

Sequence analysis

Sequence analysis

- Next generation sequencing methods have generated large amounts of sequencing data
- **Bioinformatics** uses information from protein and DNA/RNA sequences to derive information about living organisms
- Most bioinformatic approaches are based on sequence comparisons

Sequence comparison/alignment

- You have two sequences THISSEQUENCE and THATSEQUENCE

THISSEQUENCE

| | | | | | | | | |

THATSEQUENCE

- Yet another sequence THATISASEQUENCE

TH IS SEQUENCE

| | | | | | | |

THATISASEQUENCE

or THIS SEQUENCE

| | | | | | | |

THATISASEQUENCE

Which alignment is the best?- Alignment scoring

- Simplest way is % identity
- Scoring matrices e.g. Blosom62:

(A)

C	9																			
S	-1	4																		
T	-1	1	5																	
P	-3	-1	-1	7																
A	0	1	0	-1	4															
G	-3	0	-2	-2	0	6														
N	-3	1	0	-2	-2	0	6													
D	-3	0	-1	-1	-2	-1	1	6												
E	-4	0	-1	-1	-1	-2	0	2	5											
Q	-3	0	-1	-1	-1	-2	0	0	2	5										
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8									
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5								
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5							
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5						
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4					
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4				
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	3	2	1	3	1	4			
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6		
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7	
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W

Sequence comparison/alignment

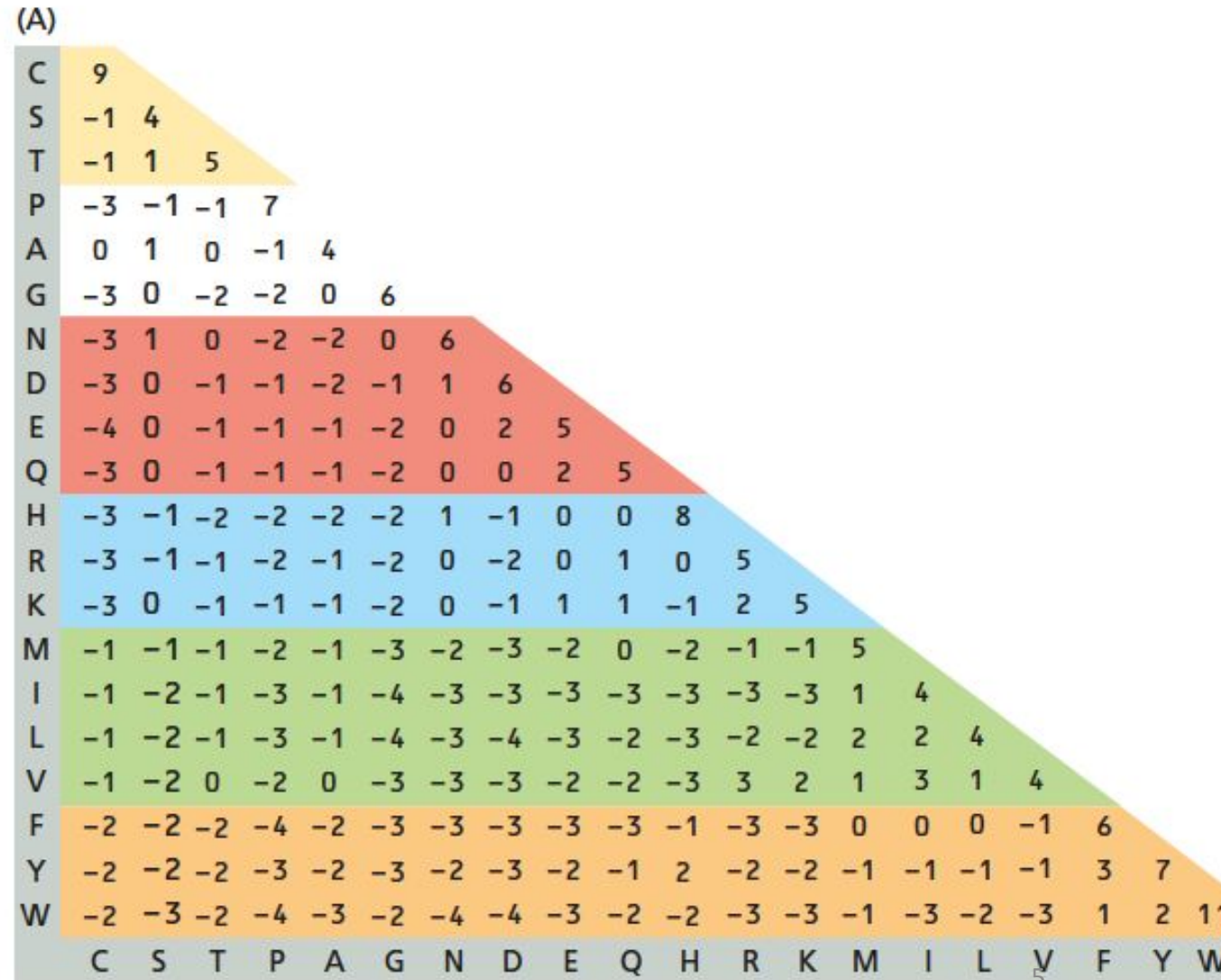
- Scoring the first two sequences

T H I S S E Q U E N C E

T H A T S E Q U E N C E

5 8-1 1 4 5 5 0 5 6 9 5

Summ is 52

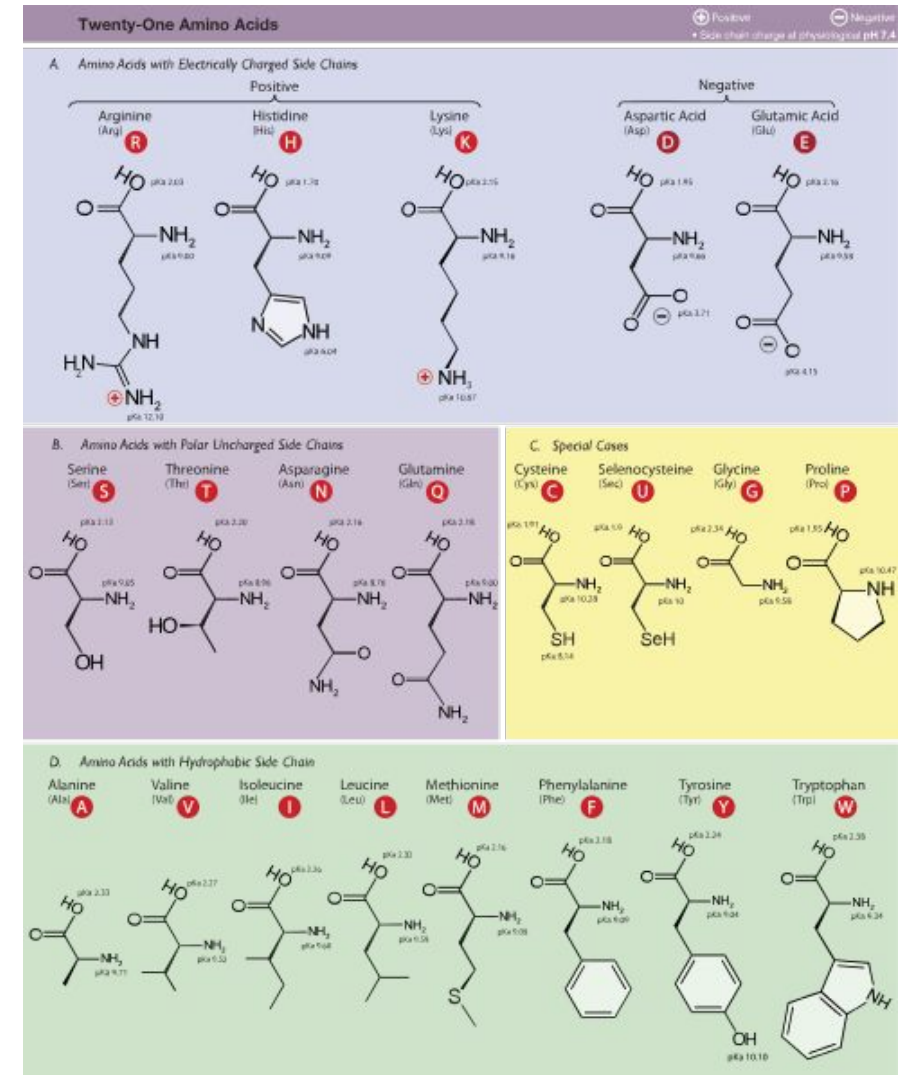


Substitution matrices

- “Various types of substitution matrices have been used over the years. Some were based on theoretical considerations, such as the ***number of mutations that are needed to convert one amino acid into another, or similarities in physicochemical properties***. The most successful, however, use ***actual evidence of what has happened during evolution***, and are based on analysis of alignments of numerous homologs of well-studied proteins from many different species.”

The genetic code – denoted as RNA triplets

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



Sequence homology – sequences derived from a common ancestor

- Sequences derived from a common ancestor during evolution share sequence similarities, how much is dependent on mutation rate and evolutionary distance
- Sequence homology can imply a common function or structure of a protein
- Sequence identity above 30% over a whole protein sequence is considered a good indication for homology

Blosum and PAM matrices

- PAM uses substitution frequencies from known closely related sequences - PAM matrix number indicates evolutionary distance thus increases with evolutionary distance and sequence divergence
- Blosum use mutation data from highly conserved local regions of sequence – Blosum matrix number refers to percentage identity thus decreases with evolutionary distance and sequence divergence

What about gaps?

- The likelihood of insertion and deletion mutations is lower than that of point mutations and often result in frame shifts thus gaps receive a penalty the standard is -11
- Insertion and deletions can have different length thus and it is likely that not just a single amino acid got deleted gap extensions thus receive a lower penalty -1

- Score the alignment of these two sequences

58

58

11

Local and global alignment

Global Alignment

<i>Target sequence</i>	1	ATCGTTGACGCACAAACACACTCTTCCAAGACCACCACATGCTGAGGTGT	50
<i>Query sequence</i>	1	ATCGTTG---ACAAACACACTCTTCCAAGAC--CCACATGCT--GGTG-	41

Local Alignment

<i>Target sequence</i>	1	ATCGTTGACGCACAAACACACTCTTCCAAGACCACCACATGCTGAGGTGT	50
<i>Query sequence</i>	1	-----TTCCAAGACCACCACATG-----	18

Basic local alignment search tool = Blast

- blastn: nucleotide blast
- tblastn: translated nucleotide blast
- blastp: protein blast
- Watch NCBI tutorial:
<https://digitalworldbiology.com/BLAST/slide1.html>