# Note on sequences and alignment matrices in exercises

- Example solutions to alignment problems will have sequences arranged like this:
  - Perform global alignment of the sequences
    - s = AGCTGCGTACT
    - t = ATGAGCGTTA

So if you want to be able to compare your solution easily against the example, use this convention.





### Tracing back optimal path

- When filling the DP matrix, one can store to each cell the list of pointers to the cells where the optimal value was computed; then tracing back all the optimal paths is easy.
- Explicit pointers are actually not needed, since C one can just reverse the computation; start from G cell (m,n), check which are the neighboring cells T where the optimum value must have been computed, and continue like this to follow all optimal paths.

-	1	G	G	1	G
0	-2	-4	-6	-8	-10
-2	-1	-3	-5	-7	-9
-4	-1	-2	-4	-4	-6
-6	-3	-2	<sup>ې</sup>	-5	-5
-8	-5	-2	-1 🗖	-3	-4
-10	-7	-4	-3	0	-2



# Overlap alignment

- Overlap matrix used by Overlap-Layout-Consensus algorithm can be computed with dynamic programming
- Initialization:  $O_{i,0} = O_{0,i} = 0$  for all i, j
- Recursion:

$$O_{i,j} = \max \{ \\ O_{i-1,j-1} + s(a_i, b_i), \\ O_{i-1,j} - \delta, \\ O_{i,j-1} - \delta, \\ \}$$



Best overlap: maximum value from rightmost column and bottom row

→ ...more on sequence assembly at the Spring 2010 seminar!



## Non-uniform mismatch penalties

- We used uniform penalty for mismatches: s('A', 'C') = s('A', 'G') = ... = s('G', 'T') = µ
- Transition mutations (A<->G, C<->T) are approximately twice as frequent than transversions (A<->C, A<->T, C<->G, G<->T)
  - use non-uniform mismatch penalties collected into a substitution matrix





# Gaps in alignment

Gap is a succession of indels in alignment

C T <u>- -</u> A A C T C G C A A

- Previous model scored a length k gap as w(k) = -kδ
- Replication processes may produce longer stretches of insertions or deletions

5

 In coding regions, insertions or deletions of codons may preserve functionality

#### Gap open and extension penalties (2)

We can design a score that allows the penalty opening gap to be larger than extending the gap:

 $w(k) = -\alpha - \beta(k-1)$ 

- Gap open cost α, Gap extension cost β
- Alignment algorithms can be extended to use w(k) as follows:

$$S_{i,j} = \max_{i'+j' < i+j} \{ S_{i-1,j-1} + S(a_i, b_j), \\ S_{i',j'} + W(j - j' + i - i') \}$$



However, this is too inefficient: O(m<sup>2</sup>n<sup>2</sup>).

# Speeding up gap open and extension alignment

$$\begin{array}{rl} S_{i,j} &=& max \ \{ & & \\ & S_{i-1,j-1} \ + \ s(a_i, \ b_i), \\ & & G_{i-1,j-1} \ + \ s(a_i, \ b_i) \\ & & \\ \end{array} \right.$$

$$\begin{array}{l} G_{i,j} = \max \ \{ \\ & S_{i-1,j} - \alpha, \\ & G_{i-1,j} - \beta, \\ & S_{i,j-1} - \alpha, \\ & G_{i,j-1} - \beta, \\ \end{array} \end{array}$$



 Equivalent result in O(mn) time.



# Demonstration of the EBI web site

- European Bioinformatics Institute (EBI) offers many biological databases and bioinformatics tools at http://www.ebi.ac.uk/
  - Sequence alignment: Tools -> Sequence Analysis -> Align



#### Sequence Alignment (chapter 6)

- The biological problem
- Global alignment
- Local alignment
- Multiple alignment



# Multiple alignment

- Consider a set of d sequences on the right
  - Orthologous sequences from different organisms
  - Paralogs from multiple duplications
- How can we study relationships between these sequences?

aggcgagctgcgagtgcta cgttagattgacgctgac ttccggctgcgac gacacggcgaacgga agtgtgcccgacgagcgaggac gcgggctgtgagcgcta aagcggcctgtgtgcccta atgctgctgccagtgta agtcgagcccgagtgc agtccgagtcc actcggtgc

# Optimal alignment of three sequences

- Alignment of A = a<sub>1</sub>a<sub>2</sub>...a<sub>i</sub> and B = b<sub>1</sub>b<sub>2</sub>...b<sub>j</sub> can end either in (-, b<sub>j</sub>), (a<sub>i</sub>, b<sub>j</sub>) or (a<sub>i</sub>, -)
- $2^2 1 = 3$  alternatives
- Alignment of A, B and C = c<sub>1</sub>c<sub>2</sub>...c<sub>k</sub> can end in 2<sup>3</sup> 1 ways: (a<sub>i</sub>, -, -), (-, b<sub>j</sub>, -), (-, -, c<sub>k</sub>), (-, b<sub>j</sub>, c<sub>k</sub>), (a<sub>i</sub>, -, c<sub>k</sub>), (a<sub>i</sub>, b<sub>j</sub>, -) or (a<sub>i</sub>, b<sub>j</sub>, c<sub>k</sub>)
- Solve the recursion using three-dimensional dynamic programming matrix: O(n<sup>3</sup>) time and space
- Generalizes to d sequences but impractical with even a moderate number of sequences

# Multiple alignment in practice

- In practice, real-world multiple alignment problems are usually solved with heuristics
- Progressive multiple alignment
  - Choose two sequences and align them
  - Choose third sequence w.r.t. two previous sequences and align the third against them

29.9-1.10/

- Repeat until all sequences have been aligned
- Different options how to choose sequences and score alignments
- Note the similarity to Overlap-Layout-Consensus

# Multiple alignment in practice

- Profile-based progressive multiple alignment: CLUSTALW
  - Construct a distance matrix of all pairs of sequences using dynamic programming
  - Progressively align pairs in order of decreasing similarity
  - CLUSTALW uses various heuristics to contribute to accuracy

# Additional material

- R. Durbin, S. Eddy, A. Krogh, G. Mitchison: Biological sequence analysis
- N. C. Jones, P. A. Pevzner: An introduction to bioinformatics algorithms
- Biological sequence analysis course, next time in 2010-2011?

# Rapid alignment methods: FASTA and BLAST (Section 7)

- The biological problem
- Search strategies
- FASTA
- BLAST



# The biological problem

- Global and local alignment algoritms are slow in practice
- Consider the scenario of aligning a *query sequence* against a large database of sequences
  - New sequence with unknown function



- NCBI GenBank size in January 2007 was 65 369 091 950 bases (61 132 599 sequences)
- Feb 2008: 85 759 586 764 bases (82 853 685 sequences)



#### Problem with large amount of sequences

- Exponential growth in both number and total length of sequences
- Possible solution: Compare against model organisms only
- With large amount of sequences, chances are that matches occur by random
  - Need for statistical analysis



### Rapid alignment methods: FASTA and BLAST

- The biological problem
- Search strategies
- FASTA
- BLAST



### FASTA

- FASTA is a multistep algorithm for sequence alignment (Wilbur and Lipman, 1983)
- The sequence file format used by the FASTA software is widely used by other sequence analysis software
- Main idea:
  - Choose regions of the two sequences (query and database) that look promising (have some degree of similarity)
  - Compute local alignment using dynamic programming in these regions



# Search strategies

- How to speed up the computation?
  - Find ways to limit the number of pairwise comparisons
- Compare the sequences at *k-mer* level:
  - k-mer (or a k-word, or a k-tuple) is a substring of length k



#### **FASTA** outline

#### • FASTA algorithm has five steps:

- I. Identify common k-mers between I and J
- 2. Score diagonals with k-mer matches, identify 10 best diagonals
- 3. Rescore initial regions with a substitution score matrix
- 4. Join initial regions using gaps, penalise for gaps
- 5. Perform dynamic programming to find final alignments

# Analyzing the k-mer content

- Example query string I: TGATGATGAAGACATCAG
- For k = 8, the set of k-mers of I is

TGATGATG GATGATGA ATGATGAA TGATGAAG

. . .

GACATCAG



# Analyzing the k-mer content

- There are n-k+1 k-mers in a string of length n
- If at least one k-mer of I is not found from another string J, we know that I differs from J
- Need to consider statistical significance: I and J might share kmers by chance only
- Let m=|I| and n=|J|



#### Word lists and comparison by content

- The k-mers of I can be arranged into a table of k-mer occurences L<sub>w</sub>(I)
- Consider the k-mers when k=2 and I=GCATCGGC:
  GC, CA, AT, TC, CG, GG, GC

AT: 3  
CA: 2  
CG: 5  
GC: 1, 7 
$$\leftarrow$$
 Start indecies of k-mer GC in I  
GG: 6  
TC: 4 Building L<sub>w</sub>(I) takes O(n) time



# Common k-mers

- Number of common k-mers in I and J can be computed using L<sub>w</sub>(I) and L<sub>w</sub>(J)
- For each k-mer w in I, there are |L<sub>w</sub>(J)| occurrences in J
- Therefore I and J have  $\sum_w |L_w(I)| |L_w(J)|$  common k-mer pairs
- This can be computed in O(m + n + 4<sup>k</sup>) time
  - O(m + n + 4<sup>k</sup>) time to build the lists
  - O(4<sup>k</sup>) time to multiply the corresponding list entry sizes (in DNA strings)

29.9-1.10/

### Common k-mers

- I = GCATCGGC
- J = CCATCGCCATCG
  - $\begin{array}{cccc} L_w(I) & & L_w(J) \\ AT: 3 & AT: 3, 9 \\ CA: 2 & CA: 2, 8 \\ CC: 1, 7 \\ CG: 5 & CG: 5, 11 \\ GC: 1, 7 & GC: 6 \\ GG: 6 \\ TC: 4 & TC: 4, 10 \end{array}$

10 in total

26

29.9-1.10/

#### Properties of the common word list

- For large k, the table size O(4<sup>k</sup>) is too large to compute the common k-mer pairs count in the previous fashion
- Instead, an approach based on merge sort can be utilised (details skipped)
- The common k-mer technique can be combined with the local alignment algorithm to yield a rapid alignment approach
- Alternative solutions, not part of FASTA:
  - Exercise: Show how generalized suffix tree can be used for counting the common k-mer pairs
  - Generalized suffix tree with some additional data structures can also be used for directly computing all maximal pairs of tuples {(i',i),(j',j)} such that a<sub>i'</sub>...a<sub>i</sub> =b<sub>j'</sub>...b<sub>j</sub> and the ranges cannot be extended left or right (see Gusfield's book).



#### **FASTA** outline

#### FASTA algorithm has five steps:

- I. Identify common k-mers between I and J
- 2. Score diagonals with k-mer matches, identify 10 best diagonals
- 3. Rescore initial regions with a substitution score matrix
- 4. Join initial regions using gaps, penalise for gaps
- 5. Perform dynamic programming to find final alignments



#### Dot matrix comparisons

- k-mer matches in two sequences I and J can be represented as a dot matrix
- Dot matrix element (i, j) has "a dot", if the k-mer starting at position i in I is identical to the k-mer starting at position j in J

29.9-1.10/

29

The dot matrix can be plotted for various k







Dot matrix (k=1,4,8,16) for two protein sequences CAB51201.1 (531 aa) CAA72681.1 (588 aa)

k=16





k=8

k=1

Shading indicates now the match score according to a score matrix (Blosum62 here)

100 200 300 400 500 50-100 -150 200-250-300-350-400-450-500-550 29.9-1.10

frmatics, Autumn 2009

- We would like to find high scoring diagonals of the dot matrix
- Lets index diagonals by the offset, I = i j



 As an example, lets compute diagonal sums for I = GCATCGGC, J = CCATCGCCATCG, k = 2

- 1. Construct k-mer list L<sub>w</sub>(J)
- 2. Diagonal sums S<sub>I</sub> are computed into a table, indexed with the offset and initialised to zero

3. Go through k-mers of I, look for matches in  $L_w(J)$  and update diagonal sums



For the first 2-mer in I, GC,  $L_{GC}(J) = \{6\}$ .

We can then update the sum of diagonal I = i - j = 1 - 6 = -5 to  $S_{-5} := S_{-5} + 1 = 0 + 1 = 1$ 

34

3. Go through k-mers of I, look for matches in  $L_w(J)$  and update diagonal sums



Next 2-mer in I is CA, for which  $L_{CA}(J) = \{2, 8\}$ .

Two diagonal sums are updated:

3. Go through k-mers of I, look for matches in  $L_w(J)$  and update diagonal sums



Next 2-mer in I is AT, for which  $L_{AT}(J) = \{3, 9\}$ .

Two diagonal sums are updated:

$$= I - J = 3 - 3 = 0$$
  
S<sub>0</sub> := S<sub>0</sub> + 1 = 1 + 1 = 2

36
#### Computing diagonal sums



#### Algorithm for computing diagonal sum of scores

## **FASTA** outline

#### • FASTA algorithm has five steps:

- I. Identify common k-mers between I and J
- 2. Score diagonals with k-mer matches, identify 10 best diagonals
- 3. Rescore initial regions with a substitution score matrix
- 4. Join initial regions using gaps, penalise for gaps
- 5. Perform dynamic programming to find final alignments

#### **Rescoring initial regions**

- Each high-scoring diagonal chosen in the previous step is rescored according to a score matrix
- This is done to find subregions with identities shorter than k
- Non-matching ends of the diagonal are trimmed



# Joining diagonals

- Two offset diagonals can be joined with a gap, if the resulting alignment has a higher score
- Separate gap open and extension are used
- Find the best-scoring combination of diagonals



## **FASTA** outline

#### FASTA algorithm has five steps:

- I. Identify common k-mers between I and J
- 2. Score diagonals with k-mer matches, identify 10 best diagonals
- 3. Rescore initial regions with a substitution score matrix
- 4. Join initial regions using gaps, penalise for gaps
- 5. Perform dynamic programming to find final alignments

## Local alignment in the highest-scoring region

- Last step of FASTA: perform local alignment using dynamic programming around the highest-scoring diagonals
- Region to be aligned covers –w and +w offset diagonal to the highest-scoring diagonals
- With long sequences, this region is typically very small compared to the whole m x n matrix



Dynamic programming matrix M filled only for the green region



## **Properties of FASTA**

- Fast compared to local alignment using dynamic programming only
  - Only a narrow region of the full matrix is aligned
- Lossy filter : may fail to find some high scoring local alignments
- Increasing parameter k decreases the number of hits:
  - Increases specificity
  - Decreases sensitivity
  - Decreases running time

## Properties of FASTA

#### FASTA looks for initial exact matches to query sequence

- Two proteins can have very different amino acid sequences and still be biologically similar
- This may lead into a lack of sensitivity with diverged sequences

# **Demonstration of FASTA at EBI**

- http://www.ebi.ac.uk/fasta/
- Note that parameter ktup in the software corresponds to parameter k in lectures



## Rapid alignment methods: FASTA and BLAST

- The biological problem
- Search strategies
- FASTA
- BLAST



## **BLAST: Basic Local Alignment Search Tool**

- BLAST (Altschul et al., 1990) and its variants are some of the most common sequence search tools in use
- Roughly, the basic BLAST has three parts:
  - 1. Find segment pairs between the query sequence and a database sequence above score threshold ("seed hits")
  - 2. Extend seed hits into *locally maximal segment pairs*
  - 3. Calculate p-values and a rank ordering of the local alignments
- Gapped BLAST introduced in 1997 allows for gaps in alignments



### Finding seed hits

- First, we generate a set of neighborhood sequences for given k, match score matrix and threshold T
- Neighborhood sequences of a k-mer w include all strings of length k that, when aligned against w, have the alignment score at least T
- For instance, let I = GCATCGGC, J = CCATCGCCATCG and k = 5, match score be 1, mismatch score be 0 and T = 4

#### Finding seed hits

- I = GCATCGGC, J = CCATCGCCATCG, k = 5, match score 1, mismatch score 0, T = 4
- This allows for one mismatch in each k-mer
- The neighborhood of the first k-mer of I, GCATC, is GCATC and the 15 sequences

$$\begin{cases} A \\ CCATC, G \\ T \end{cases} \begin{cases} A \\ GATC, GC \\ T \end{cases} \begin{cases} C \\ GTC, GCA \\ GC \\ C \\ C \\ G \end{cases} \begin{cases} A \\ CC, GCAT \\ G \\ T \end{cases} \begin{cases} A \\ G \\ T \end{cases}$$

## Finding seed hits

- I = GCATCGGC has 4 k-mers and thus 4x16 = 64 5-mer patterns to locate in J
  - Occurrences of patterns in J are called seed hits
- Patterns can be found using exact search in time proportional to the sum of pattern lengths + length of J + number of matches (Aho-Corasick algorithm)
  - Attend 58093 String processing algorithms to learn Aho-Corasick and alike algorithms.

51

Compare this approach to FASTA

#### Extending seed hits: original BLAST

- Initial seed hits are extended into locally maximal segment pairs or High-scoring Segment Pairs (HSP)
- Extensions do not add gaps to the alignment
- Sequence is extended until the alignment score drops below the maximum attained score minus a threshold parameter value
- All statistically significant HSPs reported



Altschul, S.F., Gish, W., Miller, W., Myers, E. W. and Lipman, D. J., *J. Mol. Biol.*, 215, 403-410, 1990

582606 Introduction to Bioinformatics, Autumn 2009

#### Extending seed hits: gapped BLAST

- In a later version of BLAST, two seed hits have to be found on the same diagonal
  - Hits have to be non-overlapping
  - If the hits are closer than A (additional parameter), then they are joined into a HSP
- Threshold value T is lowered to achieve comparable sensitivity
- If the resulting HSP achieves a score at least S<sub>q</sub>, a *gapped extension* is triggered



Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, and Lipman DJ, *Nucleic Acids Res.* 1;25(17), 3389-402, 1997

582606 Introduction to Bioinformatics, Autumn 2009



#### Gapped extensions of HSPs

- Local alignment is performed starting from the HSP
- Dynamic programming matrix filled in "forward" and "backward" directions (see figure)
- Skip cells where value would be X<sub>g</sub> below the best alignment score found so far



29.9-1.10 / 54

# Estimating the significance of results

- In general, we have a score S(D, X) = s for a sequence X found in database D
- BLAST rank-orders the sequences found by p-values
- The p-value for this hit is P(S(D, Y) ≥ s) where Y is a random sequence
  - Measures the amount of "surprise" of finding sequence X
- A smaller p-value indicates more significant hit
  - A p-value of 0.1 means that one-tenth of random sequences would have as large score as our result

# Estimating the significance of results

- In BLAST, p-values are computed roughly as follows
- There are mn places to begin an optimal alignment in the m x n alignment matrix
- Optimal alignment is preceded by a mismatch and has t matching (identical) letters
  - (Assume match score 1 and mismatch/indel score -∞)
- Let p = P(two random letters are equal)
- The probability of having a mismatch and then t matches is (1-p)p<sup>t</sup>



# Estimating the significance of results

- We model this event by a Poisson distribution (why?) with mean λ = nm(1-p)p<sup>t</sup>
- P(there is local alignment t or longer)
  - $\approx$  1 P(no such event)

 $= 1 - e^{-\lambda} = 1 - \exp(-nm(1-p)p^{t})$ 

- An equation of the same form is used in Blast:
- E-value =  $P(S(D, Y) \ge s) \approx 1 exp(-mn\gamma\xi^t)$  where  $\gamma > 0$  and  $0 < \xi < 1$
- Parameters γ and ξ are estimated from data



# **Properties of BLAST**

- Better sensitivity than in FASTA
- Still a lossy filter:
  - Alternative lossless filters exist with similar performance!
    - See e.g. Lam et at.: Compressed indexing and local alignment of DNA, *Bioinformatics*, 25:1754-1760, 2008.
- Has become the standard in Bioinformatics:
  - This is due to the p-value computation and ranking of results, which give a way to compare the significance of the search results.
    - However, these computations apply to any alignment algorithm not just to BLAST.
    - BLAST may fail to find real occurrences, skewing the comparison of significance.

582606 Introduction to Bioinformatics, Autumn 2009



# Amino acid sequences

- We have discussed mainly DNA sequences
- Amino acid sequences can be aligned as well
- However, the design of the substitution matrix is more involved because of the larger alphabet
- More on the topic in the course Biological sequence analysis

# Scoring amino acid alignments

- We need a way to compute the score S(D, X) for aligning the sequence X against database D
- Scoring DNA alignments was discussed previously
- Constructing a scoring model for amino acids is more challenging
  - 20 different amino acids vs. 4 bases
- Figure shows the molecular structures of the 20 amino acids



http://en.wikipedia.org/wiki/List\_of\_standard\_amino\_acids

60

29.9-1.10 /

# Scoring amino acid alignments

- Substitutions between chemically similar amino acids are more frequent than between dissimilar amino acids
- We can check our scoring model against this



http://en.wikipedia.org/wiki/List\_of\_standard\_amino\_acids



# Score matrices

- Scores s = S(D, X) are obtained from score matrices
- Let A = a<sub>1</sub>a<sub>2</sub>...a<sub>n</sub> and B = b<sub>1</sub>b<sub>2</sub>...b<sub>n</sub> be sequences of equal length (no gaps allowed to simplify things)
- To obtain a score for alignment of A and B, where a<sub>i</sub> is aligned against b<sub>i</sub>, we take the ratio of two probabilities
  - The probability of having A and B where the characters match (*match model M*)
  - The probability that A and B were chosen randomly (random model R)

# Score matrices: random model

- Under the random model, the probability of having A and B is  $P(A, B|R) = \prod_i q_{ai} \prod_i q_{bi}$ 
  - where  $\mathbf{q}_{xi}$  is the probability of occurence of amino acid type  $\mathbf{x}_i$
- Position where an amino acid occurs does not affect its type

# Score matrices: match model

- Let p<sub>ab</sub> be the probability of having amino acids of type a and b aligned against each other given they have evolved from the same ancestor c
- The probability is

$$P(A, B|M) = \prod_{i} p_{a_i b_i}$$



# Score matrices: log-odds ratio score

We obtain the score S by taking the ratio of these two probabilities

$$\frac{P(A,B|M)}{P(A,B|R)} = \frac{\prod_i p_{a_i b_i}}{\prod_i q_{a_i} \prod_i q_{b_i}} = \prod_i \frac{p_{a_i b_i}}{q_{a_i} q_{b_i}}$$

and taking a logarithm of the ratio

$$S = \log_2 rac{P(A,B|M)}{P(A,B|R)} = \sum_{i=1}^n \log_2 rac{p_{a_i b_i}}{q_{a_i} q_{b_i}} = \sum_{i=1}^n s(a_i,b_i)$$

# Score matrices: log-odds ratio score

$$S = \log_2 rac{P(A,B|M)}{P(A,B|R)} = \sum_{i=1}^n \log_2 rac{p_{a_i b_i}}{q_{a_i} q_{b_i}} = \sum_{i=1}^n s(a_i,b_i)$$

The score S is obtained by summing over character pairspecific scores:

$$s(a,b) = \log_2 \frac{p_{ab}}{q_a q_b}$$

The probabilities q<sub>a</sub> and p<sub>ab</sub> are extracted from data

# Calculating score matrices for amino acids

- Probabilities q<sub>a</sub> are in principle easy to obtain:
- $s(a,b) = \log_2 \frac{p_{ab}}{q_a q_b}$
- Count relative frequencies of every amino acid in a sequence database



# Calculating score matrices for amino acids

- To calculate p<sub>ab</sub> we can use a known pool of aligned sequences
- BLOCKS is a database of highly conserved regions for proteins
- It lists multiply aligned, ungapped and conserved protein segments
- Example from BLOCKS shows genes related to human gene associated with DNA-repair defect xeroderma pigmentosum

$$s(a,b) = \log_2 \frac{p_{ab}}{q_a q_b}$$

#### Block PR00851A

```
ID XRODRMPGMNTB; BLOCK
AC PR00851A; distance from previous block=(52,131)
DE Xeroderma pigmentosum group B protein signature
BL adapted; width=21; seqs=8; 99.5%=985; strength=1287
XPB_HUMAN | P19447 ( 74)
                         RPLWVAPDGHIFLEAFSPVYK 54
XPB MOUSE | P49135 ( 74)
                         RPLWVAPDGHIFLEAFSPVYK 54
P91579 (80)
                         RPLYLAPDGHIFLESFSPVYK 67
XPB_DROME Q02870 (84)
                         RPLWVAPNGHVFLESFSPVYK 79
RA25 YEAST Q00578 (131) PLWISPSDGRIILESFSPLAE 100
038861 ( 52)
                          RPLWACADGRIFLETFSPLYK 71
013768 ( 90)
                          PLWINPIDGRIILEAFSPLAE 100
000835 (79)
                         RPIWVCPDGHIFLETFSAIYK 86
```

http://blocks.fhcrc.org



# **BLOSUM** matrix

- BLOSUM is a score matrix for amino acid sequences derived from BLOCKS data
- First, count pairwise matches f<sub>x,y</sub> for every amino acid type pair (x, y)
- For example, for column 3 and amino acids L and W, we find 8 pairwise matches: f<sub>L,W</sub> = f<sub>W,L</sub> = 8

RPLWVAPD RPLWVAPR RPLWVAPN PLWISPSD RPLWACAD PLWINPID RPLWVCPD



# Creating a BLOSUM matrix

Probability p<sub>ab</sub> is obtained by dividing f<sub>ab</sub> with the total number of pairs (note difference with course book):

$$p_{ab} = f_{ab} / \sum_{x=1}^{20} \sum_{y=1}^{x} f_{xy}$$

We get probabilities q<sub>a</sub> by

$$q_a = \sum_{b=1}^{20} p_{ab}$$

582606 Introduction to Bioinformatics, Autumn 2009



RP LWVAPD RP LWVAPR RP LWVAPN PL WISPSD RP LWACAD PL WINPID RP IWVCPD

# Creating a BLOSUM matrix

The probabilities p<sub>ab</sub> and q<sub>a</sub> can now be plugged into

$$s(a,b) = \log_2 \frac{p_{ab}}{q_a q_b}$$

to get a 20 x 20 matrix of scores s(a, b).

- Next slide presents the BLOSUM62 matrix
  - Values scaled by factor of 2 and rounded to integers
  - Additional step required to take into account expected evolutionary distance
  - Described in Deonier's book in more detail



#### BLOSUM62

\* A R N D С Q Ε G Η I L K Μ F Ρ S Т Υ В Ζ Х 0 4 -1 -2 -2 0 -2 -1 -1 -1 -1 -2 -1 -2 -1 Α -1 -1 1 0 -3 -20 0 - 4R -1 5 1 0 -2 2 -1 -3 -2 -3 0 0 -2 -3 0 -3 -2 -1 -3 - 2 -1 -1 -4 0 0 0 -3 -2 -3 N -2 0 6 1 -3 -3 -3 0 -2 1 0 -2 3 0 -1 -4 1 -4 D -2 -2 6 -3 0 2 -3 -3 -1 0 - 1 - 4-3 1 1 -1 -1 -3 -4 -1 -3 4 -1 -4 С 0 -3 -3 -3 -4 -3 -1 -2 -3 -1 -1 -2 -3 -2 -4 -3 9 -3 -1 -3 -2 -1 -3 -1 5 2 -2 0 -1 1 0 0 -3 0 -3 -20 -3 -1 0 -1 -2 -1 -20 3 -1 -4 1 2 E -1 5 -2 0 0 2 -4 0 -3 -3 1 -2 -3 -1 0 -1 -3 - 2 -2 1 4 -1 -4 -3 -2 -2 -3 -3 -2 G 0 -2 0 6 -2 -2 0 - 2-3 -3 -1 -2 -1 -4 -4 - 2 -1 -4 н -2 0 0 -2 -2 -1 -2 0 1 -1 -3 8 -3 -3 -1 - 2 2 -3 0 0 -1 -4 I -1 -3 -3 2 -3 1 0 -3 -2 3 -3 -1 -3 -3 -4 -3 4 -3 -1 -3 -3 -4 L -1 -2 -3 -4 -1 -2 -3 -4 -3 2 4 -2 2 0 -3 -2 -1 -2 -1 1 -4 - 4 5 K -1 2 0 -1 -3 1 1 -2 -1 -3 -2 -1 -3 -1 0 -1 -3 -2 -2 0 1 - 4 M -1 -1 -2 -3 -1 0 -2 -3 -2 1 2 -1 5 0 -2 -1 -1 -1 -1 1 -3 -1 -1 -4 0 0 -3 0 F - 2 - 3 - 3-3 -2 -3 -3 -3 -1 6 -4 -2 1 3 -3 -3 -1 -4 7 -1 -1 -4 -3 -1 -1 -2 -2 -3 -3 -1 -2 -3 -2 -1 -2 -4 P -1 -2 -2 -1 -4 -2 S 1 -1 0 0 0 0 -1 -2 -2 0 -2 -1 4 -3 -2 -2 0 0 0 - 41 -1 -1 1 Т 0 0 -1 0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1 1 5 -2 -2 0 -1 -1 - 4 -2 -3 -2 -2 -3 -2 -3 -1 2 W - 3 - 3-4 -4 -2 1 -4 -3 -2 11 -3 -4 -3 -2 -4 2 3 - 3 2 7 -1 -1 -4 Y -2 -2 - 2. -3 -2 -1 -2 -3 -1 -1 -2 -1 -2 -2 -3 -2 V -3 -3 -3 -1 -2 -2 -3 -3 3 1 -2 1 -1 -2 -2 0 -3 -1 4 -3 -2 -1 -4 B -2 -1 3 4 -3 0 1 -1 0 -3 -4 0 -3 -3 -2 0 - 1 - 4-3 -3 4 1 -1 -4 Ζ 0 0 -3 3 4 -2 -3 -1 0 -1 -3 -2 -2 1 0 -3 -3 1 1 - 4 Х 0 -1 -2 -1 -1 -1 -1 -1 -2 0 0 -2 -1 -1 -1 -1 -1 -4 \_ 1 - 1 -1 -1 -1 -1 \* -4 -4 -4 -4 -4 1 29.9-1.10/ 72 582606 Introduction to Bioinformatics, Autumn 2009
## Using BLOSUM62 matrix

MQLEANADTSV | | | LQEQAEAQGEM

$$s = \sum_{i=1}^{11} s(a_i, b_i)$$
  
= 2 + 5 - 3 - 4 + 4 + 0 + 4 + 0 - 2 + 0 + 1  
= 7

582606 Introduction to Bioinformatics, Autumn 2009



## **Demonstration of BLAST at NCBI**

http://www.ncbi.nlm.nih.gov/BLAST/

582606 Introduction to Bioinformatics, Autumn 2009



## Signals in DNA (Section 9)

- Genes
- Promoter regions
- Binding sites for regulatory proteins (*transcription* factors, motifs)



# Gene finding

- Could be done by dynamic programming similar to exercise 3/6 taking into account
  - Sequencing errors
  - Exons / intron gap constraints
  - Total gene length constraints
  - Codon usage optimization
  - Start codons / stop codons.
- In practise, often just the consensus reverse translation is taken and BLAST is used for finding local alignments (exons).



### Promoter sequences

- Often immediately before the gene.
- Clear structure in prokaryotes, more complex in eukaryotes.
- An example from *E coli* is shown in next slide (taken from course book).



### Promoter example

Table 9.2. A sample of *E. coli* promoter sequences. These sequences have been aligned relative to the transcriptional start site at position +1 (boldface large letter). Sequences from -40 to +11 are shown. Close matches to consensus -35 and -10 hexamers are underlined. See also Appendix C.3 for additional examples and sources of the data.

	-35		-10	-1	
ORF83P1	I		1		
	CTCTGCTGGCA <u>TTCACA</u> A	ATGCGCAGGGG	TAAAACGT	TTCCTGTAGCA	CCG
ada			00000		001
D.4	GTTGGTTTTTGCGTGATG	FIGACCGGGCA	GCCTAAAG	GCTATUCITAA	UCA
amnP4	TTCACATTTCTCTCACAT	ACTATCCCATC	TCCCCTAN	TTGTATGGAAC	204
araFCH	IIGAGAIIIGI <u>GIGAGA</u> I	ACTATOGGATO	10000177	IIIIIAIOOANO	1144
unui (m	CTCTCCTATGGAGAATTA	ATTTCTCGCTA	AAACTATG	TCAACACAGTO	CACT
aroG	••••				
	CCCCG <u>TTTACA</u> CATTCTG/	ACGGAAGATA <u>T</u>	AGATT GGA	AGTATTGCATT	CAC
atpI					
. —	TATTGTTTGAAATCACGG	GGGCGCACCG <u>T</u>	<u>'ATAAT</u> TTG	ACCGCTTTTTC	ATG
caiT					
J= 4D1	AATCACAGAATACAGCTT	ATTGAATACC <u>C</u>	<u>ATTAT</u> GAG	TIAGCCATIA	IÇGÜ
cipAP I	TTATTCACCTCTTACAAA	8 8 TTOTTTOT	TATCATOT	AGA A COTOCAL	LCGC
crrP2-I	TTA <u>TIONOU</u> TOTIAOAAA	nn: 10111101	1111061 01	Noni Logi doni	
<b>.</b> .	GTGGTGAGCTTGCTGGCG.	ATGAACGTGCT	ACACTTCT	GTTGCTGGGG	ATGG
		-			



## Representing signals in DNA

Consensus sequence: -10 site in E coli: TATAAT GRE half-site consensus: AGAACA Simple regular expression: A(C/G)AA(C/G)(A/T) Positional weight matrix (PWM): A 1.00 0.00 1.00 1.00 0.00 0.86 C 0.00 0.14 0.00 0.00 0.86 0.00 **G** 0.00 0.86 0.00 0.00 0.14 0.00 T|0.00 0.00 0.00 0.00 0.00 0.14

GRE half-sites: AGAACA ACAACA AGAACA AGAAGA AGAACA AGAACA AGAACT <u>AGAACA</u> AGAACA



# Position-specific scoring matrix (PSSM)

- PSSM is a log-odds normalized version of PWM.<sup>1</sup>
- Calculated by log(p<sub>ai</sub>/q<sub>a</sub>), where
  - p<sub>ai</sub> is the frequency of a at column i in the samples.
  - q<sub>a</sub> is the probability of a in the whole organism (or in some region of interest).
- Problematic when some values p<sub>ai</sub> are zero.
- Solution is to use pseudocounts:
  - add 1 to all the sample counts where the frequencies are calculated.

<sup>1</sup> In the following log denotes base 2 logarithm.



### PWM versus PSSM



### **Extended representations**

- PWM representation can be extended to first-order *non-homogenous* Markov chain.
- Identical to what described in lecture 2, except that transition probabilities are position-specific.



Fig. 9.4. Contrast between the Markov chain representation for positions with identical probability distributions (panel A) and the model for positions in a signal (binding site) sequence (panel B). Numbered four-element horizontal rows are the vectors of probabilities for A, C, G, and T at each position. These vectors are transformed to vectors at the next position by matrix multiplication by a transition matrix. In panel A, a single transition matrix is employed. In panel B, a different transition matrix is required for each successive position.



## Sequence logos

- Many known transcription factor binding site PWM:s can be found from JASPAR database (http://jaspar.cgb.ki.se/).
- PWM:s are visualized as sequence logos, where the height of each nucleotide equals its proportion of the relative entropy (expected log-odds score) in that column.

83

- $E(S_i) = \sum_a p_{ai} \log(p_{ai} / q_a)$
- Height of a at column i is  $p_{ai}E(S_i)$

#### Example sequence logo





## More on signals in DNA

- Book by Durbin et al.: Biological sequence analysis.
  - Hidden Markov Model (HMM) approach to alignment.
  - Viterbi / forward-backward dynamic programming algorithms for finding most probable paths in the model.
  - Profile-HMMs for sequence family classification.
- 582653 Computational methods of systems biology (4 credits)
  - Professor Esko Ukkonen, II period
  - The course is an advanced introduction to computational methods for analysing genomic and gene expression data to find different functional units (such as genes) and regulatory structures and relations (such as gene enhancers).

