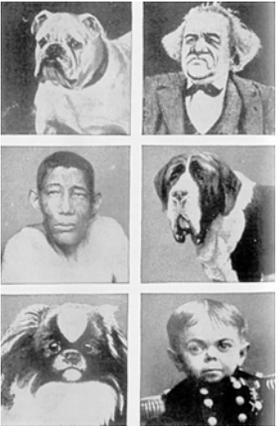


November 22, 2010

Pairwise sequence alignment

Jonathan Pevsner, Ph.D.
Bioinformatics
Johns Hopkins School Med.



Copyright notice

Many of the images in this powerpoint presentation are from *Bioinformatics and Functional Genomics* by Jonathan Pevsner (ISBN 0-471-21004-8). Copyright © 2009 by John Wiley & Sons, Inc.

These images and materials may not be used without permission from the publisher. We welcome instructors to use these powerpoints for educational purposes, but please acknowledge the source.

The book has a homepage at <http://www.bioinfbook.org> including hyperlinks to the book chapters.

Announcements

The moodle quiz from lecture 1 is due one week later—by today at noon. After then the quiz “closes” and won’t be available to you.

The quiz from today’s lecture (“opens” at 10:30 am) is due in one week later at noon. Because of the Thanksgiving break, I’m extending the deadline a day to Tuesday November 30 (5:00 pm).

Outline: pairwise alignment

- Overview and examples
- Definitions: homologs, paralog, orthologs
- Assigning scores to aligned amino acids: Dayhoff’s PAM matrices
- Alignment algorithms: Needleman-Wunsch, Smith-Waterman

Learning objectives

- Define homologs, paralog, orthologs
- Perform pairwise alignments (NCBI BLAST)
- Understand how scores are assigned to aligned amino acids using Dayhoff’s PAM matrices
- Explain how the Needleman-Wunsch algorithm performs global pairwise alignments

Pairwise alignments in the 1950s

β -corticotropin (sheep)	ala gly glu asp asp glu
Corticotropin A (pig)	asp gly ala glu asp glu

Oxytocin	CYIQNCPLG
Vasopressin	CYFQNCPRG

globins: α - β - myoglobin

Early example of sequence alignment: globins (1961)

H.C. Watson and J.C. Kendrew, "Comparison Between the Amino-Acid Sequences of Sperm Whale Myoglobin and of Human Haemoglobin." *Nature* 190:670-672, 1961.

Pairwise sequence alignment is the most fundamental operation of bioinformatics

- It is used to decide if two proteins (or genes) are related structurally or functionally
- It is used to identify domains or motifs that are shared between proteins
- It is the basis of BLAST searching (next week)
- It is used in the analysis of genomes

Page 47

Pairwise alignment: protein sequences can be more informative than DNA

- protein is more informative (20 vs 4 characters); many amino acids share related biophysical properties
- codons are degenerate: changes in the third position often do not alter the amino acid that is specified
- protein sequences offer a longer "look-back" time
- DNA sequences can be translated into protein, and then used in pairwise alignments

		Second letter				
		U	C	A	G	
First letter	U	UUU } Phe UUC } UUA } Leu UUG }	UCU } UCC } Ser UCA } UCG }	UAU } Tyr UAC } UAA } Stop UAG } Stop	UGU } Cys UGC } UGA } Stop UGG } Trp	U C A G
	C	CUU } CUC } Leu CUA } CUG }	CCU } CCC } Pro CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } CGC } Arg CGA } CGG }	U C A G
	A	AUU } AUC } Ile AUA } AUG } Met	ACU } ACC } ACA } Thr ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }	U C A G
	G	GUU } GUC } Val GUA } GUG }	GCU } GCC } Ala GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } GGC } GGA } Gly GGG }	U C A G

Page 54

Pairwise alignment: protein sequences can be more informative than DNA

- Many times, DNA alignments are appropriate
 - to confirm the identity of a cDNA
 - to study noncoding regions of DNA
 - to study DNA polymorphisms
 - example: Neanderthal vs modern human DNA

```

Query: 181 catcaactacaactccaaagacccttacaaccactagatcaacaacctaaccac 240
          |||||  |||||  |||||  |||||  |||||  |||||  |||||  |||||  |||||  |||||
Sbjct: 189 catcaactgcaacccaagaccacct-caccactagatatacaacaacctaaccac 247
    
```

Outline: pairwise alignment

- Overview and examples
- Definitions: homologs, paralogs, orthologs
- Assigning scores to aligned amino acids: Dayhoff's PAM matrices
- Alignment algorithms: Needleman-Wunsch, Smith-Waterman

Definition: pairwise alignment

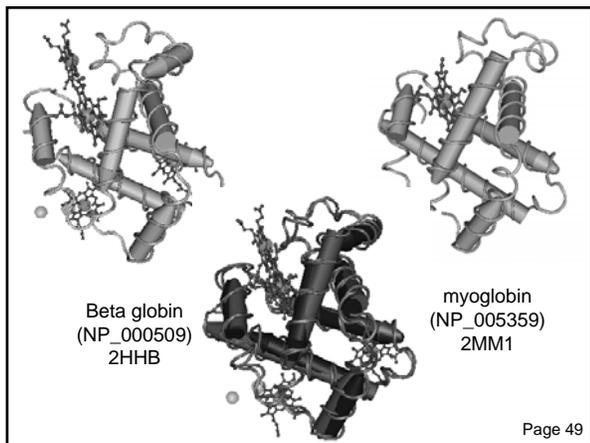
Pairwise alignment

The process of lining up two sequences to achieve maximal levels of identity (and conservation, in the case of amino acid sequences) for the purpose of assessing the degree of similarity and the possibility of homology.

Definition: homology

Homology

Similarity attributed to descent from a common ancestor.



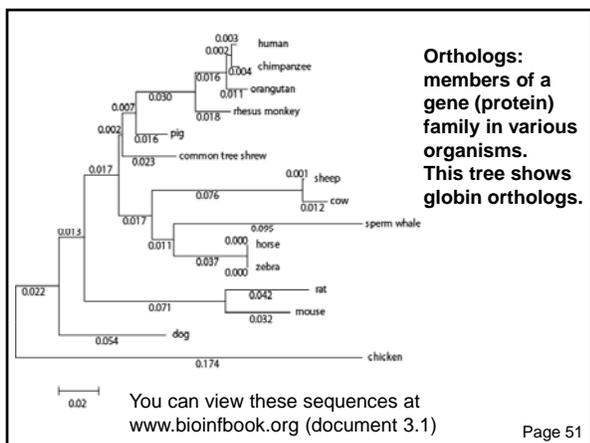
Definitions: two types of homology

Orthologs

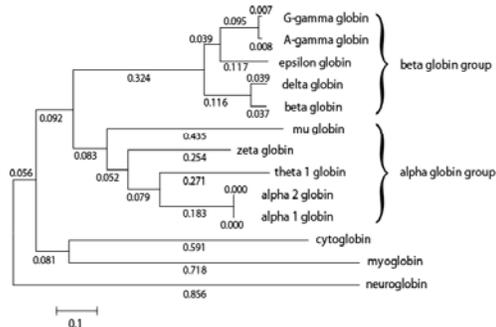
Homologous sequences in different species that arose from a common ancestral gene during speciation; may or may not be responsible for a similar function.

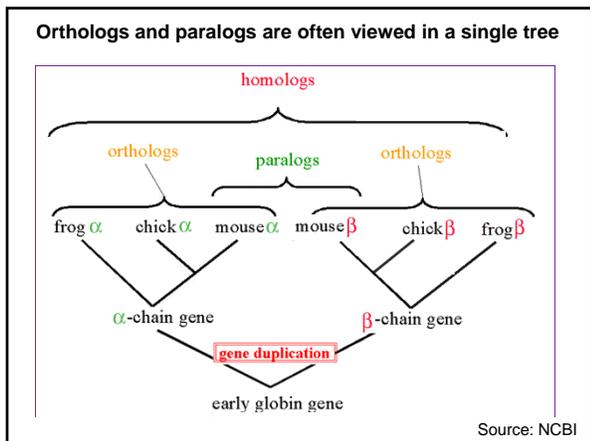
Paralogs

Homologous sequences within a single species that arose by gene duplication.



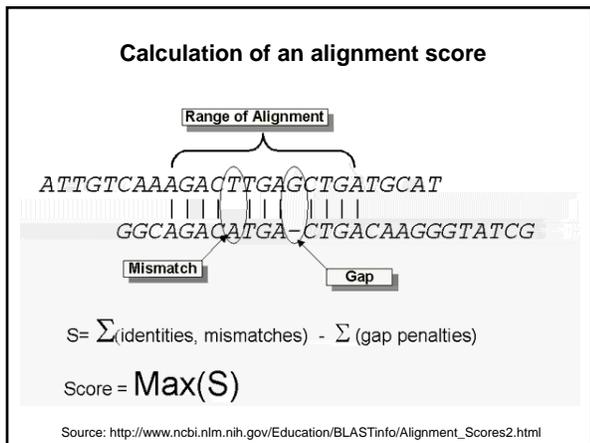
Paralogs: members of a gene (protein) family within a species. This tree shows human globin paralogs.





General approach to pairwise alignment

- Choose two sequences
- Select an algorithm that generates a score
- Allow gaps (insertions, deletions)
- Score reflects degree of similarity
- Alignments can be global or local
- Estimate probability that the alignment occurred by chance



Popular Resources

- PubMed
- PubMed Central
- NCBI Bookshelf
- BLAST**
- GenBank
- Nucleotide
- Protein
- GFO
- Conserved Domain

Find BLAST from the home page of NCBI and select protein BLAST...

Page 52

Choose align two or more sequences...

Page 52

Enter the two sequences (as accession numbers or in the fasta format) and click BLAST.

Optionally select "Algorithm parameters" and note the matrix option.

Pairwise alignment result of human beta globin and myoglobin

Myoglobin RefSeq

Information about this alignment: score, expect value, identities, positives, gaps...

```
>|ref|NP_005359.1|G| myoglobin [Homo sapiens]
ref|NP_076311.1|G| myoglobin [Homo sapiens]
ref|NP_076312.1|G| myoglobin [Homo sapiens]
>|l| myoglobin [Homo sapiens]
Length=154
GENE ID: 4151 MB | myoglobin [Homo sapiens] (Over 10 Pubmed links)
Score = 47.4 bits (144), Expect = 8e-11, Method: Compositional matrix adjust.
Identities = 37/145 (25%), Positives = 57/145 (39%), Gaps = 2/145 (1%)
Query 4  LTPEEKSAVTLNGKVVWDEVG--GEALGRLLVYFVTFURFFESFGDLSTPDVWGNPKV 61
      1  L E V +NGKV D G E L RL +S T F+ F L + D + + +
Sbjct 3  LSDGEWQLVLMVWGVKVEADIPGHGQEVLRIRLFGHFETLEKDFKFKLKSDEMKASEDL 62
Query 62  KAHGKVVLAGFSDGLAHLMLKGTFTATLSELNCDKLFVDPENFRLLGNVLCVLAHDFK 121
      K H G VL A L + + L + H K + + + ++ VL
Sbjct 63  KGHGATVLTALGGLKGGHAEAEIKFLAQSHATKIPVKYLEFISECIQVLSKHPG 122
Query 122 EFTFPVQAAYKQVAVANALAHKY 146
      +F Q A K + +A Y
Sbjct 123 DFGADAGAMHKALELFRDHAKY 147
```

Query = HBB
Subject = MB

Middle row displays identities; + sign for similar matches

Page 53

Pairwise alignment result of human beta globin and myoglobin: the score is a sum of match, mismatch, gap creation, and gap extension scores

Score = 18.1 bits (35), Expect = 0.015, Method: Composition-based stats. Identities = 11/24 (45%), Positives = 12/24 (50%), Gaps = 2/24 (8%)

```
Query 12 VTALMGKVVND--EVGGEALGRLL 33
      V +NGKV D G E L RL
Sbjct 11 VLMVWGVKVEADIPGHGQEVLRIRLF 34
```

match	4	11	5	6	6	5	4	5	sum of matches: +60
				6	4			4	
mismatch	-1	1	0	-2	-2	-4	0	0	sum of mismatches: -13
	-2		0	-3	0				
gap open								-11	sum of gap penalties: -12
gap extend								-1	
									total raw score: 60 - 13 - 12 = 35

Page 53

Pairwise alignment result of human beta globin and myoglobin: the score is a sum of match, mismatch, gap creation, and gap extension scores

Score = 18.1 bits (35), Expect = 0.015, Method: Composition-based stats. Identities = 11/24 (45%), Positives = 12/24 (50%), Gaps = 2/24 (8%)

```
Query 12 VTALMGKVVND--EVGGEALGRLL 33
      V +NGKV D G E L RL
Sbjct 11 VLMVWGVKVEADIPGHGQEVLRIRLF 34
```

match	4	11	5	6	6	5	4	5	sum of matches: +60
				6	4			4	
mismatch	-1	1	0	-2	-2	-4	0	0	sum of mismatches: -13
	-2		0	-3	0				
gap open								-11	sum of gap penalties: -12
gap extend								-1	
									total raw score: 60 - 13 - 12 = 35

V matching V earns +4
T matching L earns -1

These scores come from a "scoring matrix"!

Page 53

Definitions: homology

Homology
Similarity attributed to descent from a common ancestor.

Page 50

Definitions: identity, similarity, conservation

Identity
The extent to which two (nucleotide or amino acid) sequences are invariant.

Similarity
The extent to which nucleotide or protein sequences are related. It is based upon identity plus conservation.

Conservation
Changes at a specific position of an amino acid or (less commonly, DNA) sequence that preserve the physico-chemical properties of the original residue.

Page 51

Definition: pairwise alignment

Pairwise alignment

The process of lining up two sequences to achieve maximal levels of identity (and conservation, for amino acid sequences) for the purpose of assessing the degree of similarity and the possibility of homology.

Page 53

Outline: pairwise alignment

- Overview and examples
- Definitions: homologs, paralogs, orthologs
- Assigning scores to aligned amino acids: Dayhoff's PAM matrices
- Alignment algorithms: Needleman-Wunsch, Smith-Waterman

Substituent residue
(Percentage of total residue sites at which the substituent occurs)

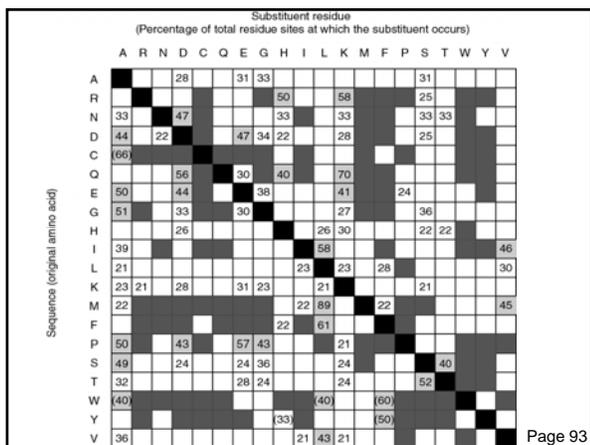
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A																				
R																				
N																				

lys found at 58% of arg sites

Emile Zuckerkandl and Linus Pauling (1965) considered substitution frequencies in 18 globins (myoglobins and hemoglobins from human to lamprey).

Black: identity
 Gray: very conservative substitutions (>40% occurrence)
 White: fairly conservative substitutions (>21% occurrence)
 Red: no substitutions observed

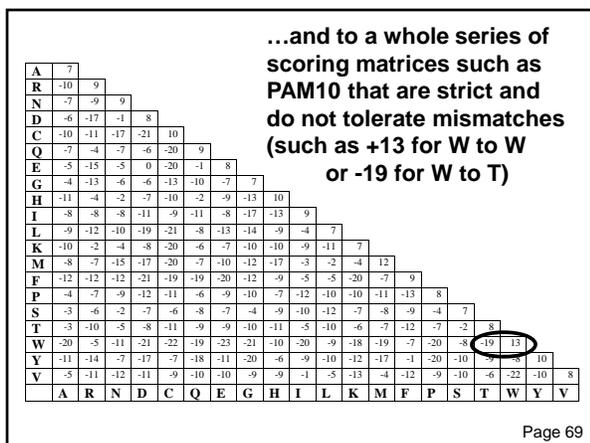
Page 93



Where we're heading:
to a PAM250 log odds scoring matrix that assigns scores and is forgiving of mismatches... (such as +17 for W to R or -5 for W to T)

A	2																			
R	-2	6																		
N	0	0	2																	
D	0	-1	2	4																
C	-2	-4	-4	-5	12															
Q	0	1	1	2	-5	4														
E	0	-1	1	3	-5	2	4													
G	1	-3	0	1	-3	-1	0	5												
H	-1	2	2	1	-3	3	1	-2	6											
I	-1	-2	-2	-2	-2	-2	-3	-2	5											
L	-2	-3	-3	-4	-6	-2	-3	-4	-2	-2										
K	-1	3	1	0	-5	1	0	-2	0	-2	-3	5								
M	-1	0	-2	-3	-5	-1	-2	-3	-2	2	4	0	6							
F	-3	-4	-3	-6	-4	-5	-5	-2	1	2	-5	0	9							
P	1	0	0	-1	-3	0	-1	0	0	-2	-3	-1	-2	-5	6					
S	1	0	1	0	0	-1	0	1	-1	-1	-3	0	-2	-3	1	2				
T	1	-1	0	0	-2	-1	0	0	-1	0	-2	0	-1	-3	0	1	3			
W	-6	2	-4	-7	-8	-5	-7	-7	-3	-5	-2	-3	-4	0	-6	-2	-5	17		
Y	-3	-4	-2	-4	0	-4	-4	-5	0	-1	-1	-4	-2	7	-5	-3	-3	0	10	
V	0	-2	-2	-2	-2	-2	-1	-2	4	2	-2	2	-1	-1	0	-6	-2	4		

Page 69



Dayhoff's 34 protein superfamilies

Protein	PAMs per 100 million years
Ig kappa chain	37
Kappa casein	33
luteinizing hormone b	30
lactalbumin	27
complement component 3	27
epidermal growth factor	26
proopiomelanocortin	21
pancreatic ribonuclease	21
haptoglobin alpha	20
serum albumin	19
phospholipase A2, group IB	19
prolactin	17
carbonic anhydrase C	16
Hemoglobin α	12
Hemoglobin β	12

Page 59

Substitution Matrix

A substitution matrix contains values proportional to the probability that amino acid *i* mutates into amino acid *j* for all pairs of amino acids.

Substitution matrices are constructed by assembling a large and diverse sample of verified pairwise alignments (or multiple sequence alignments) of amino acids.

Substitution matrices should reflect the true probabilities of mutations occurring through a period of evolution.

The two major types of substitution matrices are PAM and BLOSUM.

PAM matrices: Point-accepted mutations

PAM matrices are based on global alignments of closely related proteins.

The PAM1 is the matrix calculated from comparisons of sequences with no more than 1% divergence. At an evolutionary interval of PAM1, one change has occurred over a length of 100 amino acids.

Other PAM matrices are extrapolated from PAM1. For PAM250, 250 changes have occurred for two proteins over a length of 100 amino acids.

All the PAM data come from closely related proteins (>85% amino acid identity).

Page 63

Dayhoff's PAM1 mutation probability matrix

	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly	H His
A	9867	2	9	10	3	8	17	21	2
R	1	9913	1	0	1	10	0	0	10
N	4	1	9822	36	0	4	6	6	21
D	6	0	42	9859	0	6	53	6	4
C	1	1	0	0	9973	0	0	0	1
Q	3	9	4	5	0	9876	27	1	23
E	10	0	7	56	0	35	9865	4	2
G	21	1	12	11	1	3	7	9935	1
H	1	8	18	3	1	20	1	0	9912
I	2	2	3	1	2	1	2	0	0

Page 66

Dayhoff's PAM0 mutation probability matrix: the rules for extremely slowly evolving proteins

PAM0	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu
A	100%	0%	0%	0%	0%	0%	0%
R	0%	100%	0%	0%	0%	0%	0%
N	0%	0%	100%	0%	0%	0%	0%
D	0%	0%	0%	100%	0%	0%	0%
C	0%	0%	0%	0%	100%	0%	0%
Q	0%	0%	0%	0%	0%	100%	0%
E	0%	0%	0%	0%	0%	0%	100%
G	0%	0%	0%	0%	0%	0%	0%

Top: original amino acid
Side: replacement amino acid

Page 68

Dayhoff's PAM2000 mutation probability matrix: the rules for very distantly related proteins

PAM ∞	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly
A	8.7%	8.7%	8.7%	8.7%	8.7%	8.7%	8.7%	8.7%
R	4.1%	4.1%	4.1%	4.1%	4.1%	4.1%	4.1%	4.1%
N	4.0%	4.0%	4.0%	4.0%	4.0%	4.0%	4.0%	4.0%
D	4.7%	4.7%	4.7%	4.7%	4.7%	4.7%	4.7%	4.7%
C	3.3%	3.3%	3.3%	3.3%	3.3%	3.3%	3.3%	3.3%
Q	3.8%	3.8%	3.8%	3.8%	3.8%	3.8%	3.8%	3.8%
E	5.0%	5.0%	5.0%	5.0%	5.0%	5.0%	5.0%	5.0%
G	8.9%	8.9%	8.9%	8.9%	8.9%	8.9%	8.9%	8.9%

Top: original amino acid
Side: replacement amino acid

Page 68

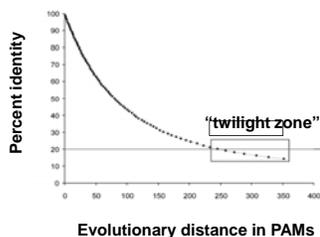
PAM250 mutation probability matrix

A	R	N	D	C	Q	E	H	I	L	K	M	F	P	S	T	W	Y	V			
A	33	6	9	9	5	8	9	15	6	8	6	7	7	4	11	11	2	4	9		
R	2	17	4	3	2	5	3	2	6	3	2	9	4	1	4	4	3	2	2		
N	4	4	6	7	2	5	6	4	6	3	2	5	3	2	4	5	4	2	3	3	
D	5	4	8	11	1	7	10	5	6	3	2	5	3	1	4	5	5	1	2	3	
C	2	1	1	1	52	1	1	2	2	2	1	1	1	1	2	3	2	1	4	2	
Q	3	5	5	6	1	10	7	3	7	2	3	5	3	1	4	3	3	1	2	3	
E	5	4	7	11	1	9	12	5	6	3	2	5	3	1	4	5	5	1	2	3	
H	2	5	5	4	2	7	4	2	15	2	2	2	2	2	3	3	2	2	3	2	
I	3	2	2	2	2	2	2	2	10	6	2	6	5	2	3	4	1	3	9		
L	6	4	4	3	2	6	4	3	5	15	34	4	20	13	5	4	6	6	7	13	
K	6	18	10	8	2	10	8	5	8	5	4	24	9	2	6	8	8	4	3	5	
M	1	1	1	1	0	1	1	1	1	1	2	3	2	6	2	1	1	1	1	2	
F	2	1	2	1	1	1	1	1	3	5	6	1	4	32	1	2	2	4	20	3	
P	7	5	5	4	3	5	4	5	5	3	3	4	3	2	20	6	5	1	2	4	
S	6	8	7	7	6	7	9	6	5	4	7	5	3	9	10	9	4	4	6		
T	5	5	6	6	4	5	5	6	4	6	4	6	5	3	6	8	11	2	3	6	
W	0	2	0	0	0	0	0	0	1	0	1	0	0	1	0	1	0	55	1	0	
Y	1	1	2	1	3	1	1	1	3	2	2	1	2	15	1	2	2	3	31	2	
V	7	4	4	4	4	4	4	4	5	4	15	10	4	10	5	5	5	7	2	4	17

Top: original amino acid
Side: replacement amino acid

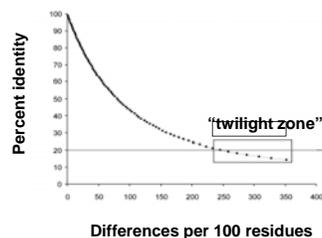
Page 68

Two randomly diverging protein sequences change in a negatively exponential fashion



Page 74

At PAM1, two proteins are 99% identical
 At PAM10.7, there are 10 differences per 100 residues
 At PAM80, there are 50 differences per 100 residues
 At PAM250, there are 80 differences per 100 residues



Page 75

PAM: "Accepted point mutation"

- Two proteins with 50% identity may have 80 changes per 100 residues. (Why? Because any residue can be subject to back mutations.)
- Proteins with 20% to 25% identity are in the "twilight zone" and may be statistically significantly related.
- PAM or "accepted point mutation" refers to the "hits" or matches between two sequences (Dayhoff & Eck, 1968)

Page 75

Outline: pairwise alignment

- Overview and examples
- Definitions: homologs, paralogs, orthologs
- Assigning scores to aligned amino acids: Dayhoff's PAM matrices
- Alignment algorithms: Needleman-Wunsch, Smith-Waterman

Two kinds of sequence alignment: global and local

We will first consider the global alignment algorithm of Needleman and Wunsch (1970).

We will then explore the local alignment algorithm of Smith and Waterman (1981).

Finally, we will consider BLAST, a heuristic version of Smith-Waterman. We will cover BLAST in detail on Monday.

Page 76

Global alignment with the algorithm of Needleman and Wunsch (1970)

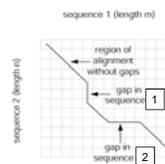
- Two sequences can be compared in a matrix along x- and y-axes.
- If they are identical, a path along a diagonal can be drawn
- Find the optimal subpaths, and add them up to achieve the best score. This involves
 - adding gaps when needed
 - allowing for conservative substitutions
 - choosing a scoring system (simple or complicated)
- N-W is guaranteed to find optimal alignment(s)

Page 76

Three steps to global alignment with the Needleman-Wunsch algorithm

- [1] set up a matrix
- [2] score the matrix
- [3] identify the optimal alignment(s)

Four possible outcomes in aligning two sequences



- [1] identity (stay along a diagonal)
- [2] mismatch (stay along a diagonal)
- [3] gap in one sequence (move vertically!)
- [4] gap in the other sequence (move horizontally!)

<p>Sequence 1</p> <table border="1" style="font-size: small;"> <tr><td></td><td></td><td>D</td><td>P</td><td>L</td><td>E</td></tr> <tr><td>D</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>P</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>L</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>E</td><td></td><td></td><td></td><td></td><td></td></tr> </table> <p>1 DPLE 2 DPLE</p>			D	P	L	E	D						P						L						E						<p>Sequence 1</p> <table border="1" style="font-size: small;"> <tr><td></td><td></td><td>D</td><td>P</td><td>M</td><td>E</td></tr> <tr><td>D</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>P</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>L</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>E</td><td></td><td></td><td></td><td></td><td></td></tr> </table> <p>1 DPLE 2 DPME</p>			D	P	M	E	D						P						L						E					
		D	P	L	E																																																								
D																																																													
P																																																													
L																																																													
E																																																													
		D	P	M	E																																																								
D																																																													
P																																																													
L																																																													
E																																																													
<p>Sequence 1</p> <table border="1" style="font-size: small;"> <tr><td></td><td></td><td>D</td><td>P</td><td>L</td><td>E</td></tr> <tr><td>D</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>P</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>L</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>E</td><td></td><td></td><td></td><td></td><td></td></tr> </table> <p>1 DPLE 2 DP-E</p>			D	P	L	E	D						P						L						E						<p>Sequence 1</p> <table border="1" style="font-size: small;"> <tr><td></td><td></td><td>D</td><td>P</td><td>L</td><td>E</td></tr> <tr><td>D</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>P</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>L</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>E</td><td></td><td></td><td></td><td></td><td></td></tr> </table> <p>1 D-L-E 2 DPLE</p>			D	P	L	E	D						P						L						E					
		D	P	L	E																																																								
D																																																													
P																																																													
L																																																													
E																																																													
		D	P	L	E																																																								
D																																																													
P																																																													
L																																																													
E																																																													

Start Needleman-Wunsch with an identity matrix

(e)

		Sequence 2 (from honeybee globin)							
		F	M	D	T	P	L	N	E
Sequence 1 (from human cytoglobin)	F	1							
	K								
	H								
	M		1						
	E								1
	D			1					
	P					1			
	L						1		
E								1	

Start Needleman-Wunsch with an identity matrix

		Sequence 2							
		F	M	D	T	P	L	N	E
Sequence 1	F	6	0	-3	-2	-4	0	-3	-3
	K	-3	-1	-1	-1	-1	-2	0	1
	H	-1	-2	-1	-2	-2	-3	1	0
	M	0	5	-3	-1	-2	2	-2	-2
	E	-3	-2	2	-1	-1	-3	0	5
	D	-3	-3	6	-1	-1	-4	1	2
	P	-4	-2	-1	-1	7	-3	-2	-1
	L	0	2	-4	-1	-3	4	-3	-3
E	-3	-2	2	-1	-1	-3	0	5	

Fill in the matrix using "dynamic programming"

(a) Sequence 2: F M D T P L N E

	F	M	D	T	P	L	N	E
F	6	0	-3	-2	-4	0	-3	-3
K	-3	-1	-1	-1	-1	-2	0	1
H	-1	-2	-1	-2	-2	-3	1	0
M	0	5	-3	-1	-2	2	-2	-2
E	-3	-2	2	-1	-1	-3	0	5
D	-3	-3	6	-1	-1	-4	1	2
P	-4	-2	-1	-1	7	-3	-2	-1
L	0	2	-4	-1	-3	4	-3	-3
E	-3	-2	2	-1	-1	-3	0	5

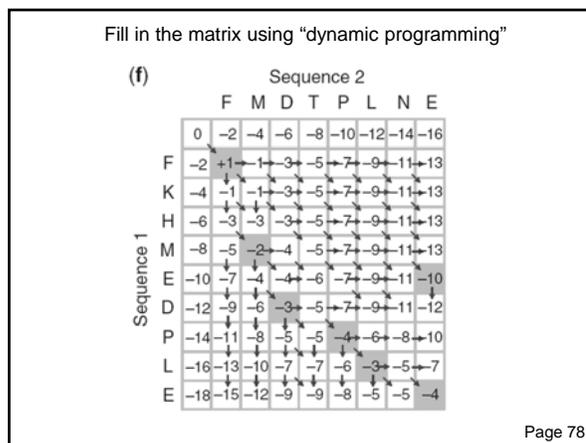
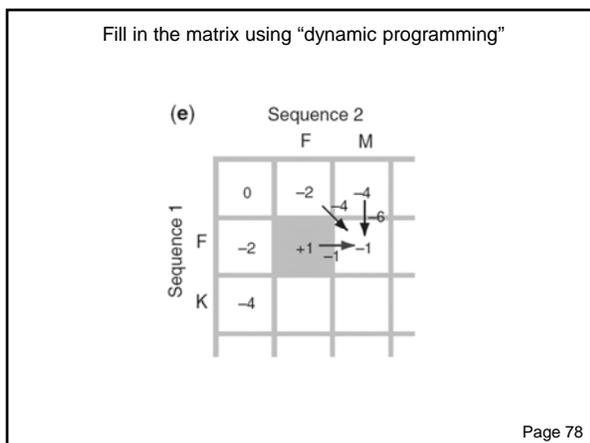
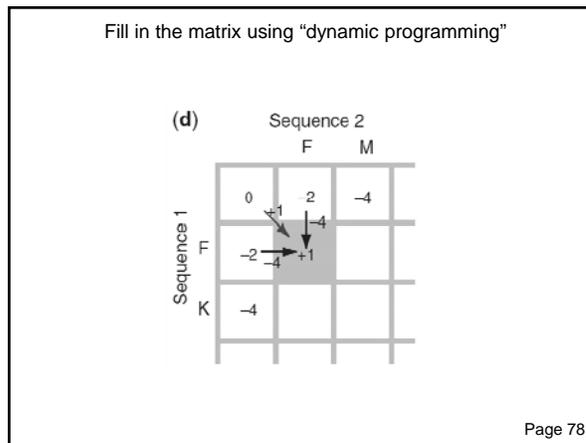
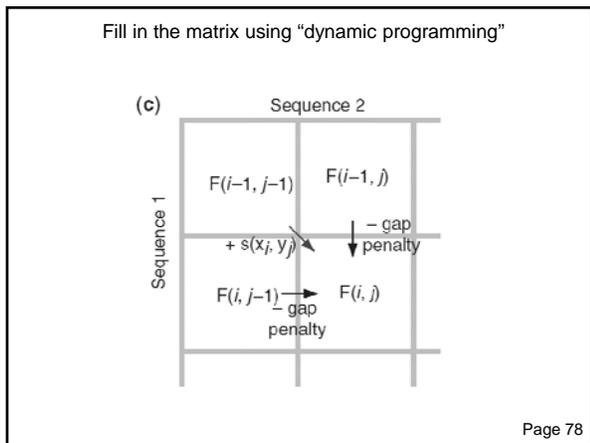
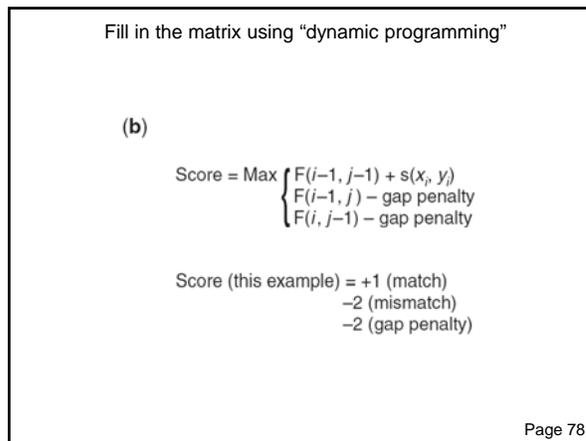
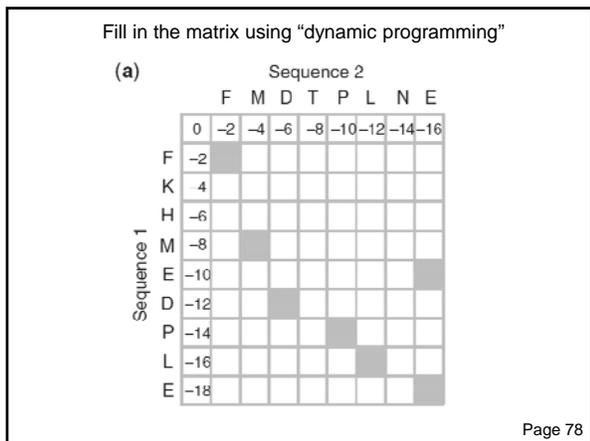
(b) Formulas:
 $F(i,j) = \max \{ F(i-1,j), F(i-1,j-1) + 1, F(i-1,j-1) - 1, F(i-1,j-1) - 2 \}$
 Score (this example) = 11 (best)
 -2 (penalty) = 9 (gap penalty)

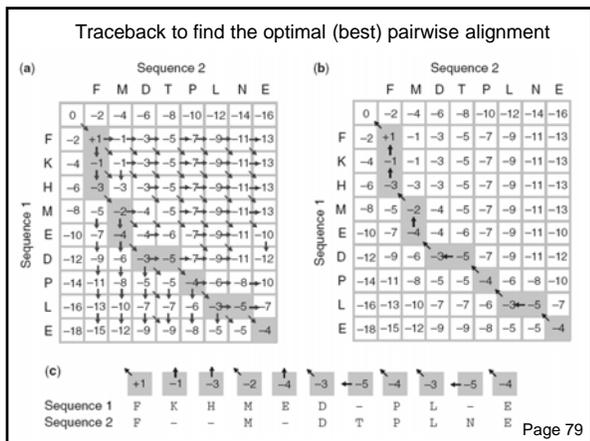
(c) Sequence 1: F M D T P L N E

(d) Sequence 2: F M D T P L N E

(e) Sequence 1: F M D T P L N E

(f) Sequence 2: F M D T P L N E





Needleman-Wunsch: dynamic programming

N-W is guaranteed to find optimal alignments, although the algorithm does not search all possible alignments.

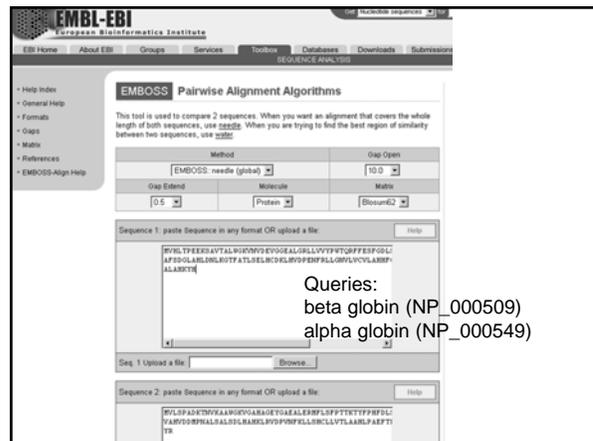
It is an example of a dynamic programming algorithm: an optimal path (alignment) is identified by incrementally extending optimal subpaths. Thus, a series of decisions is made at each step of the alignment to find the pair of residues with the best score.

Page 80

Try using needle to implement a Needleman-Wunsch global alignment algorithm to find the optimum alignment (including gaps):

<http://www.ebi.ac.uk/emboss/align/>

Page 81



```
#####
# Program: needle
# Runday: Tue Aug 22 16:29:59 2006
# Align_format: srspair
# Report_file: /ebi/entsew/old-work/needle-20060822-16295743003385_output
#####
#
# Aligned_sequences: 2
# 1: EMBOS_001
# 2: EMBOS_001
# Matrix: EMBOS62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 149
# Identity: 65/149 (43.6%)
# Similarity: 90/149 (60.4%)
# Gaps: 9/149 (6.0%)
# Score: 292.5
#
#
#####
EMBOS_001 1 NVHILTFEERSAVTALNGKY--NYDEVGGEALGELLVYVFFTORFFSFQD 48
EMBOS_001 1 NV-LSPADKTNVVAANGKGVGAGAGEYGAELERMFISPTTKTVFPHF-D 48
EMBOS_001 49 ISTPDAVWQKPKPKVAKGCVLGAFSQGLARLHWLKGTFATISLHCDKLR 98
EMBOS_001 49 IS-----RGSAGVKGKGVADALTNVAVHYDQMPNALSLSDLHAKLR 93
EMBOS_001 99 VDFENRFLIGNVLYLVAHDFKQETFFPQAAVQKTVAGFANALAHKYH 147
EMBOS_001 94 VDFVWFKLLISCLLVTLAHLPAEFTPAVRHAGLDFLASVSTVLSKYR 142
```

Global alignment versus local alignment

Global alignment (Needleman-Wunsch) extends from one end of each sequence to the other.

Local alignment finds optimally matching regions within two sequences ("subsequences").

Local alignment is almost always used for database searches such as BLAST. It is useful to find domains (or limited regions of homology) within sequences.

Smith and Waterman (1981) solved the problem of performing optimal local sequence alignment. Other methods (BLAST, FASTA) are faster but less thorough.

Page 82

Statistical significance of pairwise alignment

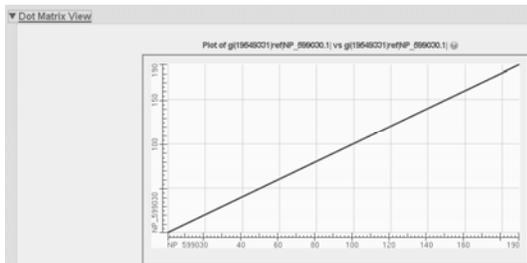
We will discuss the statistical significance of alignment scores in the next lecture (BLAST). A basic question is how to determine whether a particular alignment score is likely to have occurred by chance. According to the null hypothesis, two aligned sequences are not homologous (evolutionarily related). Can we reject the null hypothesis at a particular significance level alpha?

**Pairwise alignments with dot plots:
graphical displays of relatedness with NCBI's BLAST**

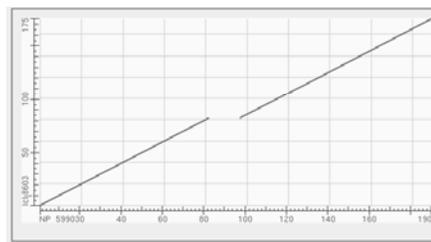
[1] Compare human cytoglobin (NP_599030, length 190 amino acids) with itself. The output includes a dot plot. The data points showing amino acid identities appear as a diagonal line.

[2] Compare cytoglobin with a globin from the snail *Biomphalaria glabrata* (accession CAJ44466, length 2,148 amino acids. See lots of repeated regions!

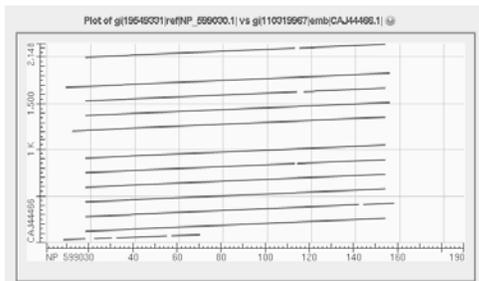
**Pairwise alignments with dot plots:
cytoglobin versus itself yields a straight line**



**Pairwise alignments with dot plots:
cytoglobin versus itself
(but with 15 amino acids deleted from one copy)**



**Pairwise alignments with dot plots:
cytoglobin versus a snail globin**



Next in the course...

Take the quiz (on pairwise alignment), due in a week (because of the Thanksgiving break, it's due TUESDAY at 5 pm).

Next Monday: BLAST