

WEBLEM 10

Introduction to Pairwise Alignment

Pairwise Alignment:

Sequence comparison lies at the heart of bioinformatics analysis. It is an important first step toward structural and functional analysis of newly determined sequences. As new biological sequences are being generated at exponential rates, sequence comparison is becoming increasingly important to draw functional and evolutionary inference of a new protein with proteins already existing in the database. The most fundamental process in this type of comparison is sequence alignment. This is the process by which sequences are compared by searching for common character patterns and establishing residue-residue correspondence among related sequences. Pairwise sequence alignment is the process of aligning two sequences and is the basis of database similarity searching and multiple sequence alignment. This chapter introduces the basics of pairwise alignment.

Sequence Similarity vs Sequence Identity:

Another set of related terms for sequence comparison are sequence similarity and sequence identity. Sequence similarity and sequence identity are synonymous for nucleotide sequences. For protein sequences, however, the two concepts are very different. In a protein sequence alignment, sequence identity refers to the percentage of matches of the same amino acid residues between two aligned sequences. Similarity refers to the percentage of aligned residues that have similar physicochemical characteristics and can be more readily substituted for each other.

Needleman-Wunsch Algorithm:

Needleman-Wunsch Algorithm developed by Saul B. Needleman and Christian D. Wunsch in 1970. It was designed to compare biological sequences and was one of the first applications of dynamic programming to the biological sequence comparison. This algorithm is usually used for global alignment of two sequences (nucleotide or amino acids). For the purpose of explanation, first summarizing the algorithm in five steps and then I will move forwards with examples and explanations.

- 1) Consider all the possible pairs of residues from two sequences, the best way is to generate a 2D matrix of two sequences. We will need 2 such matrices one for sequence and one for scores.
- 2) Initialize the score matrix. Scoring matrices are used to determine the relative score made by matching two characters in a sequence alignment. There are many flavors of scoring matrices for amino acid sequences, nucleotide sequences, and codon sequences, and each is derived from the alignment of "known" homologous sequences. These alignments are then used to determine the likelihood of one character being at the same position in the sequence as another character.
- 3) Gap penalty: Usually there are high chances of insertions and deletions (indels) in biological sequences but one large indel is more likely rather than multiple small indels in a given sequences. In order to tackle this issue we give two kind of penalties; Gap opening penalty (relatively higher) and gap extension penalty (relatively lower)
- 4) Calculate scores and fill the traceback matrix
- 5) Deduce the alignment from the traceback matrix

Smith-Waterman alignment algorithm:

Over a decade after the initial publication of the Needleman-Wunsch algorithm, a modification was made to allow for local alignments (Smith and Waterman, 1981). Today, the Smith-Waterman alignment algorithm is the one used by the Basic Local Alignment Search Tool (BLAST) which is the most cited resource in biomedical literature. In this adaptation, the alignment path does not need to reach the edges of the search graph, but may begin and end internally. In order to accomplish this, 0 was added as a term in the score calculation described by Needleman and Wunsch.

Recall that for global alignments the value at any point is:

$$M(i,j) = \text{MAX}(M(i-1,j-1) + S(A_i, B_j)$$

$$M(i-1,j) + \text{gap}, M(i,j-1) + \text{gap})$$

However for local alignments the score calculation becomes:

$$M(i,j) = \text{MAX}(M(i-1,j-1) + S(A_i, B_j)$$

$$M(i-1,j) + \text{gap}, M(i,j-1) + \text{gap}, 0)$$

The implication of this is that there are no values below zero in a local alignment scoring matrix.

WEBLEM 10/A

(URL: <https://www.ebi.ac.uk/Tools/psa/>)

Aim:

To Study the query “Actin Sequence” in Needleman-Wunsch algorithm

Introduction:

Needleman-Wunsch Algorithm developed by Saul B. Needleman and Christian D. Wunsch in 1970. It was designed to compare biological sequences and was one of the first applications of dynamic programming to the biological sequence comparison. This algorithm is usually used for global alignment of two sequences (nucleotide or amino acids). For the purpose of explanation, first summarizing the algorithm in five steps and then I will move forwards with examples and explanations.

Actin is the most abundant protein in most eukaryotic cells. It is highly conserved and participates in more protein-protein interactions than any known protein. These properties, along with its ability to transition between monomeric (G-actin) and filamentous (F-actin) states under the control of nucleotide hydrolysis, ions, and a large number of actin-binding proteins, make actin a critical player in many cellular functions, ranging from cell motility and the maintenance of cell shape and polarity to the regulation of transcription. Moreover, the interaction of filamentous actin with myosin forms the basis of muscle contraction.

Methodology:

1. Open the homepage of Pair-Wise Alignment.
2. After that open the Needle-Wusch algorithm.
3. Enter the protein sequence.
4. Open the result.
5. Interpret the result.

Observation:

Tools > Pairwise Sequence Alignment

Pairwise Sequence Alignment is used to identify regions of similarity that may indicate functional, structural and/or evolutionary relationships between two biological sequences (protein or nucleic acid).

By contrast, **Multiple Sequence Alignment (MSA)** is the alignment of three or more biological sequences of similar length. From the output of MSA applications, homology can be inferred and the evolutionary relationship between the sequences studied.

Global Alignment

Global alignment tools create an end-to-end alignment of the sequences to be aligned.

Needle (EMBOSS)
EMBOSS Needle creates an optimal global alignment of two sequences using the Needleman-Wunsch algorithm.
Launch [Needle](#)

Stretcher (EMBOSS)
EMBOSS Stretcher uses a modification of the Needleman-Wunsch algorithm that allows larger sequences to be globally aligned.
Launch [Stretcher](#)

GSEARCH2SEQ
GSEARCH2SEQ finds an optimal global alignment using the Needleman-Wunsch algorithm.
Launch [ggsearch2seq](#)

Local Alignment

Fig1: Homepage of pairwise Sequence Alignment

```

>sp|P07830|ACT1_DICDI Major actin OS=Dictyostelium discoideum OX=44689 GN=act1 PE=1 SV=2
MDGEDVQALVIDNGSGMCKAGFAGDDAPRAVFPISVGRPRHTGVMVGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIVTNWDDMEKIWHHTFYNELRVAPEEHPVLLTEAPLNPKANREKM
TQIMFETFNTPAMYVAIQAVLSLYASGRRTTGIVMDSGDGVSHTVPIYEGYALPHAILRLD
LAGRDLTDYMMKILTERGYSFTTTAEREIVRDIKEKLAYVALDFAEMQTAASSSALEKS
YELPDGQVITIGNERFRCPEALFQPSFLGMESAGIHETTYNSIMKCDVDIRKDLYGNVVL
SGGTTMFPGIADRMNKELTALAPSTMKIKIIPPERKYSVWIGGSILASLSTFQQMWISK
EEYDESGPSIVHRKCF

>sp|P63268|ACTH_MOUSE Actin, gamma-enteric smooth muscle OS=Mus musculus OX=10090 GN=Actg2 PE=1 SV=1
MCEEETALVCDNGSGLCKAGFAGDDAPRAVFPISVGRPRHQGVMMGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIITNWDDMEKIWHHSFYNELRVAPEEHTLLTEAPLNPKANREKM
TQIMFETFNVPAMYVAIQAVLSLYASGRRTTGIVLDSGDGVTNNVPIYEGYALPHAIMRLD
LAGRDLTDYLMKILTERGYSFVTTAEREIVRDIKEKLCYVALDFENEMATAASSSSLEKS
YELPDGQVITIGNERFRCPETLFQPSFIGMESAGIHETTYNSIMKCDIDIRKDLYANNVL
SGGTTMYPGIADRMQKEITALAPSTMKIKIIPPERKYSVWIGGSILASLSTFQQMWISK
PEYDEAGPSIVHRKCF

```

Fig2: Two FASTA sequence for actin from uniprot

Pairwise Sequence Alignment

EMBOSS Needle reads two input sequences and writes their optimal global sequence alignment to file.

STEP 1 - Enter your protein sequences

Enter a pair of

PROTEIN

sequences. Enter or paste your first **protein** sequence in any supported format:

```

>sp|P07830|ACT1_DICDI Major actin OS=Dictyostelium discoideum OX=44689 GN=act1 PE=1 SV=2
MDGEDVQALVIDNGSGMCKAGFAGDDAPRAVFPISVGRPRHTGVMVGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIVTNWDDMEKIWHHTFYNELRVAPEEHPVLLTEAPLNPKANREKM
TQIMFETFNTPAMYVAIQAVLSLYASGRRTTGIVMDSGDGVSHTVPIYEGYALPHAILRLD
LAGRDLTDYMMKILTERGYSFTTTAEREIVRDIKEKLAYVALDFAEMQTAASSSALEKS
YELPDGQVITIGNERFRCPEALFQPSFLGMESAGIHETTYNSIMKCDVDIRKDLYGNVVL
SGGTTMFPGIADRMNKELTALAPSTMKIKIIPPERKYSVWIGGSILASLSTFQQMWISK
EEYDESGPSIVHRKCF

```

Or, upload a file: No file chosen

Use a example sequence | | [See more example inputs](#)

AND

Enter or paste your second **protein** sequence in any supported format:

```

>sp|P63268|ACTH_MOUSE Actin, gamma-enteric smooth muscle OS=Mus musculus OX=10090 GN=Actg2 PE=1 SV=1
MCEEETALVCDNGSGLCKAGFAGDDAPRAVFPISVGRPRHQGVMMGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIITNWDDMEKIWHHSFYNELRVAPEEHTLLTEAPLNPKANREKM
TQIMFETFNVPAMYVAIQAVLSLYASGRRTTGIVLDSGDGVTNNVPIYEGYALPHAIMRLD
LAGRDLTDYLMKILTERGYSFVTTAEREIVRDIKEKLCYVALDFENEMATAASSSSLEKS
YELPDGQVITIGNERFRCPETLFQPSFIGMESAGIHETTYNSIMKCDIDIRKDLYANNVL
SGGTTMYPGIADRMQKEITALAPSTMKIKIIPPERKYSVWIGGSILASLSTFQQMWISK
PEYDEAGPSIVHRKCF

```

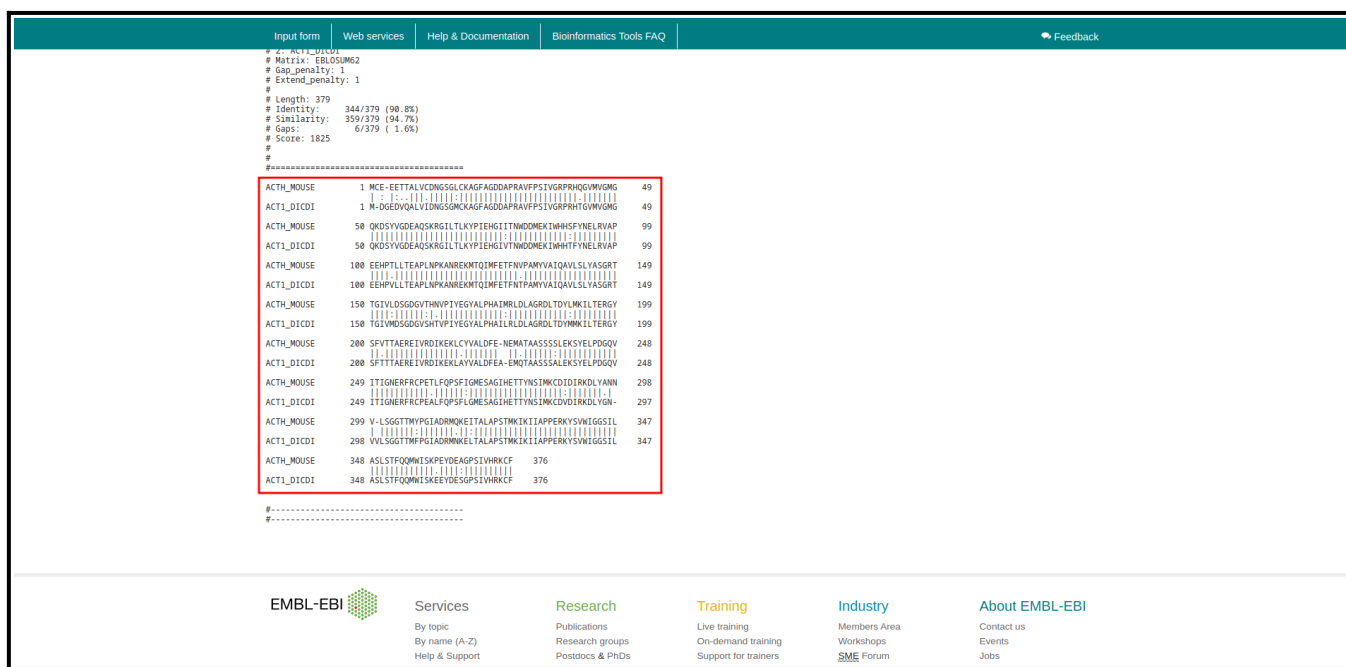
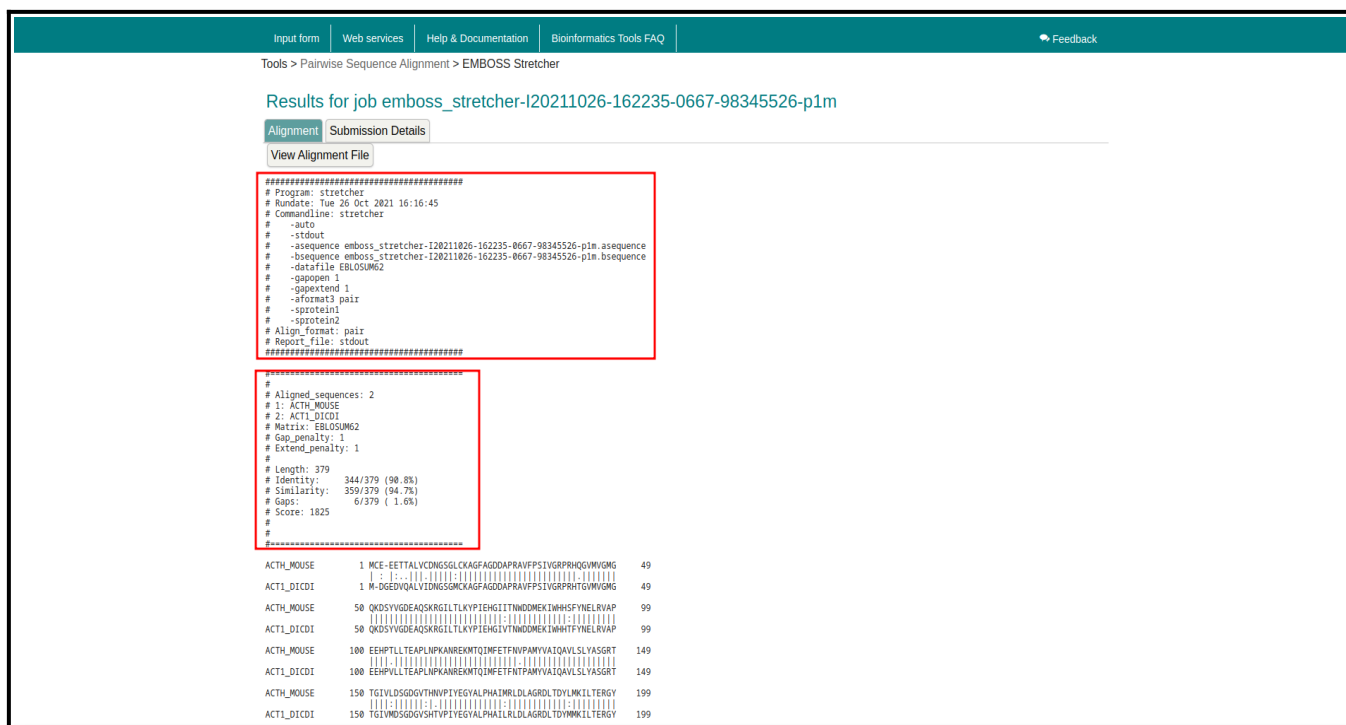
Or, upload a file: No file chosen

STEP 2 - Set your pairwise alignment options

OUTPUT FORMAT

pair

Fig3: Needleman-Wunsch algorithm search bar with two actin sequences



Conclusion:

Needleman-Wunsch algorithm in bioinformatics to align protein or nucleotide sequences. It was one of the first applications of dynamic programming to compare biological sequences. The Needleman-Wunsch algorithm is still widely used for optimal global alignment, particularly when the quality of the global alignment is of the utmost importance. The algorithm assigns a score to every possible alignment, and the purpose of the algorithm is to find all possible alignment having the highest score.

References:

1. Xiong, J. (2006). Essential Bioinformatics (1st ed.). Cambridge University Press.
2. (n.d.). Retrieved from <https://www.uniprot.org/uniprot/P63268.fasta>
3. (n.d.). Retrieved from <https://www.uniprot.org/uniprot/P07830.fasta>
4. Dominguez, R., & Holmes, K. C. (2011). Actin structure and function. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3130349/>
5. Embl-Ebi. (n.d.). Pairwise Sequence Alignment. Retrieved from <https://www.ebi.ac.uk/Tools/psa/>
6. Embl-Ebi. (n.d.). EMBOSS Needle. Retrieved from https://www.ebi.ac.uk/Tools/psa/emboss_needle/
7. Embl-Ebi. (n.d.). EMBOSS Stretcher. Retrieved from https://www.ebi.ac.uk/Tools/services/web/toolresult.ebi?jobId=emboss_stretcher-I20211026-162235-0667-98345526-p1m
8. Embl-Ebi. (n.d.). EMBOSS Stretcher. Retrieved from https://www.ebi.ac.uk/Tools/services/web/toolresult.ebi?jobId=emboss_stretcher-I20211026-162235-0667-98345526-p1m

WEBLEM 10/B

(URL: <https://www.ebi.ac.uk/Tools/psa/>)

Aim:

To study the query “Actin” in Smith-Waterman alignment algorithm

Introduction:

Over a decade after the initial publication of the Needleman-Wunsch algorithm, a modification was made to allow for local alignments (Smith and Waterman, 1981). Today, the Smith-Waterman alignment algorithm is the one used by the Basic Local Alignment Search Tool (BLAST) which is the most cited resource in biomedical literature. In this adaptation, the alignment path does not need to reach the edges of the search graph, but may begin and end internally. In order to accomplish this, 0 was added as a term in the score calculation described by Needleman and Wunsch.

Actin is the most abundant protein in most eukaryotic cells. It is highly conserved and participates in more protein-protein interactions than any known protein. These properties, along with its ability to transition between monomeric (G-actin) and filamentous (F-actin) states under the control of nucleotide hydrolysis, ions, and a large number of actin-binding proteins, make actin a critical player in many cellular functions, ranging from cell motility and the maintenance of cell shape and polarity to the regulation of transcription. Moreover, the interaction of filamentous actin with myosin forms the basis of muscle contraction.

Methodology:

1. Open the homepage of Pair-Wise Alignment.
2. After that open the Smith-Waterman alignment algorithm.
3. Enter the protein sequence.
4. Open the result.
5. Interpret the result.

Observation:

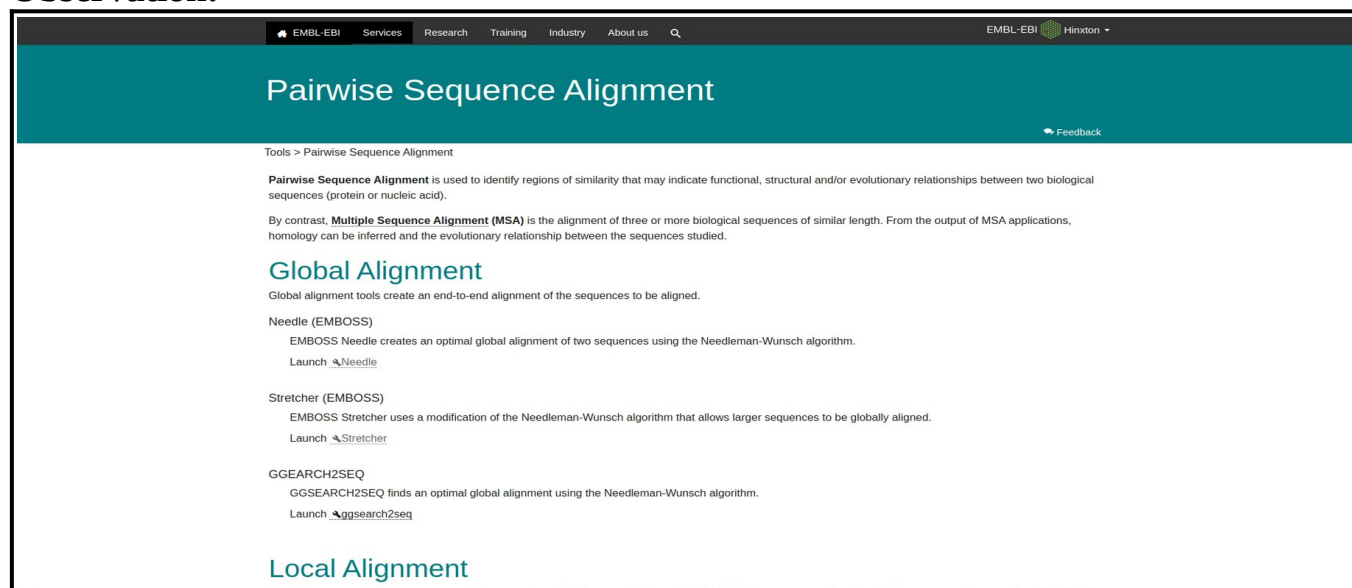


Fig1. Homepage of pairwise Sequence alignment

```

>sp|P07830|ACT1_DICDI Major actin OS=Dictyostelium discoideum OX=44689 GN=act1 PE=1 SV=2
MDGEDVQALVIDNGSGMCKAGFAGDDAPRAVFPSIVGRPRHTGVMVGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIVTNWDDMEKIWHHTFYNELRVAPEEHPVLLTEAPLNPKANREKM
TQIMFETFNTPAMYVAIQAVLSLYASGRRTGIVMDSGDGVSHTVPIYEGYALPHAILRLD
LAGRDLTDYMMKILTERGYSFTTTAEREIVRDIKEKLAYVALDFAEMQTAASSSALEKS
YELPDGQVITIGNERFRCPEALFQPSFLGMESAGIHETTYNSIMKCDVDIRKDLGNVVL
SGGTTMFGIADRMNKELTALAPSTMKIKIAPPKYSVWIGGSILASLSTFQQMWISK
EEYDESGPSIVHRKCF

>sp|P63268|ACTH_MOUSE Actin, gamma-enteric smooth muscle OS=Mus musculus OX=10090 GN=Actg2 PE=1 SV=1
MCEEETALVCDNGSGLCKAGFAGDDAPRAVFPSIVGRPRHQGVMMGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIITNWDDMEKIWHHSFYNELRVAPEEHPVLLTEAPLNPKANREKM
TQIMFETFNTPAMYVAIQAVLSLYASGRRTGIVLDSGDGVTNNVPIYEGYALPHAIMRLD
LAGRDLTDYLMKILTERGYSFVTTAEREIVRDIKEKLCYVALDFENEMATAASSSSLEKS
YELPDGQVITIGNERFRCPETLFQPSFIGMESAGIHETTYNSIMKCDIDIRKDLANNVL
SGGTTMYPGIADRMQKEITALAPSTMKIKIAPPKYSVWIGGSILASLSTFQQMWISK
PEYDEAGPSIVHRKCF

```

Fig.2. Two FASTA sequence for actin from uniprot

EMBOSS Water

[Input form](#)
[Web services](#)
[Help & Documentation](#)
[Bioinformatics Tools FAQ](#)

Feedback

Tools > Pairwise Sequence Alignment > EMBOSS Water

Pairwise Sequence Alignment

EMBOSS Water uses the Smith-Waterman algorithm (modified for speed enhancements) to calculate the local alignment of two sequences.

STEP 1 - Enter your protein sequences

Enter a pair of

PROTEIN

sequences. Enter or paste your first **protein** sequence in any supported format:

```

>sp|P07830|ACT1_DICDI Major actin OS=Dictyostelium discoideum OX=44689 GN=act1 PE=1 SV=2
MDGEDVQALVIDNGSGMCKAGFAGDDAPRAVFPSIVGRPRHTGVMVGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIVTNWDDMEKIWHHTFYNELRVAPEEHPVLLTEAPLNPKANREKM
TQIMFETFNTPAMYVAIQAVLSLYASGRRTGIVMDSGDGVSHTVPIYEGYALPHAILRLD
LAGRDLTDYMMKILTERGYSFTTTAEREIVRDIKEKLAYVALDFAEMQTAASSSALEKS
YELPDGQVITIGNERFRCPEALFQPSFLGMESAGIHETTYNSIMKCDVDIRKDLGNVVL
SGGTTMFGIADRMNKELTALAPSTMKIKIAPPKYSVWIGGSILASLSTFQQMWISK
EEYDESGPSIVHRKCF

```

Or, upload a file:

Choose File

No file chosen

Use a example sequence

Clear sequence

See more example inputs

AND

Enter or paste your second **protein** sequence in any supported format:

```

>sp|P63268|ACTH_MOUSE Actin, gamma-enteric smooth muscle OS=Mus musculus OX=10090 GN=Actg2 PE=1 SV=1
MCEEETALVCDNGSGLCKAGFAGDDAPRAVFPSIVGRPRHQGVMMGMGQKDSYVGDEAQ
SKRGILTLKYPIEHGIITNWDDMEKIWHHSFYNELRVAPEEHPVLLTEAPLNPKANREKM
TQIMFETFNTPAMYVAIQAVLSLYASGRRTGIVLDSGDGVTNNVPIYEGYALPHAIMRLD
LAGRDLTDYLMKILTERGYSFVTTAEREIVRDIKEKLCYVALDFENEMATAASSSSLEKS
YELPDGQVITIGNERFRCPETLFQPSFIGMESAGIHETTYNSIMKCDIDIRKDLANNVL
SGGTTMYPGIADRMQKEITALAPSTMKIKIAPPKYSVWIGGSILASLSTFQQMWISK
PEYDEAGPSIVHRKCF

```

Or, upload a file:

Choose File

No file chosen

Fig.3 Smith-Waterman alignment algorithm search bar with 2 actin sequences

Conclusion:

The alignment tab shows the alignment of the two sequences, with all the described parameters, used scoring matrices and Gap penalty scored values. The alignment tab has an option for the user to download the entire alignment. The alignment between the two sequences is shown in Figure 4, the gaps are represented with '-'. If a match is there between the two nucleotides there is a symbol '|' and the mismatch is represented with a dot '.'.

References:

1. Xiong, J. (2006). Essential Bioinformatics (1st ed.). Cambridge University Press.
2. (n.d.). Retrieved from <https://www.uniprot.org/uniprot/P63268.fasta>
3. (n.d.). Retrieved from <https://www.uniprot.org/uniprot/P07830.fasta>
4. Dominguez, R., & Holmes, K. C. (2011). Actin structure and function. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3130349/>
5. Embl-Ebi. (n.d.). Pairwise Sequence Alignment. Retrieved from <https://www.ebi.ac.uk/Tools/psa/>
6. Embl-Ebi. (n.d.). EMBOSS Needle. Retrieved from https://www.ebi.ac.uk/Tools/psa/emboss_needle/
7. Embl-Ebi. (n.d.). EMBOSS Stretcher. Retrieved from https://www.ebi.ac.uk/Tools/services/web/toolresult.ebi?jobId=emboss_stretcher-I20211026-162235-0667-98345526-p1m
8. Embl-Ebi. (n.d.). EMBOSS Stretcher. Retrieved from https://www.ebi.ac.uk/Tools/services/web/toolresult.ebi?jobId=emboss_stretcher-I20211026-162235-0667-98345526-p1m