



# Algorithmen auf Sequenzen

## SoSe 2022

### Projekt 3

## Paarweise Sequenzalignments

**Abgabetermin: 22.06.2022**

**Ansprechpersonen:**

- Georg Grünert: [georg.gruenert@mni.thm.de](mailto:georg.gruenert@mni.thm.de)
- Sujaya Shrestha: [sujaya.shrestha@mni.thm.de](mailto:sujaya.shrestha@mni.thm.de)
- Cornelia Meckbach: [cornelia.meckbach@mni.thm.de](mailto:cornelia.meckbach@mni.thm.de)

**Abgabe:** Beschreiben Sie Ihr Vorgehen und Ihre wichtigsten Erkenntnisse in einem Dokument. Fügen Sie dem Dokument auch Ihren kommentierten Quellcode hinzu, sowie wichtige Ausgaben und speichern Sie es als PDF. Vergessen Sie nicht Ihren Namen dazuzuschreiben. Speichern Sie den Quellcode inkl. ReadMe-file (für Ausführungsdetails), sowie das oben beschriebene PDF-Dokument in einem Ordner und geben Sie diesen entweder als .zip oder als .tar.gz ab.

**Wichtig:** Die Abgabe zu dem Projekt erfolgt in einem „Mini-Bericht“ in dem Sie logisch strukturiert und losgelöst von der Reihenfolge der Aufgaben Ihr Vorgehen und Ihre Ergebnisse präsentieren.

### Projekt 3: Paarweise Sequenzalignments

#### Aufgabe 1: Implementation

Implementieren Sie die beiden unten aufgeführten Algorithmen. Beachten Sie: Ihr Skript sollte zwei Sequenzen (im FASTA-Format), eine Substitutionsmatrix, und eine Gap-Penalty übergeben bekommen und den Score des besten Alignments, sowie das beste Alignment selbst zurückgeben.

- a) Needleman-Wunsch-Algorithmus
- b) Smith-Waterman-Algorithmus

## Aufgabe 2: Testen

Laden Sie sich die Sequenzen des Onkogens KRAS des Menschen (NM\_004985.3) und der Maus (NM\_021284.7) herunter und speichern Sie sie zusammen in einer FASTA-Datei.

- Was ist ein Onkogen?
- Erstellen Sie eine Substitutionsmatrix für Nukleotide, in der Sie Matches mit 2 belohnen und Mismatches mit 0 bewerten.
- Alignen Sie die beiden Sequenzen mit
  - dem Needleman-Wunsch-Algorithmus
  - dem Smith-Waterman-AlgorithmusVerwenden Sie dabei eine Gap-Penalty von -1.
- Berechnen Sie für beide Alignments die Sequenzidentität, d.h. wieviel Prozent der Alignmentlänge ist identisch (Matches)?
- Vergleichen Sie die beiden Alignments hinsichtlich finalem Score, Sequenzidentität und ihrer generellen Struktur.

## Aufgabe 3: Signifikanz von Ähnlichkeit

Oft möchte man wissen, ob die Ähnlichkeiten zwischen zwei Sequenzen wirklich einen biologischen Hintergrund hat, oder aber einfach durch Zufall zustande gekommen sein könnte. Um das herauszufinden, werden die Sequenzen mehrfach (~1000 mal) miteinander alignt, dabei bleibt aber nur eine Sequenz in ihrer eigentlichen Form bestehen, die andere wird permutiert, d.h. die Nukleotide dieser Sequenz tauschen zufällig ihre Positionen.

Gegeben sind Ihnen zwei Dateien: „PTCH2\_TRA2.fasta“ und „Eyeless\_Tra.fasta“.  
„PTCH2\_TRA2.fasta“ beinhaltet die Aminosäuresequenzen des menschlichen PTCH2 Genes (NP\_003729) und die Aminosäuresequenz von TRA2 von *C. elegans* (NC\_003280).  
„Eyeless\_Tra.fasta“ beinhaltet TRA2 und das ey (eyeless) Protein von *D. melanogaster*.

- Was sind die Funktionen von TRA2 in *C. elegans* und eyeless in *D. melanogaster*? (Zusatz: Welche Informationen finden Sie über PTCH2 in *h.sapiens*?)
- Besorgen Sie sich die BLOSUM62 Matrix und lesen Sie sie in Python ein.
- Alignen Sie die Sequenzen aus „PTCH2\_TRA2.fasta“ mit dem Needleman-Wunsch-Algorithmus. Verwenden Sie die BLOSUM62 Matrix und eine Gap-Penalty von -5.
- Erstellen Sie ~1000 der oben beschriebenen „Zufalls-Alignments“ mit dem Needleman-Wunsch-Algorithmus und speichern Sie deren Scores. (Wenn Ihr PC sehr langsam ist, verwenden Sie nur 500 (oder 100) Alignments).
- Erstellen Sie ein Histogramm, in dem Sie Anzahl der (Zufalls-) Alignments (y-Achse) gegen den Score (x-Achse) auftragen. Kennzeichnen Sie dabei den Score des originalen Alignments.
- Ermitteln Sie die Wahrscheinlichkeit, dass ein Alignment mit einem Score  $\geq$  dem Score des Alignments mit den originalen Sequenzen (Aufgabe 2.c) gefunden wird.
- Wiederholen Sie c)-f) für die Sequenzen in „Eyeless\_Tra.fasta“.
- Welche Rückschlüsse ziehen Sie bezüglich der Ähnlichkeit der jeweiligen Sequenzen? Begründen Sie Ihre Annahmen.



THM

CAMPUS  
GIESSEN

MNI

Mathematik, Naturwissenschaften  
und Informatik

TECHNISCHE HOCHSCHULE MITTELHESSEN

#### **Aufgabe 4: Einfluss von Alignmentparametern**

Die behandelten Algorithmen erzeugen Alignments die optimal sind, bezogen auf die Zielfunktion. Verwenden Sie den Smith-Waterman-Algorithmus um die beiden RNA-Sequenzen in der Datei rRNASeq.fasta zu alignen.

- a) Testen Sie unterschiedliche Werte als Gap-Penalties, Match- und Mismatch-Scores.
- b) Welche Auffälligkeiten finden Sie im Hinblick auf die Länge und die Sequenzidentität der resultierenden Alignments? Welchen Einfluss haben die einzelnen Parameter?