Sj; R[j+m], R[j+n]) = the frequencies of all 136 (=16\*17/2) possible di-residue pairs (doublets) in the window.

◊ 20\*20\*3\*16\*17/2=163200 pairs

- Parameter Explosion Problem: 1000 dom. struc. \* 100 res./dom. = 100k counts, over how many bins
- Dummy counts for low values (Bayes)

#### <u>Assessment</u>

THE GOR METHOD

	TABLE II
GLOBAL	<b>RESULTS FOR DATABASE PREDICTION</b>

- Q3 + other assess, 3x3
- Q3 = total number of residues predicted correctly over total number of residues
- GOR gets 65%
  - \$\\$ sum of diagonal over total number of residue -- (14K+5K+21K)/ 64K
- Under predict strands & to a lesser degree, helices: 5.9 v 4.1, 10.9 v 10.6

	shah.	Obse		
	H	Е	С	Total
Predicted	h			
н	14,460	3094	4790	22.344
E	1124	4965	2089	8178
С	6002	5546	21,496	33.044
Total	21,586	13,605	28,375	63,566
Qprd "	64.7	60.7	65.1	
Qobs b	67.0	36.5	75.8	
$Q_3^c = 64.4\%$				

" Number of correctly predicted residues/number of predicted residues.

\* Number of correctly predicted residues/number of observed residues

<sup>c</sup> Total number of correctly predicted residues/total number of residues.

Credits: Garnier et al., 1996

# Training and Testing Set

 Cross Validation: Leave one out, seven-fold



Figure 2. Comparison of prediction accuracy (correctly predicted residues as a proportion of total residues) versus effective number of parameters for linear-logistic models (number of parameters ≤640) and penalized likelihood models for crossvalidated (�) and

uncrossvalidated ( $\Box$ ) results. The values of the penalty parameter  $\lambda$  are shown.

#### Credits: Munson, 1995; Garnier et al., 1996

		_	TABLE	I		
		I	DATABASE PRO	DTEINS <sup>a</sup>		
1aaj.x	laak.x	laan a	Áchar			
lads.x	lalk.a	1807.2	laba.x	labk.x	1abm.a	ladd y
lavh.a	lavh.x	1bab a	14pa.x	lapm.e	larb.x	latr x
1bll.e	1bmd.a	1bov a	IDDh.a	1bbp.a	1bet.x	lbre a
1caj.x	lcau.a	Icau b	торь.х	1brs.d	1btc.x	1c2r a
1chm.a	1cmb.a	Icob a	Icde.x	lcdt.a	1cew.i	lcet y
lcrl.x	1cse.i	letf v	Icol.a	lcpc.a	1cpc.b	lenty
1dog.x	1dsb.a	leaf v	Ictm.x	lcus.x	lddt.x	Idhr y
1fba.a	1fdd.x	1fba y	leco.x	lede.x	lend.x	lena a
1fxi.a	lgal.x	ladl o	Ina.a	lfkb.x	1fna.x	1 fpr v
1gof.x	lgox.x	lgul.o	lgdh.a	lgky.x	lglt.x	lamfa
1hdx.a	1hiv.a	1gp1.a	Igpb.x	1gpr.x	lgsr.a	1bbo x
1hrh.a	1hsl.a	10.8	Ihle.a	1hmy.x	1hoe.x	1hoja
1129.x	lle4.x	llen o	lifc.x	lipd.x	1isu.a	lith a
1lts.a	lits.d	1mda v	liga.a	1lis.x	Illa.x	11mh 3
1mpp.x	1mup x	1mac.x	1mgn.x	1min.a	1min.b	1mio.s
1nxb.x	lofy.x	lolb a	Inba.a	1ndk.x	1noa.x	Inch o
1pda.x	lpfk.a	1010.a	lomf.x	lomp.x	lonc.x	1000.4
1plf.a	1poc.x	1pgb.x	1pgd.x	1phh.x	1php.x	1084.8
1ppn.x	lprc.c	1pon.x	lpox.a	1ppa.x	1ppf.e	1pn.x
1pya.b	1pvd a	iprc.n	lprc.l	1prc.m	1pts.a	
1rve.a	1s01.x	100.8	lrec.x	1rib.a	1rnd.x	Ipya.a
1shf.a	1sim x	Isac.a	lsbp.x	1ses.a	1set.x	lebo o
1tca.x	Itie x	ISILD	lsnc.x	1spa.x	1stf i	15112.2
ltro.a	1tth a	1tml.x	1tnd.a	ltpl.a	ltrb x	Itoe.a
1wht.b	1wsv a	lutg.x	lvaa.a	lvaa.b	lymo a	Turk.a
2aza.a	2bon a	Twsy.b	1yhb.x	1zaa.c	256b a	Iwnt.a
2cpl.x	200p.a	2ccy.a	2cdv.x	2chs.a	2cmd x	2001.0
2hbg.x	2hhm a	ZCIC.X	2cts.x	2cyp.x	2dnia	2cp4.x
2mhr.x	2mpr v	2nip.a	2hpd.a	2ihl.x	2lh2.x	2er/.e
2pia.x	2nol a	2msb.a	2mta.c	2mta.h	2mta I	211V.X
2sas.x	2scn a	2por.x	2reb.x	2rn2.x	2rsl.a	2p11.x
2tpr.a	2tsc.a	Zsga.x	2sn3.x	2spc.a	2tgi x	2sar.a
3chy.x	3cla x	Jaan.a	3aah.b	3adk.x	3b5c.x	Zund.a
3ink.c	3rub I	SCOX.X	3dfr.x	3eca.a	3gan a	3cha u
4enl.x	4fef x	Srub.s	3sdh.a	3tgl.x	451c.x	Ablm a
5tim.a	6fab b	4gcr.x	4ts1.a	4xis.x	5fbn a	40im.a
8atc.b	8cat a	OIAD.I	6taa.x	8abp.x	8acn x	Sp21.x
9wga.a		611D.X	8rxn.a	8tln.e	9ldt a	oalc.a
5						STULX

<sup>a</sup> he database was prepared by J. M. Levin and checked for homologous sequences with the help of V. Di Francesco. This database has been modified to restore the total length of the sequences as defined in the SEQRES field of the Protein Data Bank (PDB) file (the DSSP program omits residues whose coordinates are missing in the PDB file, and thus if this occurs in the middle of the polypeptide chain it is split into two or more chains). Residues having no coordinates were assigned the conformation X and were not taken into account for the prediction accuracy although the prediction was done with the whole sequence length. The PDB code is followed by the chain name a, b, c, d, h (heavy), l (light), x (one chain only), e (enzyme), or i (inhibitor).

#### Is 100% Accuracy Possible?

#### **Quoted from Barton (1995):**

One problem that has arisen is how to evaluate secondary structure predictions. For prediction of a single protein sequence one might expect the best residue by residue accuracy to be 100%. It is not possible to define the secondary structure of a protein exactly, however. There is always room for alternative interpretations of where a helix or strand begins or ends so failure of a prediction to match exactly the secondary structure definition is not a disaster [24]. The problem of evaluation is more complicated for prediction from multiple sequences, as the prediction is a consensus for the family and so is not expected to be 100% in agreement with any single family member. The expected range in accuracy for a perfect consensus prediction is a function of the number, diversity and length of the sequences. Russell and I have calculated estimates of this range [11].

Simple residue by residue percentage accuracy has long been the standard method of assessment of secondary structure predictions. Although a useful guide, high percentage accuracies can be obtained for predictions of structures that are unlike proteins. For example, predicting myoglobin to be entirely helical (no strand or coil) will give over 80% accuracy but the prediction is of little practical use. Rost *et al.* [25] and Wang [26] explore these problems and suggest some alternative measures of predictive success based on secondary structure segment overlap. Although such measures help in an objective assessment of the prediction, there is no complete substitute for visual inspection. By eye, serious errors stand out and predictions of structures that are unlike proteins are usually recognizable. By eye, it is also straightforward to weight the importance of individual secondary structures. For example, prediction of what is in fact a core strand to be a helix would seriously hamper attempts to generate the correct tertiary structure of the protein from the predicted secondary structure, whereas prediction of a non-core helix as coil may have little impact on the integrity of the tertiary structure.

## <u>'Predictable' regions of secondary</u> <u>structure</u>

#### **Quoted from Barton (1995):**

When recent predictions are examined in the light of the corresponding experimentally determined structures, the results look good. In general, the regions predicted with the highest confidence measure are also the most accurate. For example, Livingstone and I [13] assigned 41% of the tyrosine phosphatase structure with high confidence. Within these regions 88% of the residues were correctly predicted. Interestingly, these figures agree with Rost and Sander's observation that 40% of a sequence will be predicted with >88% accuracy by their method [1]. This agreement suggests that there is a core of 'predictable' regions in a protein. Examination of six blind predictions shows that the most accurately predicted regions are those that have clear periodicity in conservation, where conserved positions either alternate (beta-strand) or have a 1, 4, 5, 8 pattern characteristic of one face of an alpha-helix (CD Livingstone, personal communication). Problems remain with buried alpha-helices that comprise short runs of conserved hydrophobic amino acids. These often look like potential beta-strands and can mislead both automatic and manual predictive methods.

### <u>Types of Secondary Structure</u> <u>Prediction Methods</u>

- Parametric Statistical
  - ◊ struc. = explicit numerical func. of the data (GOR)
- Non-parametric
  - ◊ struc. = NON- explicit numerical func. of the data
  - ◊ generalize Neural Net, seq patterns, nearest nbr, &c.
- Semi-parametric: combine both
- single sequence
- multi sequence
  - ◊ with or without multiple-alignment

## <u>GOR Semi-</u> <u>parametric</u> <u>Improvements</u>

 $[\neg a, c, \neg c, a, a, c, \neg a] \rightarrow c$   $[\neg a, \neg a, c, b, *, \neg b] \rightarrow c$   $[\neg a, a, c, c, a, a, \neg b, \neg a] \rightarrow c$   $[\neg a, *, *, a, b] \rightarrow b$   $[\neg a, *, *, a, c] \rightarrow c$   $[a, *, *, a, c, *, \neg c] \rightarrow c$   $[a, *, *, a, a, a, c, \neg a] \rightarrow c$   $[c, b, t, a, a, a, a] \rightarrow c$   $[c, *, a, a, \neg a, a] \rightarrow c$ 

a =  $\alpha$ -helix, b  $\beta$ -strand, c = coil, \* = vildcard ( $\alpha$ -helix or  $\beta$ -strand or coil)  $\neg$  = not.

If the pattern ( ) the left is met in a predict on, then the secondary structure in bold on the left is rewritten as the secondary structure on the right of the rule. For example:

 Filtering GOR to regularize

 $[b, b, b, \mathbf{a}, \mathbf{c}] \rightarrow [b, b, b, \mathbf{c}, \mathbf{c}]$  $[b, b, \mathbf{c}, \mathbf{a}, \mathbf{c}] \rightarrow [b, b, \mathbf{c}, \mathbf{c}, \mathbf{c}]$  $[b, b, b, \mathbf{a}, b, b, b] \rightarrow [b, b, b, \mathbf{b}, \mathbf{b}, \mathbf{b}, \mathbf{b}]$ 

Illustration Credits: King & Sternberg, 1996

# <u>Multiple</u> <u>Sequence</u> <u>Methods</u>

- Average GOR over multiple seq. Alignment
- The GOR method only uses single sequence information and because of this achieves lower accuracy (65 versus >71 %) than the current "state-ofthe-art" methods that incorporate multiple sequence information (e.g. King & Sternberg, 1996; Rost, 1996; Rost & Sander, 1993).

Illustration Credits: Livingston & Barton, 1996



FiG. 5. Conservation analysis of the 17 flavodoxin sequences clustered in Fig. 3. The Taylor Venn diagram was used (Fig. 1) with a threshold of T = 7. See text for details.







- GOR parms
- + simple linear discriminant analysis on:
  - Ist from C-term, N-term
  - ◊ insertions/deletes
  - ◊ overall composition
  - hydrophobic moments
  - ◊ autocorrelate: helices
  - conservation moment

#### **Predator**

Figure 5. Pairwise local alignments of the query sequence 0 with the related sequences m=1,2...M. Every alignment is characterized by its length and residue percentage identity.

Sequence 0	 :	<b>_</b>	:	. ·	:	:
			:			
Sequence 1	 		 			•
Sequence 2						1
Sequence 3						
Sequence M						

#### Predator

- Frishman & Argos (1997) proposed an alternate way to utilize the additional information contained in a set of related sequences. A careful pairwise alignment of the query sequence with all related sequences is performed.
- NNSSP, Salamov & Solovyev, 1995
  - Issues in segments, Dist metric betw. them (use eisenberg env.), k-closest neighbors, >71% accuracy
- Yi & Lander, 1994; Presnell et al., 1992

Illustration Credits: D Frishman handout

## <u>Neural</u> <u>Networks</u>

Figure 1. Function of a perceptron, the simplest neural network. A simple perceptron has only 1 output unit (black). Each of the left nodes receives a certain input signal (e.g. binary, i.e. =0 or 1). All units are connected to the output node by the junctions  $J^1$ , with e.g.  $J^1_{1J}$ connecting input unit j with output unit 1. The contribution of each left node (e.g. the jth) to the signal arriving at the right one is a product of the strength of the junction connecting the 2 units, and the input: e.g. J1/47. All products (here 3) are summed by the right node (here s1). This sum is then evaluated by a non-linear trigger function. The resulting map of the sum onto an interval between 0 and 1 is the actual output of the network. The broken-line nodes show a potential extension of the perceptron to a 2-layered feed-forward network. Stippled circles, input units, signal = 1 or 0. Black circle, output unit. Step 1, the input to this unit is summed according

$$h_{i}^{1} = \sum_{j=1}^{N^{2}-1} J_{ij}s_{j}^{0} \quad (\text{here}, \ i=1)$$

Step 2, the output from this unit is computed by a sigmoid trigger function:

$$a_i^z = \frac{1}{1 + \exp((-b_i^z))}$$
.

Broken-line circles, the potential extension to a 2-layered feed-forward network.



- Somehow generalize and learn patterns
- Black Box
- Rost, Kneller, Qian....
- Perceptron (above) is Simplest network
  - Multiply junction \* input, sum, and threshold


- Hidden Layer
- Learning
  - Steepest descent to minimize an error function
- Jury Decision
  - $\diamond\,$  Combine methods
  - Escape initial conditions

#### Yet more methods....

#### struc class predict

- ◊ Vect dist. between composition vectors
- threading via pair pot
- seq comparison
- ab initio from md
- ab initio from pair pot.

## Mail Servers and Web Forms

			Source
			code
Method	URL	Institution	Availability
ANTHE-		Institute of Biology and	
DDOT		Chemistry of Proteins	
PROT	nttp://www.ibcp.fr/antheprot.ntml (currently unreachable)	(Lion) Revier Cellege of	YES
PSSP	http://dot.imgen.bcm.tmc.edu:9331/pssprediction/pssp.html	Medicine (Houston)	NO
		Imperial Cancer	
		Research Center	
DSC	http://bonsai.lif.icnet.uk/bmm/dsc/dsc_form_align.html	(London)	YES
		University of	
GOR	nttp://molbiol.soton.ac.uk/compute/GOR.ntml	Southampton	NO
		Linivorsity of California	
nnPredict	http://www.cmpharm.ucsf.edu/~nomi/nnpredict.html	(San Francisco)	NO
Dradiat		(Garri Taricisco)	NO
Predict-			
Protein	http://www.embl-heidelberg.de/predictprotein/predictprotein.html	EMBL (Heidelberg)	NO
PRED-			
ATOR	http://www.embl-heidelberg.de/argos/predator/predator_form.html	EMBL (Heidelberg)	YES
PSA	http://bmerc-www.bu.edu/psa/	BioMolecular Engineering Research Center, Boston	NO
SSPPED	http://www.ombl.boidolborg.do/coprod/coprod_info.html	EMBL (Heidelborg)	NO
SOFILED			NO
GOR and			
DSC	http://genome.imb-jena.de/cgi-bin/GDEWWW/menu.cgi	IMB (Jena)	NO
GOR	http://absalpha.dcrt.nih.gov:8008/gor.html	DCRT/NIH (Washington)	NO
GOR	ftp://ftp.virginia.edu/pub/fasta	University of Virginia	YES
Mult-		Ludwig Institute for	
Prodict	http://keatrol.ludwig.uol.ee.uk/zpred.html	Cancer Research	NO
Fredici	jnup.//kesirei.iuuwig.uci.ac.uk/zpreu.nimi		NU UNI

Argos P. (1976) Prediction of the secondary structure of mouse nerve growth factor and its comparison with insulin. *Biochemical and Biophysical Research Communications* 3:805-811.

Bairoch A and Apweiler R. (1996) The SWISS-PROT protein sequence data bank and its new supplement TREMBL. Nucleic Acids Res 24:21-25.

Barton GJ. (1995) Protein secondary structure prediction. Curr Opinion Struct Biol 5:372-376.

Benner SA, Gerloff DL, and Jenny TF. (1994) Predicting protein crystal structures. Science 265:1642-1644.

Benner SA. (1995) Predicting the conformation of proteins from sequences. Progress and future progress. J Mol Recogn 8:9-28.

Boyd, D., Schierle, C. & Beckwith, J. (1998). How many membrane proteins are there? Prot. Sci. 7, 201-205.

Crawford IP, Niermann T, and Kirschner K. (1987) Prediction of secondary structure by evolutionary comparison: application to the alpha subunit of thryptophan synthase. *Proteins: Struct Func Genet* 2:118-129.

Deleage G and Roux B. (1987) An algorithm for protein secondary structure prediction based on class prediction. Protein Engineering 4:289-294.

Eigenbrot C, Randal M, and Kossiakoff AA. (1992) Structural Effects Induced by Mutagenesis Affected by Crystal Packing Factors: the Structure of a 30-51 Disulfide Mutant of Basic Pancreatic Trypsin Inhibitor. *Proteins* 14:75.

Fasman, G. D. & Gilbert, W. A. (1990). The prediction of transmembrane protein sequences and their conformation: an evaluation. Trends Biochem Sci 15, 89-92.

Frishman D, and Argos P. (1995) Knowledge-Base Protein Secondary Structure Assignment. Proteins: Structure, Function, and Genetics 23:566-79.

Frishman D, and Argos P. (1996) Incorporation of Non-Local Interactions in Protein Secondary Structure Prediction From the Amino Acid Sequence. *Protein Engineering* 2:in the press.

Frishman D, and Argos P. (1997) The Future of Protein Secondary Structure Prediction Accuracy. Folding & Design 2:159-62.

Frishman, D, and P Argos. (1996) 75% Accuracy in Protein Secondary Structure Prediction. Proteins 1997 Mar;27(3):329-335.

Garnier J and Levin JM. (1991) The protein structure code: what is its present status. Comput Appl Biosci 7:133-142.

Garnier, J. (1990). Protein structure prediction. Biochimie 72, 513-24.

Garnier, J., Gibrat, J. F. & Robson, B. (1996a). GOR method for predicting protein secondary structure from amino acid sequence. Meth. Enz. 266, 540-553.

Garnier, J., Gibrat, J. F. & Robson, B. (1996b). GOR method for predicting protein secondary structure from amino acid sequence. *Methods Enzymol* 266, 540-53.

Garnier, J., Osguthorpe, D. & Robson, B. (1978). Analysis of the accuracy and implications of simple methods for predicting the secondary structure of globular proteins. *J. Mol. Biol.* **120**, 97-120.

Geourjon C, and Deléage G. (1995) SOPMA: Significant Improvements in Protein Secondary Structure Prediction by Consensus Prediction From Multiple Sequences. *Comput Appl Biosci* 11:681-84.

Gibrat, J., Garnier, J. & Robson, B. (1987). Further developments of protein secondary structure prediction using information theory. J. Mol. Biol. 198, 425-443.

Gilbert RJ. (1992) Protein structure prediction from predicted residue prperties utilizing a digital encoding algorithm. J Mol Graph 10:112-119.

Holley LH, and Karplus M. (1989) Protein Secondary Structure Prediction With a Neural Network. Proc Natl Acad Sci USA 86:152-56.

Hunt NG, Gregoret LM, and Cohen FE. (1994) The origins of protein secondary structure. Effects of packing density and hydrogen bonding studied by a fast conformational search. *J Mol Biol* 241:214-225.

Kabsch W and Sander C. (1984) On the use of sequence homologies to predict protein structure. Identical pentapeptides can have completely different conformation. *Proc Natl Acad Sci USA* 81:1075-1078.

Kabsch W, and Sander C. (1983) Dictionary of Protein Secondary Structure: Pattern Recognition of Hydrogen-Bonded and Geometrical Features. *Biopolymers* 22:2577-637.

Kendrew JC, Klyne W, Lifson S, Miyazawa T, Nemethy G, Phillips DC, Ramachandran GN, and Sheraga HA. (1970). Biochemistry 9:3471-79.

King, R. D. & Sternberg, M. J. E. (1996). Identification and application of the concepts important for accurate and reliable protein secondary structure prediction. *Prot. Sci.* 5, 2298-2310.

King, R. D., Saqi, M., Sayle, R. & Sternberg, M. J. (1997). DSC: public domain protein secondary structure predication. Comput Appl Biosci 13, 473-4.

Levin J, Pascarella S, Argos P, and Garnier J. (1993) Quantification of secondary structure prediction improvement using multiple alignments. *Protein Engineering* 6:849-854.

Levin JM, Robson B, and Garnier J. (1986) An algorithm for secondary structure determination in proteins based on sequence similarity. FEBS Lett 205:303-308.

Levitt M, and Greer J. (1977) Automatic Identification of Secondary Structure in Globular Proteins. J Mol Biol 114:181-239.

Lim VI. (1974) Algorithms for prediction of alpha-helical and beta-structural regions in globular proteins. J Mol Biol 88:873-894.

Livingstone CD, Barton GJ (1996). Identification of functional residues and secondary structure from protein multiple sequence alignment. Methods Enzymol 266:497-512

Lupas A, Koster AJ, Walz J, and Baumeister W. (1994) Predicted secondary structure of the 20 S proteasome and model structure of the putative peptide channel. *FEBS Lett* 354:45-49.

Matthews B. (1975) Comparison of the predicted and observed secondary structure of T4 phage lysozyme. Biochem Biophys Acta 405:442-451.

Mehta PK, Heringa J, and Argos P. (1995) A simple and fast approach to prediction of protein secondary structure from multiply aligned sequences with accuracy above 70%. *Protein Science* 4:2517-2525.

Muggleton S, King RD, and Sternberg MJE. (1992) Protein Secondary Structure Prediction Using Logic-Based Machine Learning. *Protein Engineering* 5:647-57. Nishikawa K, and Ooi T. (1986) Amino Acid Sequence Homology Applied to the Prediction of Protein Secondary Structures, and Joint Prediction With Existing Methods. *Biochimica Et Biophysica Acta* 871:45-54.

Pauling L, Corey RB, and Branson HR. (1951) The structure of proteins: two hydrogen-bonded helical configurations of the polypeptide chain. *Proc Natl Acad Sci USA* 37:205-211.

Persson, B. & Argos, P. (1997). Prediction of membrane protein topology utilizing multiple sequence alignments. J Protein Chem 16, 453-7.

Presnell SR, Cohen BI, and Cohen FE. (1992) A segment-based approach to protein secondary structure prediction. Biochemistry 31:983-993.

Ptitsyn OB and Finkelstein AV. (1983) Theory of protein secondary structure and algorithm of its prediction. Biopolymers 22:15-25.

Qian N, and Sejnowski TJ. (1988) Predicting the Secondary Structure of Globular Proteins Using Neural Network Models. J Mol Biol 202:865-84.

Rackovsky S. (1993) On the nature of the protein folding code. Proc Natl Acad Sci U S A 90:644-648.

Ramakrishnan C, and Soman KV. (1982) Identification of Secondary Structures in Globular Proteins - a New Algorithm. Int J Pept Protein Res 20:218-37.

Rao S, Zhu Q-L, Vaida S, and Smith T. (1993) The local information content of the protein structural database. FEBS Lett 2:143-146.

Rice CM, Fuchs R, Higgins DG, Stoehr PJ, and Cameron G N. (1993) The EMBL data library. Nucleic Acids Res 21:2967-2971.

Richards FM, and Kundrot CE. (1988) Identification of Structural Motifs From Protein Coordinate Data: Secondary Structure and First-Level Supersecondary Structure. *Proteins: Struct Func Genet* 3:71-84.

Robson B, and Garnier J. (1993) Protein Structure Prediction. Nature 361:506.

Rost B, and Sander C. (1993) Prediction of Protein Secondary Structure at Better Than 70% Accuracy. J Mol Biol 232:584-99.

Rost B, Sander C, and Schneider R. (1994) Redefining the goals of protein secondary structure prediction. J Mol Biol 235:13-26.

Rost B, Sander C. (1993) Improved prediction of protein secondary structure by use of sequence profiles and neural networks. *Proc Natl Acad Sci USA* 90:7558-7562

Rost, B., Schneider, R. & Sander, C. (1993). Progress in protein structure prediction? Trends Biochem Sci 18, 120-3.

Rumelhart DE, Hinton GE, and Williams R. (1986) Learning representations by back-propagating errors. Nature 323:533-536.

Salamov AA, and Solovyev VV. (1995) Prediction of Protein Secondary Structure by Combining Nearest-Meighbour Akgorithms Amd Multiple Sequence Alignments. *J Mol Biol* 247:11-15.

Salamov AA, and Solovyev VV. (1997) Protein Secondary Structure Prediction Using Local Alignments. Journal of Molecular Biology 268:31-36.

Sayle RA, and Milner-White EJ. (1995) RASMOL: Biomolecular Graphics for All. Trends in Biochemical Sciences 20:374-76.

Sayle RA, and Milner-White EJ. (1995) RASMOL: Biomolecular Graphics for All. Trends in Biochemical Sciences 20:374-76.

Sklenar H, Etchebest C, and Lavery R. (1989) Describing Protein Structure: a General Algorithm Yielding Complete Helicoidal Parameters and a Unique Overall Axis. *Proteins: Struct Func Genet* 6:46-60.

Solovyev VV and Salamov AA. (1994) Predicting alpha-helix and beta-strand segments of globular proteins. Comput Appl Biosci 10:661-669.

Stolorz P, Lapedes A, and Xia Y. (1992) Predicting Protein Secondary Structure Using Neural Net and Statistical Methods. J Mol Biol 225:363-77.

Sumpter BG, Getino C, and Noid DW. (1994) Theory and applications of neural computing in chemical science. Ann Rev phys Chem 45:439-481.

Sumpter BG, Getino C, and Noid DW. (19949 Theory and applications of neural computing in chemical science. Ann Rev phys Chem 45:439-481.

Taylor WR and Thornton JM. (1984) Recognition of super-secondary structure in proteins. J Mol Biol 173:487-5141984.

Thompson JD, Higgins DG, and Gibson TJ. (1994) CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* 22:4673-4680.

Thornton JM, Flores TP, Jones DT, and Swindells MB. (1991) Prediction of progress at last. Nature 354:105-106.

von Heijne, G. (1992). Membrane protein structure prediction. Hydrophobicity analysis and the positive-inside rule. J Mol Biol 225, 487-94.

Wasserman PD. (1989) Neural Computing. Theory and Practice. New York.

Zhang X, Mesirov JP, and Waltz DL. (1992) Hybrid System for Protein Secondary Structure Prediction. J Mol Biol 225:1049-63.

Zvelebil MJ, Barton GJ, Taylor WR, and Sternberg MJ. (1987) Prediction of Protein Secondary Structure and Active Sites Using the Alignment of Homologous Sequences. *J Mol Biol* 195:957-61.

## Subject II:

#### Packing

## Other Aspects of Structure, Besides just Comparing Atom Positions







Atom Position, XYZ triplets

Lines, Axes, Angles Surfaces, Volumes

#### What is Protein Geometry?

- Coordinates (X, Y, Z's)
- Derivative Concepts
  - Distance, Surface Area,
    Volume, Cavity, Groove,
    Axes, Angle, &c
- Relation to
  - Function,
    Energies (E(x)),
    Dynamics (dx/dt)



# <u>Close-Packing</u> of Spheres

- Efficiency
  - Volume Spheres / Volume of space
- Close packed spheres
  - ◊ 74% volume filled
  - Occordination of 12
  - Two Ways of laying out
- Fcc
  - ◊ cubic close packing
  - ♦ ABC layers
- hcp
  - Hexagonally close packed
  - ♦ ABABAB





gives rise to face-centred unit cells, and so may also be denoted cubic F (or



Fig. 21.20 The close-packing of identical spheres. (a) The first layer of close-packed spheres. (b) The second layer of close-packed apheres occupies the dips of the first layer. The two layers are the AB component of the structure.

Fig. 21.21 The third layer of close-packed spheres might occupy the dips lying directly above the spheres in the first layer, resulting in an ABA structure (a) which corresponds to hexagonal close-packing (b). This hop structure is possessed by the elements Be, Cd, Co, He, Mg, Ti, and Zh.

Illustration Credits: Atkins, Pchem, 634





Fig. 21.22 Alternatively, the third layer might lie in the dips that are not above the apheres in the first layer, resulting in an ABC structure (a) which correspond to cubic close-packing (b). This cop (or foc) structure is possessed by the elements Ag, Al, Ar, Au, Ca, Cu, Ne, Ni, Pb, Pt, and Xe.

47

## <u>Other Well Known</u> <u>Sphere</u> <u>Arrangements</u>

- Simple cubic packing
  - ◊ 8 nbrs
  - ♦ 52% efficiency
- bcc cubic packing
  - one sphere sits in middle of 8 others (body-centered)
  - ◊ 8 nbrs
  - ◊ 68% efficiency
- fcc -> bcc -> simple
  - ◊ apx 3/4, 2/3, 1/2





Space between spheres



## **Optimal Packing Finally Proved**

#### After Four Centuries, an Answer

What's the best way to stack a bunch of round objects? The answer, whether they are cannonballs or oranges, seems to be an extension of the familiar pyramid-shaped stack seen in grocery stores everywhere.



Stacking efficiency = volume of the spheres / (volume of the spheres + the space between the spheres)

Illustration Credits: Singh, New York Times

#### Water v. Argon



More Complex Systems -- what to do?

#### Voronoi Volumes

- Each atom surrounded by a single convex polyhedron and allocated space within it
  - Allocation of all space (large V implies cavities)
- 2 methods of determination
  - Find planes separating atoms, intersection of these is polyhedron
  - Locate vertices, which are equidistant from 4 atoms



#### **Classic Papers**

- Lee, B. & Richards, F. M. (1971). "The Interpretation of Protein Structures: Estimation of Static Accessibility," *J. Mol. Biol.* 55, 379-400.
- Richards, F. M. (1974). "The Interpretation of Protein Structures: Total Volume, Group Volume Distributions and Packing Density," *J. Mol. Biol.* 82, 1-14.
- Richards, F. M. (1977). "Areas, Volumes, Packing, and Protein Structure," Ann. Rev. Biophys. Bioeng. 6, 151-76.

## Voronoi Volumes, the Natural Way to Measure Packing

**Packing Efficiency** 

= Volume-of-Object

Space-it-occupies

= V(VDW) / V(Voronoi)

- Absolute v relative eff. V1 / V2
- Other methods
  - Measure Cavity Volume (grids, constructions, &c)



## Calclating Volumes with Voronoi polyhedra

- In 1908 Voronoi found a way of partitioning all space amongst a collection of points using specially constructed polyhedra. Here we refer to a collection of "atom centers" rather than "points."
- In 3D, each atom is surrounded by a unique limiting polyhedron such that all points within an atom's polyhedron are closer to this atom than all other atoms.
- Likewise, points equidistant from 2 atoms form planes (lines in 2D). Those equidistant from 3 atoms form lines, and those equidistant form 4 centers form vertices.

## Delauney Triangulation, the Natural Way to Define Packing Neighbors

- Related to Voronoi polyhedra (dual)
- What "coordination number" does an atom have? Doesn't depend on distance
- alpha shape
- threading





## Properties of Voronoi Polyhedra

- If Voronoi polyhedra are constructed around atoms in a periodic system, such as in a crystal, all the volume in the unit cell will be apportioned to the atoms. There will be no gaps or cavities as there would be if one, for instance, simply drew spheres around the atoms.
- Voronoi volume of an atom is a weighted average of distances to all its neighbors, where the weighting factor is the contact area with the neighbor.

## Voronoi diagrams are generally useful, beyond proteins

- Border of D.T. is Convex Hull
- D.T. produces "fatest" possible triangles which makes it convenient for things such as finite element analysis.
- Nearest neighbor problems. The nearest neighbor of a query point in center of the Voronoi diagram in which it resides
- Largest empty circle in a collection of points has center at a Voronoi vertex
- Voronoi volume of "something" often is a useful weighting factor. This fact can be used, for instance, to weight sequences in alignment to correct for over or under-representation

## Surfaces (slides 1-10, 20-40 from website)

These are detailed slides on how to do Voronoi construction. Go to http://bioinfo.mbb.yale.edu/geometry and follow links to "HyperTalk" tutorial on surfaces and volumes

#### Atoms have different sizes

- Difficulty with Voronoi Meth. Not all atoms created equal
- Solutions
  - Bisection -- plane midway between atoms
  - Method B (Richards)
    Positions the dividing plane according to ratio
  - ◊ Radical Plane
- VDW Radii Set



## Set of VDW Radii

- Great differences in a sensitive parameter (Radii for carbon 1.87 vs 2.00)
- Complex calculation: minimizing SD, iterative procedure, from protein structures

Atom	Bondi	New
C4	1.87	1.88
C3H1	1.76	1.76
СЗНО	1.76	1.61
Olho	1.40	1.42
O2H1	1.40	1.46
N	1.65	1.64
S	1.85	1.77

- Look for common distances in CCD
- Preliminary Solution

### Different Sets of Radii

Atom Type & Symbol		Bondi 1968	Lee & Richards 1971	Shrake & Rupley 1973	Richards	Chothia 1975	Rich- mond & Richards 1978	Gelin & Karplus 1979	Dunfield et al. 1979	ENCAD derived 1995	CHARMM derived 1995	Tsai et al. 1998
$-CH_3$	Aliphatic, methyl	2.00	1.80	2.00	2.00	1.87	1.90	1.95	2.13	1.82	1.88	1.88
$-CH_2-$	Aliphatic, methyl	2.00	1.80	2.00	2.00	1.87	1.90	1.90	2.23	1.82	1.88	1.88
>CH-	Aliphatic, CH	-	1.70	2.00	2.00	1.87	1.90	1.85	2.38	1.82	1.88	1.88
=CH	Aromatic, CH	-	1.80	1.85	*	1.76	1.70	1.90	2.10	1.74	1.80	1.76
>C=	Trigonal, aromatic	1.74	1.80	*	1.70	1.76	1.70	1.80	1.85	1.74	1.80	1.61
$-NH_3+$	Amino, protonated	-	1.80	1.50	2.00	1.50	0.70	1.75		1.68	1.40	1.64
$-NH_2$	Amino or amide	1.75	1.80	1.50	-	1.65	1.70	1.70		1.68	1.40	1.64
>NH	Peptide, NH or N	1.65	1.52	1.40	1.70	1.65	1.70	1.65	1.75	1.68	1.40	1.64
=0	Carbonyl Oxygen	1.50	1.80	1.40	1.40	1.40	1.40	1.60	1.56	1.34	1.38	1.42
-OH	Alcoholic hydroxyl	-	1.80	1.40	1.60	1.40	1.40	1.70	ĺ	1.54	1.53	1.46
-OM	Carboxyl Oxygen	-	1.80	1.89	1.50	1.40	1.40	1.60	1.62	1.34	1.41	1.42
-SH	Sulfhydryl	-	1.80	1.85	-	1.85	1.80	1.90		1.82	1.56	1.77
-S-	Thioether or -S-S-	1.80	-	-	1.80	1.85	1.80	1.90	2.08	1.82	1.56	1.77

#### **Standard Residue Volumes**

- Database of many hi-res structures (~100, 2 Å)
- Volumes statistics for buried residues (various selections, resample, &c)
- Standard atomic volumes harder... parameter set development...

G 64 c 105 T 120 V 139 H 159 M 168 R 194 A 90 C 113 P 124 E 140 L 165 K 170 Y 198 S 94 D 117 N 128 N 150 I 165 F 193 W 233

## Standard Core Volumes (Prelim.)

Atom Types		Num.	Volume (Å <sup>3</sup> )	Error (%)
Mainchain Atoms				
carbonyl carbon (except G)	С	8361	9.2	.08
alpha carbon (except G)	CA	7686	13.4	.09
nitrogen (except P)	Ν	9042	13.9	.09
carbonyl oxygen	0	7831	15.8	.10
Gly C		811	10.2	.27
Gly CA		522	23.5	.39
Pro N		334	8.6	.39
Sidechain atoms				
trigonal or aromatic carbon	>C=	3026	10.3	.13
aromatic CH (H,F,W,Y)	-CH=	4333	21.1	.14
aliphatic CH	>CH-	3411	14.6	.14
methylene group	-CH2-	5427	23.7	.12
methyl group (A,V,L,I)	-CH3	5273	36.7	.11
hydroxyl oxygen (S,T)	-OH	851	17.2	.36
carbonyl oxygen (N,Q)	=0	272	16.8	.76
carboxyl oxygen (D,E)	-0	517	16.0	.53
2° amine (R,H,W)	-NH-	530	15.6	.53
<pre>1° amine or amide (R,N,Q)</pre>	-NH2	355	23.4	.52
tetrahedral nitrogen (K)	-NH3	31	20.0	1.40
thioether or disulfide (C,M)	-S-	1242	19.3	1.22
sulfhydryl (C)	-SH	67	37.8	1.33

#### Packing at Interfaces

- Voronoi volumes (and D. triangulation) to measure packing
- Tight core packing v. Loose surface packing
- Grooves & ridges: closepacking v. H-bonding
- How packing defines a surface (hydration surface)
- Implications for Motions



<u>Richards'</u> <u>Molecular</u> <u>and</u> <u>Accessible</u> <u>Surfaces</u>



	Probe Radius	Part of Probe Sphere	Type of Surface
-	Radius		
	0	Center (or Tangent)	Van der Waals Surface (vdWS)
-	1.4 Å	Center	Solvent Accessible Surface (SAS)
		Tangent (1 atom)	Contact Surface (CS, from parts of
			atoms)
		Tangent (2 or 3 atoms)	Reentrant Surface (RS, from parts of
			Probe)
	""	Tangent (1,2, or 3 atoms)	Molecular Surface ( $MS = CS + RS$ )
	10 Å	Center	A Ligand or Reagent Accessible Surface
	~	Tangent	Minimum limit of MS (related to convex
			hull )
		Center	Undefined

## Packing defines the "Correct Definition" of the Protein Surface

- Voronoi polyhedra are the Natural way to study packing!
- How reasonable is a geometric definition of the surface in light of what we know about packing
- The relationship between
  - ◊ accessible surface
  - ◊ molecular surface
  - Oblauney Triangulation (Convex Hull)
  - ◊ polyhedra faces
  - ◊ hydration surface

## <u>Surface and Volume</u> <u>Definitions Linked</u>



## Problem of Protein Surface for Voronoi Construction



## <u>Defining Surfaces from Packing:</u> <u>Convex Hull and Layers of Waters</u>



## Defining a Surface from the Faces of Voronoi Polyhedra





Protein

## <u>Accessible Surface</u> as a Time-averaged Water Layer



## <u>The Hydration Surface:</u> <u>Trying to Model Real Water</u>





Protein
## **Small Packing Changes Significant**

- Exponential dependence
- Bounded within a range of 0.5 (.8 and .3)
- Many observations in standard volumes gives small error about the mean (SD/sqrt(N))



## Packing ~ VDW force

- Longer-range isotropic attractive tail provides general cohesion
- Shorter-ranged repulsion determines detailed geometry of interaction
- Billiard Ball model, WCA Theory



## **Close-packing is Default**

- No tight packing when highly directional interactions (such as H-bonds) need to be satisfied
- Packing spheres (.74), hexagonal
- Water (~.35), "Open" tetrahedral, H-bonds

