

Sequence alignment for similarity in biological data.

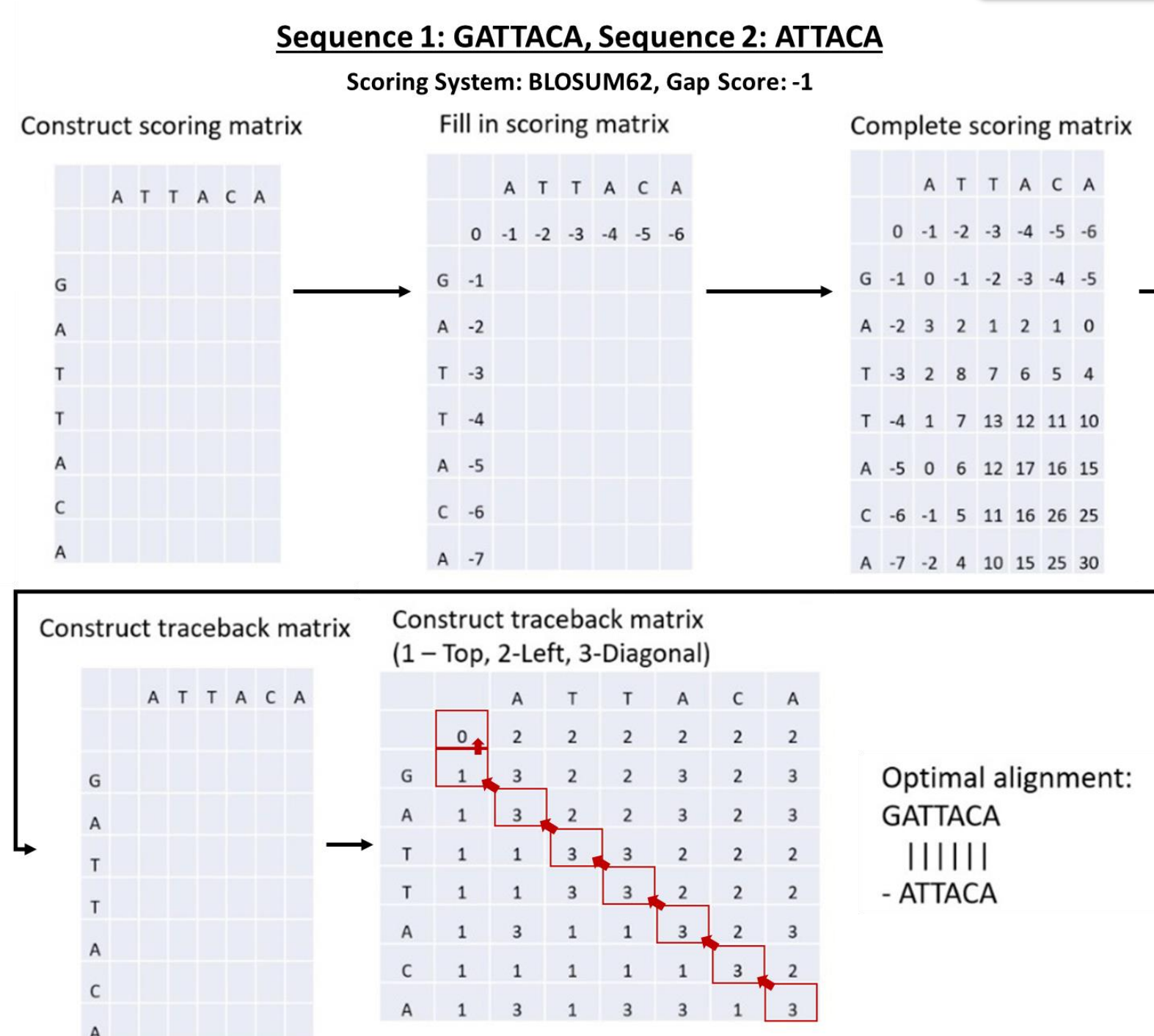
Background

The basic concept involves pattern recognition or pattern matching for DNA, RNA and protein sequences. To put things into perspective, when a new gene is discovered, the biologists use information available from previously known genes to find some characteristics of the new gene. There are different algorithms which can be utilized for sequence alignment and our team have chosen the Needleman- Wunsch algorithm.

Abstract

Every living organism's cells have their own Deoxyribonucleic Acid (DNA) that carries genetic information. Using algorithms for pattern recognition and matching, biologists can learn more about new DNA, RNA and Protein sequences by comparing them to existing known sequences. An optimal algorithm will be used to provide the users with analysed information of the sequences, so as to aid in learning more about a new sequence and its characteristics.

Our Algorithm



Needleman-Wunsch algorithm steps.

Time Complexity: $O(mn)$.

Our Application

The Needleman–Wunsch algorithm is an algorithm developed by Saul B. Needleman and Christian D. Wunsch and published in 1970, commonly used in bioinformatics to globally align protein or nucleotide sequences with the highest quality of alignment. The algorithm uses dynamic programming to compare the sequences, essentially divides a large problem (e.g. the full sequence) into a series of smaller problems (one pair of amino acid) and uses the solutions to the smaller problems to reconstruct a solution to the larger problem.

Assessing the relationships between two sequences can be done by counting the number of identical and similar amino acids. The number of identical and similar amino acids may then be compared to the total number of amino acids in the protein, giving the percentage of sequence identity and similarity.

Algorithm Walkthrough

1. Build the first matrix for containing the scores.
2. For each pair of sequence, lookup the BLOSUM62 matrix to get the match/mismatch score.
3. Build the traceback matrix to get the optimal alignment
4. Based on the scoring allocated, decide if an indel (gap) or a match/mismatch would be optimal.
5. Build the optimal alignment based on the traceback matrix.

Sequence Alignment for Similarity in Biological Data

PAIRWISE SEQUENCE ALIGNMENT
 Scoring using BLOSUM62 with linear gap score

Enter 1st protein sequence:
 MDQREILQKFLDEAQSCKITKEEFANE-FLKLRQ-S-TK-YKADKTYPT
 TTLLDFWRMIWEYSVLIIIVMACMEYEMGKKKCYWAEPMQLEFGFFSVSCAEKRRKSDYIIRTLKVKFNSERTRTIYQFHYKNWPDHVPSSIDPILFIWDVRC
 YQEDDSVPICIHCSAGCGRTGVICAIDYTWMILKDGIIIPENFSVFLIREMRTQRPISLVOTQEYELVYNVAVLELFRKQMDVIRDKHSGTESQAKHCIPKRNHTLQA
 DSYSNLPKSTTKGAQRMMNQORTKMEIKESSEDFRTSEISAKEELVLHPAKSSTSFDFLELNYSDKQADTMRKQTKAFFIVGEPLQKHQSLDLGSLFEGCSNS
 KPVNAAGRYFNSKVPITRTKSTPFELIQORETKEVDKSNFYSLEQPHDSCFVEMQAQKVMHVSSEALNYSLEPYDSKHQIRNASNVKHHDSALGVYSIPLVENP
 YFSSWPPSGTSSKMSLDLPEKQDGVTFPSSLLPTSTSLFSYNSHDSLSLNSPTNISLLNQESAVLATAPRIDDEIPPLPVRTPESTFVVEEAGEFSPNVKSL
 SSAVKVYKIGTSLEWGTSEPKKFDSDVILRPSKSVKLRSPKSELHQDRSSPPPLPERTLESFFLADEDCMQAQSIIETYSTSYPTDMENSTSSKQTLTKPGKSFTR
 KSLKILRNMRKSCICNSCPNPKFAESVQSNSSSFLNFGFANRFSKPKGRNPPPTWNI

or Upload a text file [Browse](#)

Enter 2nd protein sequence:
 MPIGSKERTFFEIFKAEKDCSTTRCNKADLGPISLWFEELSSEAPPYNSEPAEESSEKNNNYEKNLFTPQRKPSYNQLASTPIIFKEQGLTLPLYQSPVKELDK
 FKLLDGRNVNSRHKSLRTVTKRMDQDDVSCPLLNSCLSESPVVLQCTHTVTPQRDKSVVCGSLFHTPKFVKGRQTPKHISESLGAEVDPDMWSSSLATPTLSST
 VLIVRNEASETVFPHTDTANVKSYSFNHDESLKNDRFIASVDSNTNQREAAASHGFGKTSNGSKFVNSCKDHIGKSMNPVLEDEVYTVVDTSEEDSFLCSFK
 CRTQNLQKVRTSKTRKKIFHEANADECEKSKNQKYSFVSEVEPNDTDPDLSNVANQKPFESGSKISKEVVPVSLACWSQLTSLGNGAQMEKIPLLHSSCDQ
 NISEKDLDTENKRRKDFLTSENSLPRISSLPKSEKPLNEETVNNKRDEEQHLESHTDCILAVKQAISETSPVASSFQGIKKSIFRIRSPKETFNASFSGHMTDPN
 FKETEASESGLEIHTVCSQKEDSLCPNLIDNGSWPATTQNSVALKNAGLISLTKKTNKFIYAIHDETSYKGGKIPKQKSELINCSAQFEANAFAPLTFANAD
 SGLLHSSVRSQSCQNDSEPTLSLTSSFTILRKRCSNRTCSNNTVISQDDLYEAKCNKEKQLFITPEADSLCSLQEGQCQNDPKSKKYSIDKEVEALAAACHPVQ
 HSKVEYSDTDFQSKSLLYDHENASTLILITPTKDVLSLNMVMSRGKESYKMSDLKLG

or Upload a text file [Browse](#)

Set gap score:

[Submit](#)

Program GUI – Main Screen.

Sequence Alignment for Similarity in Biological Data

PAIRWISE SEQUENCE ALIGNMENT
 Scoring using BLOSUM62 with linear gap score

Optimal Alignment:

```
MDQREILQKFLDEAQSCKITKEEFANE-FLKLRQ-S-TK-YKADKTYPT
|||||
M---PIGSK--E---RP-T---FF-EIF-KAEKDCSTTRCNKAD-LGP-

TVAEKPNK-ISK-NRYKDILPYDYSR-VELSITSD-EDSSYINAN-FIK
|||
-IS---LWFEELS--SEAPPYN-SEPAE-E---SEKNNNY-EPNLF-K

GVYGP--K-AVIATQGPL-STTLDFWRMIWEYSVLIIIVMACMEYEMGKK
|||||
-T---PQRKPSY--NQ--LASTPII-F-K---EQG-L--TLP-L-YQ-SPV

KCERYWAEPMQLEFGFFSVSCAEKRRKSDYIIRTLKVKFNSERTRTIYQ
|||||
K-E---LD--KFKLDLGR-NVP-NS-RHKS---LRT--VK---TK-MDQ

FHYKNWPDHVP--P--SS-I-D-PIELFIWDVRC-Y---QEDDSVPIC--
```

Alignment Results with Gap Penalty of -2:

Identity: 249/995 (25.03%)
 Similarity: 619/995 (62.21%)
 Match Count: 249 (25.03%)
 Mismatch Count: 370 (37.19%)
 Gap Count: 376 (37.79%)
 Alignment Path Score: 604

Sequence 1 Length: 807
 Sequence 2 Length: 807
 Aligned Sequence Length: 995

Alignment runtime: 1.18700004s

[Back](#) [Save Matrices](#) [Save Results](#)

Program GUI – Results Screen.



www.tinyurl.com/seqaligner