

**Selection of slides from:
Multiple Sequence Alignment**

BMI/CS 576

www.biostat.wisc.edu/bmi576.html

Colin Dewey

cdewey@biostat.wisc.edu

Fall 2010

Multiple Sequence Alignment: Task Definition

- Given
 - a set of more than 2 sequences
 - a method for scoring an alignment
- Do:
 - determine the correspondences between the sequences such that the alignment score is maximized

Motivation for MSA

- establish input data for phylogenetic analyses
- determine evolutionary history of a set of sequences
 - At what point in history did certain mutations occur?
- discovering a common motif in a set of sequences (e.g. DNA sequences that bind the same protein)
- characterizing a set of sequences (e.g. a protein family)

Multiple Alignment of SH3 Domain

```
GGWWRG d y . g g k k q L W F P S N Y V
IGWLN G y n e t t g e r r G D F P G T Y V
PNWWE G q l . . n n r r r G I F P S N Y V
DEW W Q A r r . . d e q i G I V P S K - -
GEW W K A q s . . t g q e G F I P F N F V
GDW W L A r s . . s g q t G Y I P S N Y V
GDW W D A e l . . k g r r r G K V P S N Y L
- D W W E A r s l s s g h r r G Y V P S N Y V
GDW W Y A r s l i t n s e G Y I P S T Y V
GEW W K A r s l a t r k e G Y I P S N Y V
GDW W L A r s l v t g r e G Y V P S N F V
GEW W K A k s l s s k r e G F I P S N Y V
GEW C E A q t . k n g q . G W V P S N Y I
SDW W R V v n l t t r q e G L I P L N F V
LPW W R A r d . k n g q e G Y I P S N Y I
RDW W E F r s k t v y t p G Y Y E S G Y V
EHW W K V k d . a l g n v G Y I P S N Y V
IHW W R V q d . r n g h e G Y V P S S Y L
KDW W K V e v . . n d r q G F V P A A Y V
VGW M P G l n e r t r q r G D F P G T Y V
PDW W E G e l . . n g q r G V F P A S Y V
ENW W N G e i . . g n r k G I F P A T Y V
EEW L E G e c . . k g k v G I F P K V F V
GGW W K G d y . g t r i q Q Y F P S N Y V
DGW W R G s y . . n g q v G W F P S N Y V
QGW W R G e l . . y g r v G W F P A N Y V
GRW W K A r r . a n g e t G I I P S N Y V
GGW T Q G e l . k s g q k G W A P T N Y L
GDW W E A r s n . t g e n G Y I P S N Y V
NDW W T G r t . . n g k e G I F P A N Y V
```

Figure from A. Krogh, An Introduction to Hidden Markov Models for Biological Sequences

Scoring a Multiple Alignment

- key issue: how do we assess the quality of a multiple sequence alignment?
- usually, the assumption is made that the individual *columns* of an alignment are independent

$$\text{Score}(m) = G + \sum_i S(m_i)$$

gap function score of i^{th} column

- we'll discuss two methods
 - sum of pairs (SP)
 - minimum entropy

Scoring an Alignment: Sum of Pairs

- compute the sum of the pairwise scores

$$S(m_i) = \sum_{k < l} s(m_i^k, m_i^l)$$

m_i^k = character of the k th sequence in the i th column

S = substitution matrix

Scoring an Alignment: Minimum Entropy

- basic idea: try to minimize the *entropy* of each column
- another way of thinking about it: columns that can be communicated using few bits are good
- information theory tells us that an optimal code uses $-\log_2 p$ bits to encode a message of probability p

Scoring an Alignment: Minimum Entropy

- the messages in this case are the characters in a given column
- the entropy of a column is given by:

$$S(m_i) = -\sum_a c_{ia} \log_2 p_{ia}$$

m_i = the i th column of an alignment m

c_{ia} = count of character a in column i

p_{ia} = probability of character a in column i

Dynamic Programming Approach

- can find optimal alignments using dynamic programming
- generalization of methods for pairwise alignment
 - consider k -dimension matrix for k sequences (instead of 2-dimensional matrix)
 - each matrix element represents alignment score for k subsequences (instead of 2 subsequences)
- given k sequences of length n
 - space complexity is

$$O(n^k)$$

Heuristic Alignment Methods

- since time complexity of DP approach is exponential in the number of sequences, heuristic methods are usually used
- *progressive alignment*: construct a succession of pairwise alignments
 - star approach
 - tree approaches, like CLUSTALW
 - etc.
- iterative refinement
 - given a multiple alignment (say from a progressive method)
 - remove a sequence, realign it to profile of other sequences
 - repeat until convergence

Multiple Alignment Case Study: The Cystic Fibrosis Gene

- cystic fibrosis (CF)
 - recessive genetic disease caused by a defect in a single-gene
 - causes the body to produce abnormally thick mucus that clogs the lungs and the pancreas
- the cystic fibrosis conductance regulator (CFTR) gene
 - gene and its role in CF identified in 1989 [Riordan et al., *Science*]
 - most common mutation is called $\Delta F508$; a deletion of a phenylalanine (F) at position 508 in the CFTR protein
 - the CFTR protein controls the movement of salt and water into and out of cells; mutations in CFTR block this movement, causing mucus problem

So What Does CFTR Do? A CFTR Multiple Alignment

			*	
CFTR (N)	FSLGLGTPVLKDNFVKIERGQLLAVAGSTGAGKTSLLMMIMG	ISFCSQFSWIMPGTIK-ENIIFGVSYD		GEGGITLSGGQRARISLARAVYKDADLYLLDSPFGYLDVLTEK
CFTR (C)	YTEGGNAILENISFSISPGQRVGLLGRGTSGSGKSTLLSAFLR	DSITLQQRKAFGVIPOKVFIFSGTFR		VDGGCVLSHGKQLMCLARSVLSKAKILLLDEPSAHLDPVTYQ
hmdr1 (N)	PSRKEVKILKGLNLKLVQSGQTVALVGNSSGCGKSTTVQLMQR	IGVVSQEPVLFATTI-AENIRYGRENV		GERGAQLSGGQKQRIAIARALVRNPKILLLDEATSALDTESEA
hmdr1 (C)	PTRPDIPVLQGLSLEVKKGQTLALVGNSSGCGKSTVVQLLER	LGIVSQEPILFDCSI-AENIAYGDNSR		GDKGTLLSGGQKQRIAIARALVRPHILLLDEATSALDTESEK
mmdr1 (N)	PSRSEVQILKGLNLKVKSGQTVALVGNSSGCGKSTTVQLMQR	IGVVSQEPVLFATTI-AENIRYGRENV		GERGAQLSGGQKQRIAIARALVRNPKILLLDEATSALDTESEA
mmdr1 (C)	PTRPNIPVLQGLSLEVKKGQTLALVGNSSGCGKSTVVQLLER	LGEVVSQEPILFDCSI-AENIAYGDNSR		GDKGTQLSGGQKQRIAIARALVRPHILLLDEATSALDTESEK
mmdr2 (N)	PSRANIKILKGLNLKVKSGQTVALVGNSSGCGKSTTVQLLQR	IGVVSQEPVLSFTTI-AENIRYGRGNV		GDRGAQLSGGQKQRIAIARALVRNPKILLLDEATSALDTESEA
mmdr2 (C)	PTRANVPVLQGLSLEVKKGQTLALVGNSSGCGKSTVVQLLER	LGIVSQEPILFDCSI-AENIAYGDNSR		GDKGTQLSGGQKQRIAIARALVRPHILLLDEATSALDTESEK
pfmdr (N)	DTRKDVEIYKDLSEFTLLKEGKTYAFVGSSECGKSTILKLE	IGVVSQEPVLSFTTI-AENIRYGRGNV		GSNASKLSGGQKQRIAIARALVRNPKILLLDEATSALDTESEK
pfmdr (C)	ISRPNVPIYKNLSFTCDSSKTTAIVGETSGSGKSTFMNLLLR	FSIVSQEPMLFNMSI-YENIKFGREDA		PYGKS-LSGGQKQRIAIARALLREPILLLDEATSALDTESEK
STE6 (N)	PSRPSEAVLKNVSLNFSAGQFTFIVGKSGSGKSTLSNLLLR	ITVVEQRCTLFNDTL-RKNILLGSTDS		GTGGVTLSSGGQQRVAIARAFIRDTPILFLDEAVSALDIVHRN
STE6 (C)	PSAPTAFVYKNMNFDMFCGQTLGIIGESGTGKSTLVLLTK	ISVVEQKPLLFNGTI-RDNLYTGLQDE		RIDTTLSSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
hlyB	YKPDSPVILDNINISIKQGEVIGVGRSGSGKSTLIKLIQR	VGVVLQDNVLLNRSI-IDNISLAPGMS		GEQGAGLSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
White	IPAPRKHLKKNVCGVAYPGELLAVMGSSGAGKTTLLNALAF	RCAYVQDDDLFIGLIAREHLIFQAMVR		PGRVKGLSGGERKRLAFASEALTDPPLLICDEPTSLDVSFTAH
MbpX	KSLGNLKILDRVSLYVPKFSLIALLGPGSGSGKSSLLRILAG	MSFVVFQHYALFKHMTVYENISFGLRLR		FEYPAQLSGGQKQRVALARSLAIQPDLL-DEPFGALDDELRR
BtuD	QDVAESTRLGPLSGEVRAGRILHLVGPNGAGKSTLLARIAG	YLSQQQTPPFATPVVWHYLTLDHQDKTR		GRSTNQLSSGGEQQRVRLAAVVLTLLLLDEPMNSLDVAQOSA
PstB	FYYGKFHALKNINLDTAKNQVTAFIGPSCCGKSTLRFTFNK	VGMVFQKPTFPFMSI-YDNIAFGVRLF		HQSGYLSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
hisP	RRYGGHEVLKGVSLQARAGVISIIGSSSGKSTFLRCINF	GIMVFQHPNLWSHMTVLENVMEAPIQV		GKYPVHLSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
malk	KAWGEVVVSKDINIDIEHEGEFVVVFGPSCGCGKSTLLRMIAG	VGMVFQSYALYPHLSVAENMSFGLKPA		DRKPKALSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
oppD	TPDGDVTAVNDLNFTRAGETLGIVGESGSGKSTAFALMG	ISMIFQDPMTSLNPNYMRVGEQLMEVLM		KMPYHEFSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
oppF	QPPKTLKAVDGVTLRLYEGETLGVVGESGCGKSTFARAIIG	IQMIFQDPLASLNPRMTIGEIIAEPLR		NRYPHFSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
RbsA (N)	KAVPGVKALSGAALNVYPGRVMALVGENGAGKSTMMKVLTG	AGIIHQELNLIPLQTLIAENIFLGRFV		DKLVGDLSIGDQQMVEIAKVLVSFESKVIIMDEPTCALIDTETE
RbsA (C)	VDNLCPGVNDVVSFTLRKGEILGVSLMGAGRTELMKVLVYG	ISEDKRKDLVLVGMVKENMSLTALRY		EQAIGLLSGGQQRVAIARALVRNPKILLLDEATSALDTESEK
UvrA	LTGARGNNLKDVTLLPVLGFTCITGVSSGKSTLINDTLF	TYTGVTVPVRELFAVPESTRARGYTPG		GQSATTLSSGGEAQRVKLARELSKRGVILDEPTTGLHFDIQQ
NodI	KSYGKIVVNDLSFTTAAGECFLLGPNAGKSTIIRMILG	IGIVSQEPNLDLEFTVRENILVYGRYF		NTRVADLSGGMQRRLTLAGALINDPQLLILDEPTTGLDPHARH
FtsE	AYLGGRQALQGVTFHMQPGEMAFLLTGHSAGKSTLLKLCIG	IGMIFQDHHLLMDRTVYDNVAIPLIIA		KNFPIQLSSGGEQQRVAIARALVRNPKILLLDEATSALDTESEK

Figure from Riordan et al, *Science* 245:1066-1073, 1989.

Multiple Alignment Case Study: the Cystic Fibrosis Gene

- two key features of the protein made apparent in multiple sequence alignment (and other analyses)
 - membrane-spanning domains
 - ATP-binding motifs
- these features indicated that CFTR is likely to be involved in transporting ions across the cell membrane

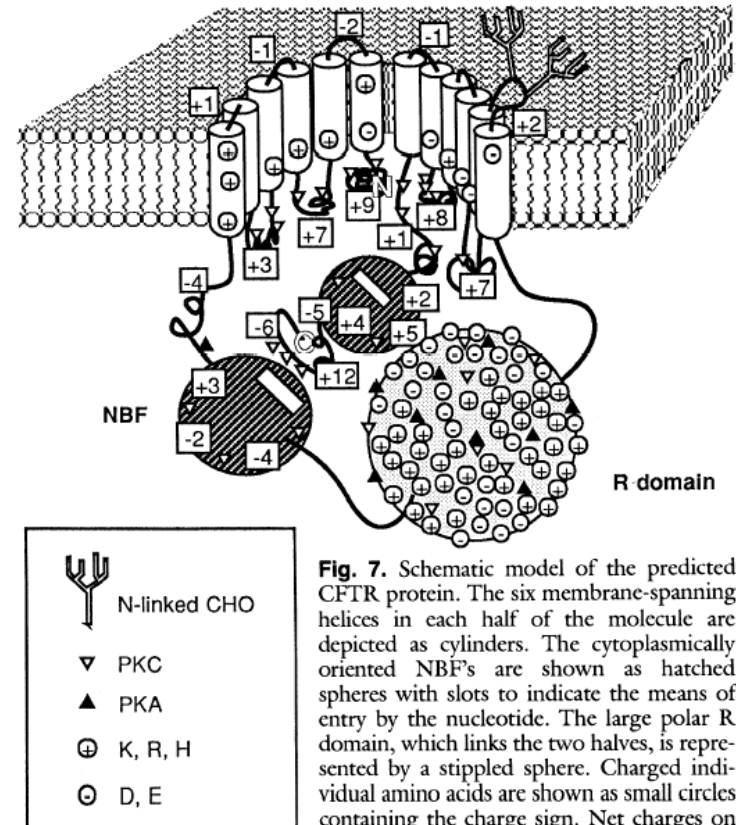


Fig. 7. Schematic model of the predicted CFTR protein. The six membrane-spanning helices in each half of the molecule are depicted as cylinders. The cytoplasmically oriented NBF's are shown as hatched spheres with slots to indicate the means of entry by the nucleotide. The large polar R domain, which links the two halves, is represented by a stippled sphere. Charged individual amino acids are shown as small circles containing the charge sign. Net charges on the internal and external loops joining the membrane cylinders and on regions of the NBF's are contained in open squares. Potential sites for phosphorylation by protein kinases A or C (PKA or PKC) and N-glycosylation (N-linked CHO) are as indicated. K, Lys; R, Arg; H, His; D, Asp; and E, Glu.

Notes on Multiple Alignment

- as with pairwise alignment, can compute *local* and *global* multiple alignments
- dynamic programming is not feasible for most cases -- heuristic methods usually used instead