

Algorithms in Bioinformatics FOUR Pairwise Sequence Alignment

| | 0 | 1 | 2 | 3 | 4 | |
|---|---|-----|----|----|----|----|
| 0 | - | 5 | -2 | -4 | -6 | -8 |
| 1 | C | 2 | 4 | -3 | -3 | -5 |
| 2 | A | -4 | 5 | 2 | -4 | -2 |
| 3 | C | -6 | 3 | 2 | 4 | -3 |
| 4 | T | -8 | -5 | 3 | 3 | -2 |
| 5 | A | -10 | -7 | -4 | 3 | -2 |
| 6 | G | -12 | -9 | -6 | -5 | 4 |

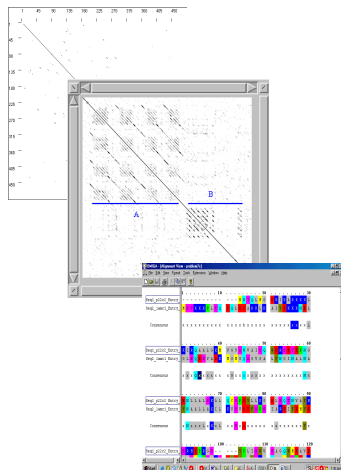
Sami Khuri
Department of Computer Science
San José State University

sami.khuri@sjsu.edu

| | 0 | 1 | 2 | 3 | 4 | |
|---|---|-----|----|----|----|----|
| 0 | - | 5 | -2 | -4 | -6 | -8 |
| 1 | C | 2 | 4 | -3 | -3 | -5 |
| 2 | A | -4 | 5 | 2 | -4 | -2 |
| 3 | C | -6 | 3 | 2 | 4 | -3 |
| 4 | T | -8 | -5 | 3 | 3 | -2 |
| 5 | A | -10 | -7 | -4 | 3 | -2 |
| 6 | G | -12 | -9 | -6 | -5 | 4 |

©2018 Sami Khuri

Pairwise Sequence Alignment



- Homology
- Similarity
- Global string alignment
- Local string alignment
- Dot matrices
- Dynamic programming
- Scoring matrices
- BLAST

©2018 Sami Khuri

Sequence Alignment

- **Sequence alignment** is the procedure of comparing sequences by searching for a series of individual characters or character patterns that are in the same order in the sequences.
 - Comparing two sequences gives us a **pairwise sequence alignment**.
 - Comparing more than two sequences gives us **multiple sequence alignment**.

©2018 Sami Khuri

Why Do We Align Sequences?

- The basic idea of aligning sequences is that **similar DNA sequences** generally produce **similar proteins**.
- To be able to predict the characteristics of a protein using only its sequence data, the **structure** or **function** information of known proteins with similar sequences can be used.
- To be able to check and see whether two (or more) genes or proteins are evolutionarily related to each other.

©2018 Sami Khuri

Importance of Alignments

- Alignment methods are at the core of many of the software tools used to search the databases.
- Alignment is the task of locating equivalent regions of two or more sequences to maximize their similarity.
- In order to assess the similarity of two sequences it is necessary to have a quantitative measure of their alignment, which includes the degree of similarity of two aligned residues as well as accounting for insertions and deletions.

Understanding Bioinformatics by Zvelebil and Baum

©2018 Sami Khuri

Query Sequence

If a query sequence is found to be significantly similar to an already annotated sequence (DNA or protein), we can use the information from the annotated sequence to possibly infer **gene structure** or **function** of the query sequence.

©2018 Sami Khuri

Similarity and Difference

- The **similarity** of two DNA sequences taken from different organisms can be explained by the theory that all contemporary genetic material has one common ancestral DNA.
- **Differences** between families of contemporary species resulted from mutations during the course of evolution.
 - Most of these changes are due to local mutations between nucleotide sequences.

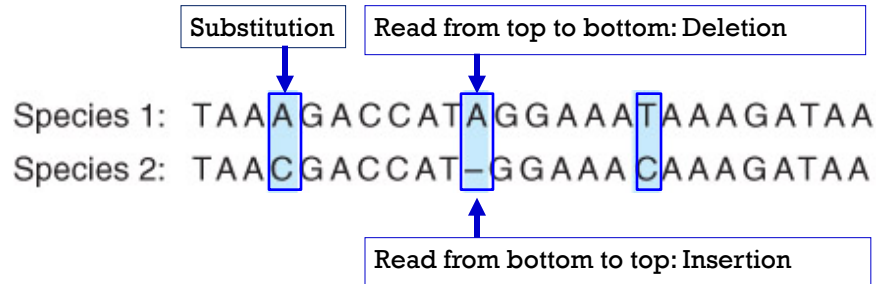
©2018 Sami Khuri

Evolution and Alignments

- Alignments reflect the **probable** evolutionary history of two sequences.
- Residues that align and that are not identical represent **substitutions**.
- Sequences without correspondence in aligned sequences are interpreted as **indels** and in an alignment are **gaps**.

©2018 Sami Khuri

How do we Compare Sequences?

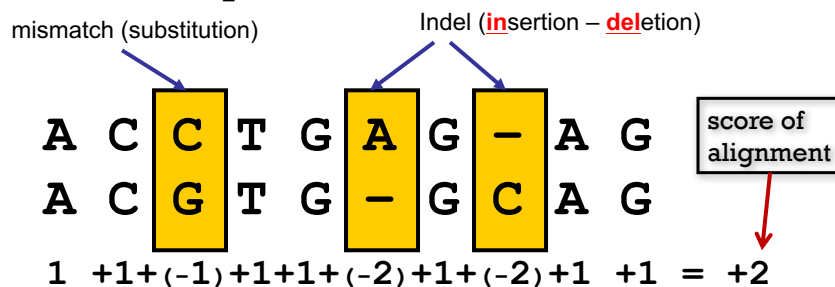


Determining the similarity of two genes by aligning their nucleotide sequences as well as possible; the differences due to mutation are shown in boxes.

©2018 Sami Khuri

Scoring a Pairwise Alignment

- The two sequences are 70% identical



- Score of the alignment where:
Match \rightarrow +1 Mismatch \rightarrow -1 Indel \rightarrow -2

©2018 Sami Khuri

Problem Definition

Given:

- Two sequences.
- A scoring system for evaluating match or mismatch of two characters.
- A penalty function for gaps in sequences.

Find:

- An **optimal pairing** of sequences that retains the order of characters in each sequence, perhaps introducing gaps, such that the total score is optimal.

©2018 Sami Khuri

Local and Global Alignments

- **Global alignment**
 - find alignment in which the **total score** is highest, perhaps at the expense of areas of great local similarity.
- **Local alignment**
 - find alignment in which the **highest scoring subsequences** are identified, at the expense of the overall score.
 - Local alignment can be obtained by performing minor modifications to the global alignment algorithm.

©2018 Sami Khuri

Local Sequence Alignment

- The **optimal local alignment** of two sequences is the one that finds the longest segment of high sequence similarity between the two sequences.

©2018 Sami Khuri

Example: Local and Global Alignments

(A) local **Local alignment**

PI3-kinase **DRHNSN**IMVKDDGGQLFHI**DFG**
cAMP PK **DLKPFN**LLIDQQGYIQVT**DFG**

(B) global **Global alignment**

PI3-kinase GQLGNLR--LEECRI--MSSAKRPLWLNWENPDIMSELLFQNNELIFKNGDDLRRQDMLT
cAMP PK GNAAAARKGGEQESVKEFLAKAKEDFLKKWENPAQNTAHLDDQFERIKTLGTGSEGRVML-

PI3-kinase LQIIRIME--NIWQNGGLDLRMLPYGCLSIGDCVGLIEVVRNSHTIMQ-IQCKKGLKGL
cAMP PK ---VKHMETGNHYAMKILDKQKVVK-----LKGIEHTLNEKRILQAVNFPFLVKLEF

PI3-kinase QFNSHT-LHQWLKDKNKGEIYDAA--IDLFTRSCAGYCVATFILGIG**DRHNSN**IMVKD-D
cAMP PK SFKDNSNLYMVMYVPGGEMFSLRRIIGRFSEPHAREFYAAQIVLTEEYLSLDLIYR**DLK**

PI3-kinase GQLEFH**IDFG**HELDHKKKKFGYKRERVP-----EVLTDQDFL---IVISKGAQECTKTREE
cAMP PK **FN**ELLIDQQGYI--QVT**DFG**FAK-RVKGRTWXLECGTPEYLAPE**ILSKG**YNKAVDWWALG

PI3-kinase RF-QEMC--YKAYLAIRQHANLFINLFSMMLGSGMPQLQSFDDIAYIRKKTALDKTEQEA
cAMP PK VLIYEMAAGYPPFA--DQPIQIYKIVSGKVR--FSSHSSDLKDLLRNLLEQVDLTKR--

PI3-kinase LEYFMKQ**MND**AHHGGWTITK**DWI**-----FHTIKQH**ALN**----
cAMP PK EGNLKNGV**NBI**KNHKWFATT**DWI**AIYQRKVEAPPFKFKGPGDTS**N**EDDYEEEEIRVXIN

Understanding Bioinformatics by Zvelebil and Baum

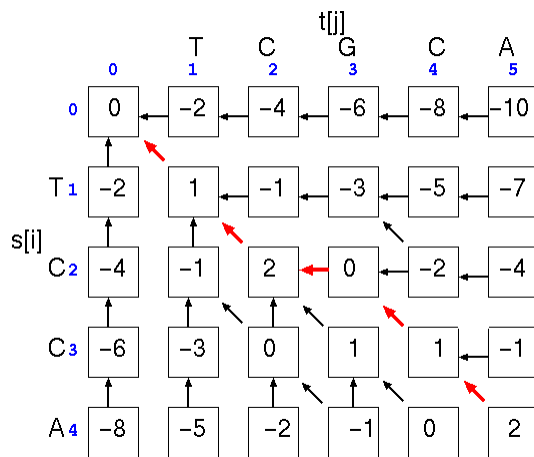
©2018 Sami Khuri

Dynamic Programming

- Dynamic programming provides a reliable and optimal computational method for aligning DNA and protein sequences.
- The **optimal alignments** provide useful information to researchers, who make **functional, structural, and evolutionary predictions** of the sequences.

©2018 Sami Khuri

Needleman Wunsch: Example



Find the optimal alignment between:
TCCA and TCGCA

Scoring Function:
+1 for match
-1 for mismatch
-2 for gap

Solution:

$$\begin{array}{cccc}
 \text{T} & \text{C} & - & \text{C} & \text{A} \\
 : & : & & : & : \\
 \text{T} & \text{C} & \text{G} & \text{C} & \text{A} \\
 \hline
 1 & +1 & -2 & +1 & +1 & = & 2
 \end{array}$$

©2018 Sami Khuri

Local Sequence Alignment

- Their dynamic algorithm gives a **global alignment** of sequences.
- A modification of the dynamic programming algorithm for sequence alignment provides a **local sequence alignment** giving the highest-scoring local match between two sequences (Smith and Waterman 1981).
- **Local alignments** are usually more meaningful than global matches because they include patterns that are conserved in the sequences.

©2018 Sami Khuri

Scoring Systems

- Use of the **dynamic programming** method requires a scoring system for
 - the comparison of symbol pairs (**nucleotides** for DNA sequences & **amino acids** for protein sequences),
 - a scheme for insertion/deletion (gap) penalties.
- The most commonly used scoring systems for protein sequence alignments are the log odds form
 - of the **PAM250** matrix and
 - the **BLOSUM62** matrix.
- A number of other choices are available.

©2018 Sami Khuri

Scoring Matrices

- The alignment algorithm needs to know if it is more likely that a given amino acid pair has occurred **randomly** or that it has occurred as a result of an **evolutionary** event.
- Similar amino acids are defined by high-scoring matches between the amino acid pairs in the substitution matrix.

©2018 Sami Khuri

Amino Acid Substitution Matrices (I)

- For proteins, an **amino acid substitution matrix**, such as the Dayhoff **P**ercent **A**ccepted **M**utation matrix 250 (**PAM250**) or **B**lock **S**ubstitution **M**atrix 62 (**BLOSUM62**) is used to score matches and mismatches.
- Similar matrices are available for aligning DNA sequences.

©2018 Sami Khuri

Amino Acid Substitution Matrices (II)

- In the **amino acid substitution matrices**, amino acids are listed both across the top of a matrix and down the side, and each matrix position is filled with a score that reflects how often one amino acid would have been paired with the other in an alignment of related protein sequences.

©2018 Sami Khuri

PAM Matrices

Point Accepted Mutation

- An **accepted mutation** is any mutation that doesn't kill the protein or organism; that is, amino acid changes “accepted” by natural selection.

One PAM (PAM1) = 1% of the amino acids have been changed.

©2018 Sami Khuri

Constructing More PAM Matrices

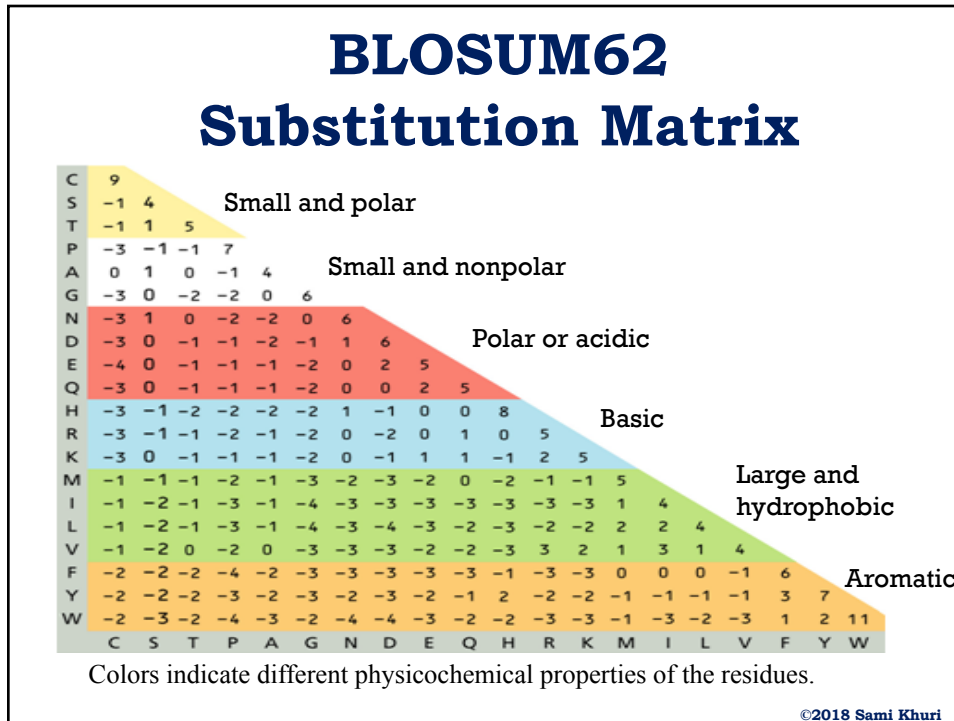
- The **PAM1** Matrix is best used for comparing sequences where 1% or less of the amino acids have changed.
- What do you do with sequences that are more divergent?
- You multiply the PAM1 matrix by itself N times to get a new matrix that works best with sequences that have PAM2, PAM20, PAM100, PAM200, etc.
- For example $PAM20 = (PAM1)^{20}$

©2018 Sami Khuri

BLOSUM

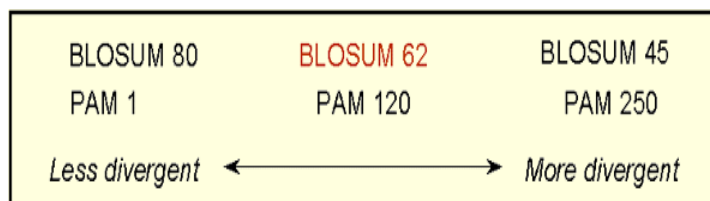
- **Block Substitution Matrix**
 - created from BLOCKS database.
- Currently the most widely used comparison matrix.
- More sensitive than PAM or other matrices
- Finds more sequences that are related
- The BLOSUM matrices are based on an entirely different type of sequence analysis and a much larger data set than the Dayhoff PAM Matrices.

©2018 Sami Khuri



Comparison: PAM and BLOSUM Matrices

The **PAM** model is designed to track the evolutionary origins of proteins, whereas the **BLOSUM** model is designed to find their conserved domains.



©2018 Sami Khuri

Approximate Methods

BLAST

- **B**asic **L**ocal **A**lignment **S**earch **T**ool
– Altschul et al. 1990, 1994, 1997
- Heuristic method for local alignment
- Designed specifically for database searches
- Idea: Good alignments contain short lengths of exact matches.

©2018 Sami Khuri

The BLAST Family

- **blastp**: compares an amino acid query sequence against a protein sequence database.
- **blastn**: compares a nucleotide query sequence against a nucleotide sequence database.
- **blastx**: compares the six-frame conceptual translation products of a nucleotide query sequence (both strands) against a protein sequence database.

©2018 Sami Khuri

Chance or Homology?

- In all methods of sequence comparison, the fundamental question is whether the similarities perceived between two sequences are due to chance, and are thus of little biological significance, or whether they are due to the derivation of the sequences from a common ancestral sequence, and are thus homologous.

Understanding Bioinformatics by Zvelebil and Baum

©2018 Sami Khuri

The Expected Value

| | | | |
|----------------------|------------------------------|-----------------|----------|
| SW:P11A BOVIN P32871 | PHOSPHATIDYLINOSITOL 3-KINAS | (1068) 2228 493 | 1.2e-138 |
| SW:P11A HUMAN P42336 | PHOSPHATIDYLINOSITOL 3-KINAS | (1068) 2216 490 | 7.4e-138 |
| SW:P11A MOUSE P42337 | PHOSPHATIDYLINOSITOL 3-KINAS | (1068) 2204 488 | 4.5e-137 |
| SW:P11B HUMAN P42338 | PHOSPHATIDYLINOSITOL 3-KINAS | (1070) 1126 254 | 1.1e-66 |

The **e-value** tells us how likely it is that the similarity between the query sequence and the database sequence is due to chance.

The lower the **e-value**, the more likely it is that the two sequences are truly similar and not just chance matches

©2018 Sami Khuri