

Ein graphtheoretischer Ansatz für das *multiple sequence Alignment*-Problem

„Effiziente Algorithmen“ in einer Bachelorarbeit

1. Motivation
2. Das DIALIGN-Verfahren mit einem Min-Cut-Ansatz für das Konsistenzproblem
3. Umsetzung in der Bachelorarbeit

1. Motivation

2. Das DIALIGN-Verfahren mit einem Min-Cut-Ansatz für das Konsistenzproblem

3. Umsetzung in der Bachelorarbeit

Das *multiple sequence Alignment*-Problem

- ▶ Gegeben ist eine Menge von Zeichenketten.
- ▶ Beispiel:
 1. TCGTCTGCACGCGCTCTGCGAT
 2. AGTCGTCTGCACGGGATCTGCGA
 3. AATAGTCATGGACGCGTGCTCTA
 4. ATAGTCATGGACGCGTGCGCGAT
- ▶ Frage: Wie kann man einzelne Symbole oder ganze Abschnitte dieser Sequenzen einander zuordnen, sodass diejenigen Bereiche übereinander stehen, die sich möglichst ähnlich sind?

Das *multiple sequence Alignment*-Problem

- ▶ Gegeben ist eine Menge von Zeichenketten.
- ▶ Beispiel:
 1. TCGTCTGCACGCGCTCTGCGAT
 2. AGTCGTCTGCACGGGATCTGCGA
 3. AATAGTCATGGACGCGTGCTCTA
 4. ATAGTCATGGACGCGTGCGCGAT
- ▶ Frage: Wie kann man einzelne Symbole oder ganze Abschnitte dieser Sequenzen einander zuordnen, sodass diejenigen Bereiche übereinander stehen, die sich möglichst ähnlich sind?
- ▶ In unserem Beispiel könnte das beispielsweise so aussehen:
 1. --TCGTC-TGCACGC--GCTCTGCGAT
 2. AGTCGTC-TGCACG-G-GATCTGCGA-
 3. AATAGTCATGGACGCGTGCTC---TA-
 4. -ATAGTCATGGACGCGTGCGC---GAT

Wofür wird MSA benötigt?

PALÄOANTHROPOLOGIE

Per DNA-Verlust zum Menschen?

Das große Gehirn, der aufrechte Gang und unser Paarungsverhalten: Dies alles hätten wir vermutlich nicht ohne die Einbuße einiger DNA-Abschnitte, die für andere Primaten wichtig sind.

Philip L. Reno



[Reno18]

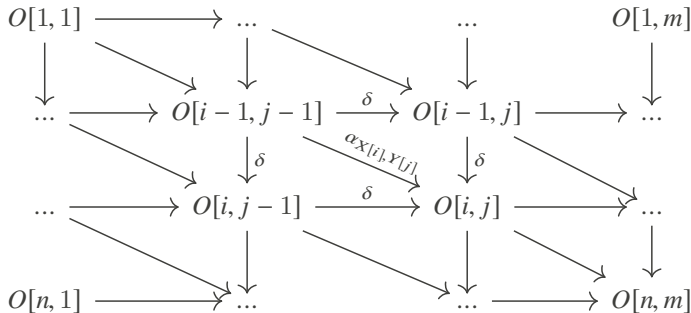
- ▶ Multiples Sequenzalignment wird hauptsächlich in der Bioinformatik angewendet. [MDW96]
- ▶ Verglichen werden üblicherweise DNA- oder Proteinsequenzen.
- ▶ Wird zum Beispiel benutzt, wenn man einen gemeinsamen evolutionären Ursprung zwischen mehreren Sequenzen vermutet oder allgemein, um mögliche Resultate von Mutationen zu finden.
- ▶ Gefunden werden können beispielsweise Punktmutationen (Änderungen einzelner Basen oder Aminosäuren), sowie eingefügte oder gelöschte Abschnitte.

Komplexität I

- ▶ Problem: Unter allen auch nur halbwegs realistischen Annahmen ist das *multiple sequence Alignment*-Problem NP-schwer mit Laufzeit $O(\text{length}^n)$, wobei n die Anzahl der alignierten Sequenzen ist. [WJ94]
- ▶ Beweisskizze mit Hilfe des Algorithmus von Needleman-Wunsch:
- ▶ Kosten für alignieren von Zeichen a mit Zeichen b kostet $\alpha_{a,b}$, wobei typischerweise $\alpha_{a,a} = 0$ gilt.
- ▶ *Gap Penalty* δ , wenn wir die beiden Zeichen nicht miteinander alignieren und stattdessen in eine der Sequenzen eine Lücke einfügen.
- ▶ Gesamtkosten lassen sich mit dynamischer Programmierung und dieser Rekursionsgleichung berechnen: [NW70]

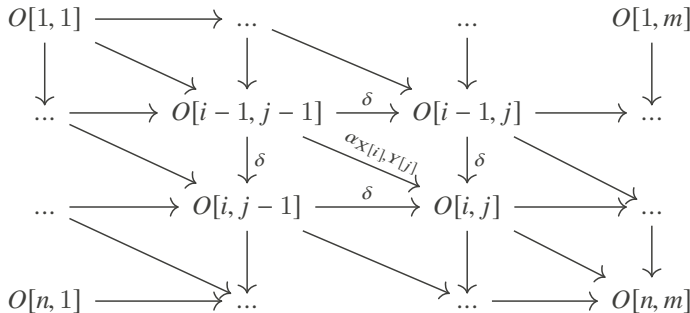
$$O[i, j] := \min\{\alpha(X[i], Y[j]) + O[i - 1, j - 1], \delta + O[i - 1, j], \delta + O[i, j - 1]\} \quad (1)$$

Komplexität II



- Alignment mit minimalen Kosten entspricht Pfad mit geringsten Kosten durch diese Matrix, der in $O(n * m)$ berechenbar ist.
- Vorgehen lässt sich auch mit n-vielen Sequenzen durchführen → Pfad durch n-dimensionale Matrix

Komplexität II



- ▶ Alignment mit minimalen Kosten entspricht Pfad mit geringsten Kosten durch diese Matrix, der in $O(n * m)$ berechenbar ist.
- ▶ Vorgehen lässt sich auch mit n-vielen Sequenzen durchführen → Pfad durch n-dimensionale Matrix
- ▶ große multiple Alignments lassen sich nicht exakt, sondern nur heuristisch berechnen

1. Motivation

2. Das DIALIGN-Verfahren mit einem Min-Cut-Ansatz für das Konsistenzproblem

3. Umsetzung in der Bachelorarbeit

Grundidee hinter DIALIGN

1. Weise den zu alignierenden Abschnitten Gewichte zu.

Grundidee hinter DIALIGN

1. Weise den zu alignierenden Abschnitten Gewichte zu.
2. Maximiere die Summe der Gewichte bei paarweisen Alignments durch dynamische Programmierung.

Grundidee hinter DIALIGN

1. Weise den zu alignierenden Abschnitten Gewichte zu.
2. Maximiere die Summe der Gewichte bei paarweisen Alignments durch dynamische Programmierung.
3. Sortiere alle so entstandenen paarweisen Zuweisungen (Diagonalen in der Matrix von Needleman-Wunsch, daher DIALIGN für „Diagonal ALIGNment“) anhand ihrer Gewichte.

Grundidee hinter DIALIGN

1. Weise den zu alignierenden Abschnitten Gewichte zu.
2. Maximiere die Summe der Gewichte bei paarweisen Alignments durch dynamische Programmierung.
3. Sortiere alle so entstandenen paarweisen Zuweisungen (Diagonalen in der Matrix von Needleman-Wunsch, daher DIALIGN für „Diagonal ALIGNment“) anhand ihrer Gewichte.
4. Wähle gierig nacheinander Diagonalen aus, die zu allen bisher gewählten Diagonalen *konsistent* sind. [MDW96]

Beispielsequenzen und Gewichtsfunktionen

1. TCGTCTGCACG
2. AGTCGTCCGC
3. AATAGTCAT

- ▶ Wahrscheinlichkeit, dass es in einem Abschnitt der Länge l mindestens m viele Übereinstimmungen gibt:

$$P(l, m) := \sum_{i=m}^l \binom{l}{i} * p^i * (1 - p)^{l-i} \quad (2)$$

- ▶ Gewichtsfunktion mit Korrekturterm $K := \ln l_1 + \ln l_2$:

$$w^* := -\ln(P(l, m)) - K \quad (3)$$

[MAHD98]

1. paarweises Alignment

1. TCGTCTGCACG

2. AGTCGTCGCG

$$\begin{matrix} \text{GT} \\ \text{GT} \end{matrix} \left. \vphantom{\begin{matrix} \text{GT} \\ \text{GT} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTC}\textcolor{red}{TGC}\textcolor{red}{ACG} \\ \text{GTC}\textcolor{red}{GTC}\textcolor{red}{GCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC}\textcolor{red}{TGC}\textcolor{red}{ACG} \\ \text{GTC}\textcolor{red}{GTC}\textcolor{red}{GCG} \end{matrix}} \right\} 2,1$$

$$\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix}} \right\} 1,9$$

$$\begin{matrix} \text{TCGTC} \\ \text{TCGTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTC} \\ \text{TCGTC} \end{matrix}} \right\} 4,6$$

$$\begin{matrix} \textcolor{red}{C}\textcolor{red}{ACG} \\ \textcolor{red}{C}\textcolor{red}{GCG} \end{matrix} \left. \vphantom{\begin{matrix} \textcolor{red}{C}\textcolor{red}{ACG} \\ \textcolor{red}{C}\textcolor{red}{GCG} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTC}\textcolor{red}{TGC} \\ \text{GTC}\textcolor{red}{GTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC}\textcolor{red}{TGC} \\ \text{GTC}\textcolor{red}{GTC} \end{matrix}} \right\} 0,7$$

$$\begin{matrix} \text{GTC}\textcolor{red}{TGCA} \\ \text{GTC}\textcolor{red}{GTCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC}\textcolor{red}{TGCA} \\ \text{GTC}\textcolor{red}{GTCG} \end{matrix}} \right\} 0,2$$

$$\begin{matrix} \text{GTC}\textcolor{red}{TGCAC} \\ \text{GTC}\textcolor{red}{GTCGC} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC}\textcolor{red}{TGCAC} \\ \text{GTC}\textcolor{red}{GTCGC} \end{matrix}} \right\} 1,1$$

$$\begin{matrix} \text{CG} \\ \text{CG} \end{matrix} \left. \vphantom{\begin{matrix} \text{CG} \\ \text{CG} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{TC} \\ \text{TC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TC} \\ \text{TC} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \textcolor{red}{CGTC} \\ \textcolor{red}{AGTC} \end{matrix} \left. \vphantom{\begin{matrix} \textcolor{red}{CGTC} \\ \textcolor{red}{AGTC} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTCT} \\ \text{GTCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCT} \\ \text{GTCG} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{TCGTCT} \\ \text{TCGTGC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTCT} \\ \text{TCGTGC} \end{matrix}} \right\} 2,8$$

$$\begin{matrix} \text{TCGTCT}\textcolor{red}{G} \\ \text{TCGTGC}\textcolor{red}{G} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTCT}\textcolor{red}{G} \\ \text{TCGTGC}\textcolor{red}{G} \end{matrix}} \right\} 1,8$$

$$\begin{matrix} \text{TCGTCT}\textcolor{red}{TGC} \\ \text{TCGTGC}\textcolor{red}{GCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTCT}\textcolor{red}{TGC} \\ \text{TCGTGC}\textcolor{red}{GCG} \end{matrix}} \right\} 1,1$$

1. paarweises Alignment

1. TCGTCTGCACG

2. AGTCGTCGCG

$$\begin{matrix} \text{GT} \\ \text{GT} \end{matrix} \left. \vphantom{\begin{matrix} \text{GT} \\ \text{GT} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTCTGCACG} \\ \text{GTCGTCGCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCTGCACG} \\ \text{GTCGTCGCG} \end{matrix}} \right\} 2,1$$

$$\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix}} \right\} 1,9$$

$$\begin{matrix} \text{TCGTC} \\ \text{TCGTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTC} \\ \text{TCGTC} \end{matrix}} \right\} 4,6$$

$$\begin{matrix} \text{CACG} \\ \text{CGCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{CACG} \\ \text{CGCG} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTCTGC} \\ \text{GTCGTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCTGC} \\ \text{GTCGTC} \end{matrix}} \right\} 0,7$$

$$\begin{matrix} \text{GTCTGCA} \\ \text{GTCGTCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCTGCA} \\ \text{GTCGTCG} \end{matrix}} \right\} 0,2$$

$$\begin{matrix} \text{GTCTGCAC} \\ \text{GTCGTCGCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCTGCAC} \\ \text{GTCGTCGCG} \end{matrix}} \right\} 1,1$$

$$\begin{matrix} \text{CG} \\ \text{CG} \end{matrix} \left. \vphantom{\begin{matrix} \text{CG} \\ \text{CG} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{TC} \\ \text{TC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TC} \\ \text{TC} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{CGTC} \\ \text{AGTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{CGTC} \\ \text{AGTC} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTCT} \\ \text{GTCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCT} \\ \text{GTCG} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{TCGTCT} \\ \text{TCGTCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTCT} \\ \text{TCGTCG} \end{matrix}} \right\} 2,8$$

$$\begin{matrix} \text{TCGTCTG} \\ \text{TCGTCGC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTCTG} \\ \text{TCGTCGC} \end{matrix}} \right\} 1,8$$

$$\begin{matrix} \text{TCGTCTGC} \\ \text{TCGTCGCG} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTCTGC} \\ \text{TCGTCGCG} \end{matrix}} \right\} 1,1$$

► Resultierendes Alignment:

--TCGTCTgcaCG

agTCGTCg---CG

► Score: 5,1

2. paarweises Alignment

1. TCGTCTGCACG

2. AATAGTCAT

$$\begin{matrix} \text{TC} \\ \text{TC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TC} \\ \text{TC} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GT} \\ \text{GT} \end{matrix} \left. \vphantom{\begin{matrix} \text{GT} \\ \text{GT} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{CA} \\ \text{CA} \end{matrix} \left. \vphantom{\begin{matrix} \text{CA} \\ \text{CA} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix}} \right\} 1,9$$

$$\begin{matrix} \text{TCGT} \\ \text{TAGT} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGT} \\ \text{TAGT} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{TCGTCT} \\ \text{TAGTCT} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGTCT} \\ \text{TAGTCT} \end{matrix}} \right\} 1,6$$

$$\begin{matrix} \text{GTCT} \\ \text{GTCA} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCT} \\ \text{GTCA} \end{matrix}} \right\} 0,5$$

2. paarweises Alignment

1. TCGTCTGCACG

2. AATAGTCAT

$$\begin{array}{ccccc}
 \begin{array}{l} \text{TC} \\ \text{TC} \end{array} \left. \vphantom{\begin{array}{l} \text{TC} \\ \text{TC} \end{array}} \right\} 0,5 &
 \begin{array}{l} \text{GT} \\ \text{GT} \end{array} \left. \vphantom{\begin{array}{l} \text{GT} \\ \text{GT} \end{array}} \right\} 0,5 &
 \begin{array}{l} \text{CA} \\ \text{CA} \end{array} \left. \vphantom{\begin{array}{l} \text{CA} \\ \text{CA} \end{array}} \right\} 0,5 &
 \begin{array}{l} \text{GTC} \\ \text{GTC} \end{array} \left. \vphantom{\begin{array}{l} \text{GTC} \\ \text{GTC} \end{array}} \right\} 1,9 &
 \begin{array}{l} \text{TCGT} \\ \text{TAGT} \end{array} \left. \vphantom{\begin{array}{l} \text{TCGT} \\ \text{TAGT} \end{array}} \right\} 0,5 \\
 \begin{array}{l} \text{TCGTCT} \\ \text{TAGTC} \end{array} \left. \vphantom{\begin{array}{l} \text{TCGTCT} \\ \text{TAGTC} \end{array}} \right\} 1,6 &
 \begin{array}{l} \text{GTCT} \\ \text{GTCA} \end{array} \left. \vphantom{\begin{array}{l} \text{GTCT} \\ \text{GTCA} \end{array}} \right\} 0,5 & & &
 \end{array}$$

► Resultierendes Alignment:

```
--tcGTCTgcacg
aataGTCat----
```

► Score: 1,9

3. paarweises Alignment

1. AGTCGTCGCG

2. AATAGTCAT

$$\begin{matrix} \text{TC} \\ \text{TC} \end{matrix} \left. \vphantom{\begin{matrix} \text{TC} \\ \text{TC} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTC} \\ \text{GTC} \end{matrix}} \right\} 1,9$$

$$\begin{matrix} \text{AGTC} \\ \text{AGTC} \end{matrix} \left. \vphantom{\begin{matrix} \text{AGTC} \\ \text{AGTC} \end{matrix}} \right\} 3,2$$

$$\begin{matrix} \text{AGTCG} \\ \text{AGTCA} \end{matrix} \left. \vphantom{\begin{matrix} \text{AGTCG} \\ \text{AGTCA} \end{matrix}} \right\} 1,6$$

$$\begin{matrix} \text{AGTCGT} \\ \text{AGTCAT} \end{matrix} \left. \vphantom{\begin{matrix} \text{AGTCGT} \\ \text{AGTCAT} \end{matrix}} \right\} 2,8$$

$$\begin{matrix} \text{GTCGT} \\ \text{GTCAT} \end{matrix} \left. \vphantom{\begin{matrix} \text{GTCGT} \\ \text{GTCAT} \end{matrix}} \right\} 1,6$$

$$\begin{matrix} \text{TCGT} \\ \text{TCAT} \end{matrix} \left. \vphantom{\begin{matrix} \text{TCGT} \\ \text{TCAT} \end{matrix}} \right\} 0,5$$

$$\begin{matrix} \text{AGT} \\ \text{AGT} \end{matrix} \left. \vphantom{\begin{matrix} \text{AGT} \\ \text{AGT} \end{matrix}} \right\} 1,9$$

3. paarweises Alignment

1. AGTCGTCGCG

2. AATAGTCAT

$$\begin{array}{ccccc}
 \begin{array}{l} \text{TC} \\ \text{TC} \end{array} \left. \vphantom{\begin{array}{l} \text{TC} \\ \text{TC} \end{array}} \right\} 0,5 &
 \begin{array}{l} \text{GTC} \\ \text{GTC} \end{array} \left. \vphantom{\begin{array}{l} \text{GTC} \\ \text{GTC} \end{array}} \right\} 1,9 &
 \begin{array}{l} \text{AGTC} \\ \text{AGTC} \end{array} \left. \vphantom{\begin{array}{l} \text{AGTC} \\ \text{AGTC} \end{array}} \right\} 3,2 &
 \begin{array}{l} \text{AGTCG} \\ \text{AGTCA} \end{array} \left. \vphantom{\begin{array}{l} \text{AGTCG} \\ \text{AGTCA} \end{array}} \right\} 1,6 &
 \begin{array}{l} \text{AGTCGT} \\ \text{AGTCAT} \end{array} \left. \vphantom{\begin{array}{l} \text{AGTCGT} \\ \text{AGTCAT} \end{array}} \right\} 2,8 \\
 \\
 \begin{array}{l} \text{GTCGT} \\ \text{GTCAT} \end{array} \left. \vphantom{\begin{array}{l} \text{GTCGT} \\ \text{GTCAT} \end{array}} \right\} 1,6 &
 \begin{array}{l} \text{TCGT} \\ \text{TCAT} \end{array} \left. \vphantom{\begin{array}{l} \text{TCGT} \\ \text{TCAT} \end{array}} \right\} 0,5 &
 \begin{array}{l} \text{AGT} \\ \text{AGT} \end{array} \left. \vphantom{\begin{array}{l} \text{AGT} \\ \text{AGT} \end{array}} \right\} 1,9 & &
 \end{array}$$

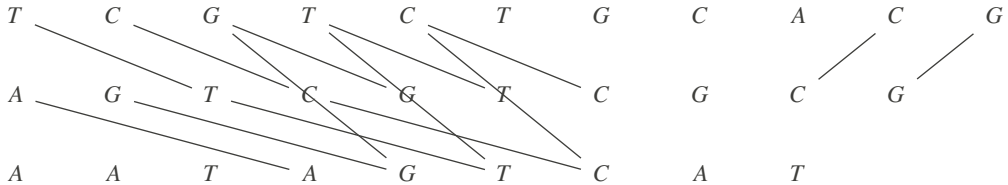
► Resultierendes Alignment:

```

---AGTCgtcgcg
aatAGTCac----
```

► Score: 3,2

Konsistenz und Zwischenstand



- ▶ nicht *konsistent*, denn „TC“ aus der dritten Sequenz sind transitiv zwei „TC“-Abschnitte aus der ersten Sequenz zugeordnet und es gibt Überkreuzungen.

Gieriges multiples Alignment

- ▶ Sortiere alle Abschnitte aus den paarweisen Alignments nach ihrem Score.
- ▶ Wähle solange den Abschnitt mit dem höchsten Score, der mit allen bisher gewählten *konsistent* ist, bis es keinen solchen mehr gibt.

Gieriges multiples Alignment

- ▶ Sortiere alle Abschnitte aus den paarweisen Alignments nach ihrem Score.
- ▶ Wähle solange den Abschnitt mit dem höchsten Score, der mit allen bisher gewählten *konsistent* ist, bis es keinen solchen mehr gibt.

	Sequenzen	Abschnitte	Score	zu den bisher gewählten <i>konsistent</i> ?
	1	TCGTC	4,6	ja
	2	TCGTC		
	2	AGTC	3,2	ja
▶	3	AGTC		
	1	GTC	1,9	nein
	3	GTC		
	1	CG	0,5	ja
	2	CG		

Gieriges multiples Alignment

- Sortiere alle Abschnitte aus den paarweisen Alignments nach ihrem Score.
- Wähle solange den Abschnitt mit dem höchsten Score, der mit allen bisher gewählten *konsistent* ist, bis es keinen solchen mehr gibt.

	Sequenzen	Abschnitte	Score	zu den bisher gewählten <i>konsistent</i> ?
	1	TCGTC	4,6	ja
	2	TCGTC		
	2	AGTC	3,2	ja
▶	3	AGTC		
	1	GTC	1,9	nein
	3	GTC		
	1	CG	0,5	ja
	2	CG		

- Resultierendes gieriges Alignment:


```

-----TCGTctgcaCG
---AGTCGTCg---CG
aatAGTCat-----

```

Erweiterung mit Min-Cut-Ansatz

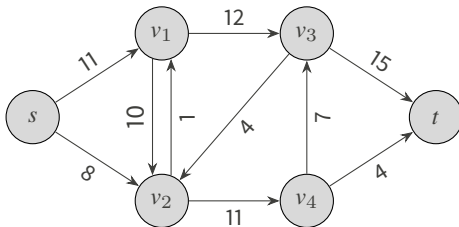
- ▶ Problem: Das gierige Konstruieren des multiplen Alignments zieht nur paarweise ähnliche Abschnitte in Betracht. [CPM10]
- ▶ Wir wollen aber Ähnlichkeiten über möglichst viele Sequenzen hinweg.

Erweiterung mit Min-Cut-Ansatz

- ▶ Problem: Das gierige Konstruieren des multiplen Alignments zieht nur paarweise ähnliche Abschnitte in Betracht. [CPM10]
- ▶ Wir wollen aber Ähnlichkeiten über möglichst viele Sequenzen hinweg.
- ▶ Alternativer Ansatz:
 1. Konstruiere Graph über alle alignierten *Sites* aus dem paarweisen Alignment (Inzidenzgraph).
 2. Entferne mit Hilfe von Flussnetzen Kanten aus den Zusammenhangskomponenten, um Inkonsistenzen aufzulösen.
 3. Konstruiere Graph über die Zusammenhangskomponenten des Inzidenzgraphen (Sukzessionsgraph) und entferne geschickt Kanten und *Sites* aus dem Graph, bis alle Inkonsistenzen aufgelöst sind.

Flussnetzwerke I

- Ein Flussnetzwerk ist ein gerichteter Graph $G = (V, E)$ mit zwei ausgewiesenen Knoten $s, t \in V$, genannt Quelle und Senke. Zusätzlich hat jede Kante $(u, v) \in E$ ein Kapazität genanntes, nichtnegatives Kantengewicht. [CLRS92]

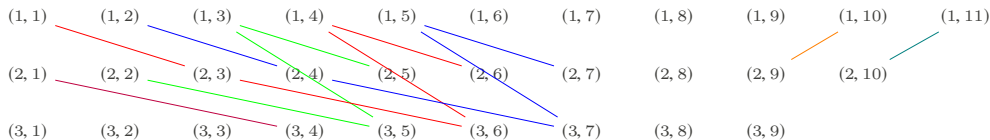


- Ziel: Berechnung des maximalen Flusses von der Quelle zur Senke.
- Jeder Kante wird ein Wert zugeordnet, der die Kapazität nicht überschreiten darf (Kapazitätsbeschränkung) und für jeden Knoten muss die Summe der einfließenden Flüsse der Summe der ausfließenden entsprechen (Flusserhaltung).

Flussnetzwerke II

- ▶ Ein $(s-t)$ -Schnitt ist eine Partitionierung der Knotenmenge V in zwei Mengen A und B , bei der die Quelle $s \in A$ und die Senke $t \in B$ ist.
- ▶ Die Kapazität eines $(s-t)$ -Schnitts ist die Summe der Kapazitäten aller Kanten, die von Knoten aus A nach B gehen.
- ▶ Der *Min-Cut-Max-Flow-Satz* besagt, dass die Kapazität des minimalen Schnitts dem Wert des maximalen Flusses entspricht.
- ▶ Es lässt sich zeigen, dass nach Berechnung des maximalen Flusses auch der minimale Schnitt bekannt ist. [CLRS92]

Inzidenzgraph



► Konstruiere Inzidenzgraph:

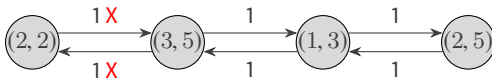
Alle *Sites* entsprechen Knoten und zwischen zwei Knoten existiert genau dann eine Kante, wenn die entsprechenden Sites in einem der paarweisen DIALIGN-Fragmenten einander zugeordnet waren.

► Zusammenhangskomponenten sind farbig markiert.

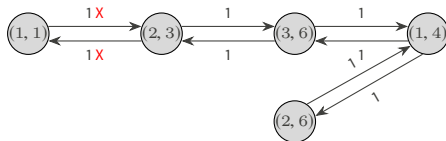
- Eine Zusammenhangskomponente wird als *mehrdeutig* bezeichnet, wenn in ihr zwei oder mehr *Sites* aus der selben Sequenz existieren → Inkonsistenz. Genauso werden diese *Sites* als *mehrdeutig* bezeichnet. [CPM10]

Minimaler Schnitt auf Zusammenhangskomponenten I

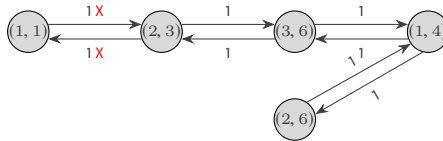
- ▶ Fasse Zusammenhangskomponenten des Inzidenzgraphen als Flussnetzwerke auf.
- ▶ Ersetze ungerichtete Kanten durch zwei gerichtete mit Kapazität 1.
- ▶ Setze jeweils zwei *mehrdeutige* Knoten aus der selben Sequenz als Quelle und Senke.
- ▶ Berechne den minimalen Schnitt und lösche die entsprechenden Kanten.
- ▶ Wiederhole dies, bis es keine *mehrdeutigen* Zusammenhangskomponenten mehr gibt.



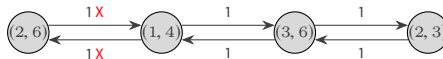
Minimaler Schnitt auf Zusammenhangskomponenten II



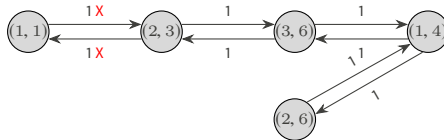
Minimaler Schnitt auf Zusammenhangskomponenten II



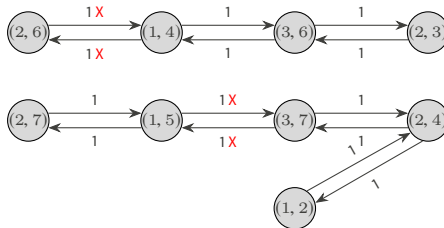
- ▶ Nach wie vor zwei *mehrdeutige* Sites in der Zusammenhangskomponente: (2,3) und (2,6).
- ▶ Berechne nächsten Min-Cut nach Entfernen der ersten Kante.



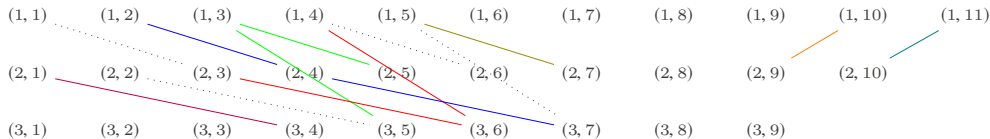
Minimaler Schnitt auf Zusammenhangskomponenten II



- ▶ Nach wie vor zwei *mehrdeutige* Sites in der Zusammenhangskomponente: $(2,3)$ und $(2,6)$.
- ▶ Berechne nächsten Min-Cut nach Entfernen der ersten Kante.

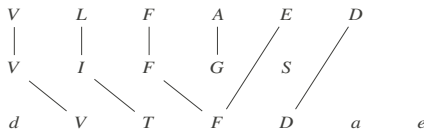


Unser Graph mit gelöschten Kanten

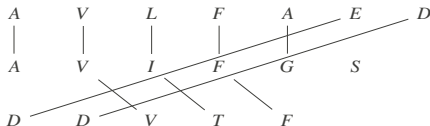


- ▶ Resultat: immer noch nicht *konsistent*.
- ▶ Grund: Überkreuzungen z.B. bei $(1,2) \rightarrow (2,4) \rightarrow (3,7)$ und $(1,3) \rightarrow (3,5)$.

Zwei Arten nichtkonsistenter Zuweisungen



Implizite, transitive Mehrfachzuweisung von F: kann mit Flussnetzen und Inzidenzgraphen behoben werden.



Überkreuzungen: können wir noch nicht auflösen.

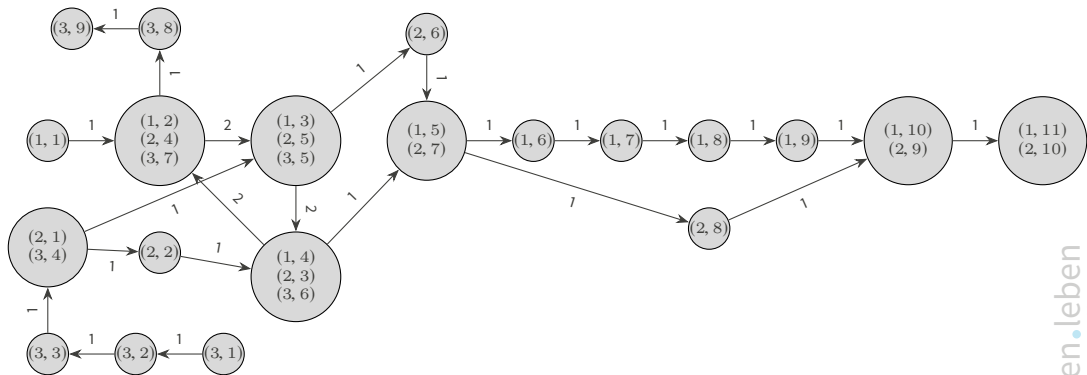
Sukzessionsgraph I

- ▶ Im Sukzessionsgraph sind die Zusammenhangskomponenten des Inzidenzgraphen Knoten. [PDC10]
- ▶ Es existiert eine Kante zwischen zwei Knoten C und C' genau dann, wenn eine Sequenz S_i existiert, sodass $Sites(i, p) \in C$ und $(i, p') \in C'$ existieren mit $p < p'$, ohne dass C'' existiert mit $(i, p'') \in C''$ und $p < p'' < p'$.
- ▶ Das Gewicht jeder Kante ist die Anzahl der Sequenzen für die die obige Bedingung gilt.

Sukzessionsgraph I

- ▶ Im Sukzessionsgraph sind die Zusammenhangskomponenten des Inzidenzgraphen Knoten. [PDC10]
- ▶ Es existiert eine Kante zwischen zwei Knoten C und C' genau dann, wenn eine Sequenz S_i existiert, sodass $Sites(i, p) \in C$ und $(i, p') \in C'$ existieren mit $p < p'$, ohne dass C'' existiert mit $(i, p'') \in C''$ und $p < p'' < p'$.
- ▶ Das Gewicht jeder Kante ist die Anzahl der Sequenzen für die die obige Bedingung gilt.
- ▶ Wenn der Sukzessionsgraph keine Zyklen enthält, ist das Alignment *konsistent*.
- ▶ Falls doch, müssen wir sukzessive *Sites* und Kanten entfernen, bis das der Fall ist.

Sukzessionsgraph II



Zyklen im Sukzessionsgraph und wie wir diese entfernen

- ▶ Herstellen von *Konsistenz* durch Löschen von Kanten und *Sites* aus dem Sukzessionsgraph. [PDC10]
- ▶ Optimal wäre Löschen von Kanten mit minimalem Kantengewicht, um alle Zyklen aufzulösen.
- ▶ Problem:

Zyklen im Sukzessionsgraph und wie wir diese entfernen

- ▶ Herstellen von *Konsistenz* durch Löschen von Kanten und *Sites* aus dem Sukzessionsgraph. [PDC10]
- ▶ Optimal wäre Löschen von Kanten mit minimalem Kantengewicht, um alle Zyklen aufzulösen.
- ▶ Problem: Das ist das NP-schwere *minimal weighted feedback arc set*-Problem.

Zyklen im Sukzessionsgraph und wie wir diese entfernen

- ▶ Herstellen von *Konsistenz* durch Löschen von Kanten und *Sites* aus dem Sukzessionsgraph. [PDC10]
- ▶ Optimal wäre Löschen von Kanten mit minimalem Kantengewicht, um alle Zyklen aufzulösen.
- ▶ Problem: Das ist das NP-schwere *minimal weighted feedback arc set*-Problem.
- ▶ Stattdessen löschen wir solange die Kanten mit dem geringsten Kantengewicht, bis der Graph zyklensfrei ist:

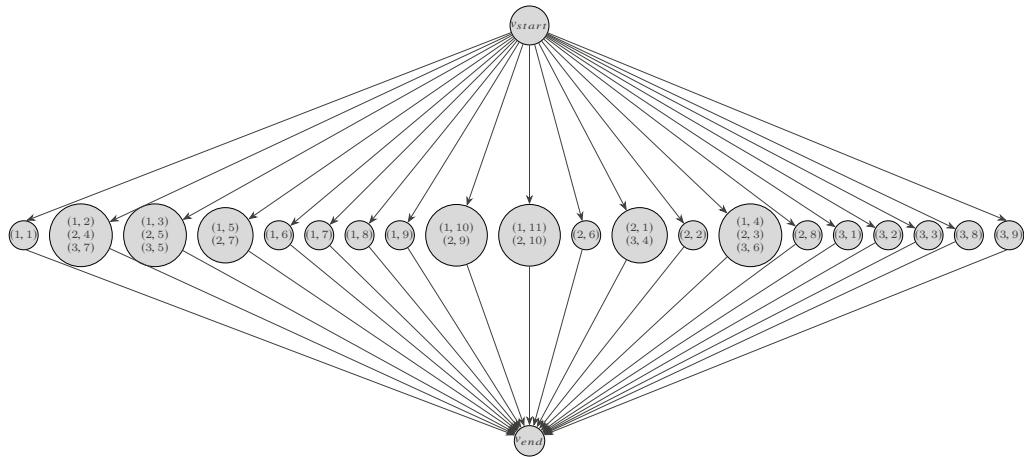
$$E_k := \{(u, v) \in E \mid w(u, v) > k\} \text{ und } k^* := \min\{k > 0 \mid (V, E_k) \text{ ist azyklisch}\}$$
- ▶ Zusätzlich fügen wir Start- und Endknoten ein, die mit den „Anfangs-“ und „Abschlussknoten“ aller Zusammenhangskomponenten nach dem Löschen verbunden sind.

Zyklen im Sukzessionsgraph und wie wir diese entfernen

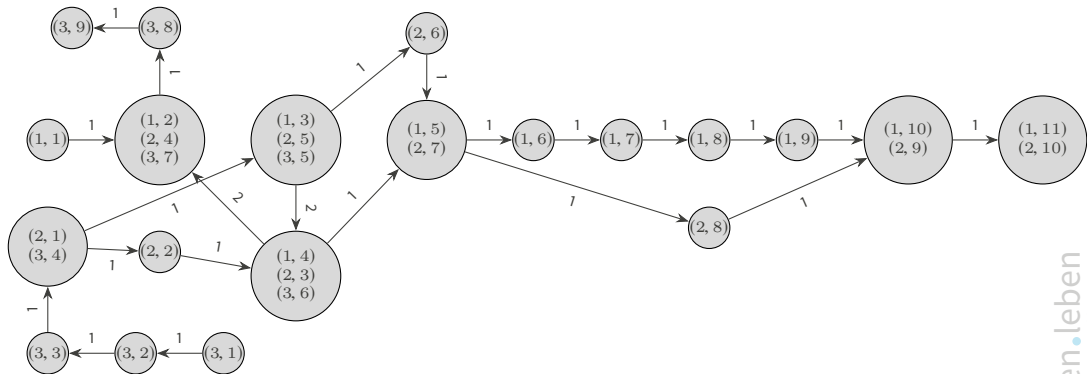
- ▶ Herstellen von *Konsistenz* durch Löschen von Kanten und *Sites* aus dem Sukzessionsgraph. [PDC10]
- ▶ Optimal wäre Löschen von Kanten mit minimalem Kantengewicht, um alle Zyklen aufzulösen.
- ▶ Problem: Das ist das NP-schwere *minimal weighted feedback arc set*-Problem.
- ▶ Stattdessen löschen wir solange die Kanten mit dem geringsten Kantengewicht, bis der Graph zyklensfrei ist:

$$E_k := \{(u, v) \in E \mid w(u, v) > k\} \text{ und } k^* := \min\{k > 0 \mid (V, E_k) \text{ ist azyklisch}\}$$
- ▶ Zusätzlich fügen wir Start- und Endknoten ein, die mit den „Anfangs-“ und „Abschlussknoten“ aller Zusammenhangskomponenten nach dem Löschen verbunden sind.
- ▶ Leider gilt bei uns $k^* = 2 \rightarrow$ alle Kanten zwischen den ursprünglichen Knoten wurden gelöscht.

Das Dilemma



Erschummelter DAG

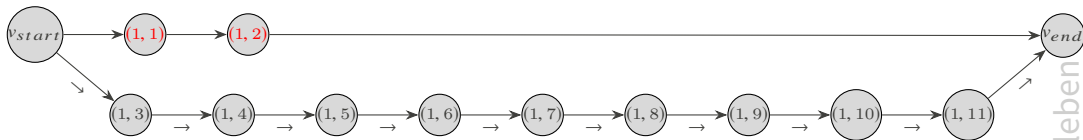


Folge der längsten Kette I

- ▶ Erstelle vom modifizierten Sukzessionsgraphen die transitive Hülle.
- ▶ Folge der längsten Kette für jede Sequenz S_i ; dabei können nur Kanten (C, C') besucht werden, falls $(i, p) \in C$ und $(i, p') \in C'$ mit $p < p'$.
- ▶ Lösche alle *Sites* aus den Knoten, die nicht besucht wurden. [PDC10]

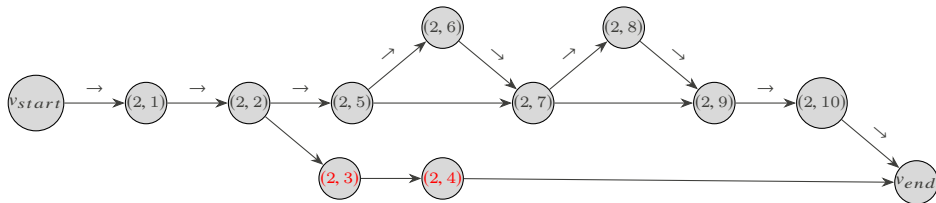
Folge der längsten Kette I

- ▶ Erstelle vom modifizierten Sukzessionsgraphen die transitive Hülle.
- ▶ Folge der längsten Kette für jede Sequenz S_i ; dabei können nur Kanten (C, C') besucht werden, falls $(i, p) \in C$ und $(i, p') \in C'$ mit $p < p'$.
- ▶ Lösche alle *Sites* aus den Knoten, die nicht besucht wurden. [PDC10]



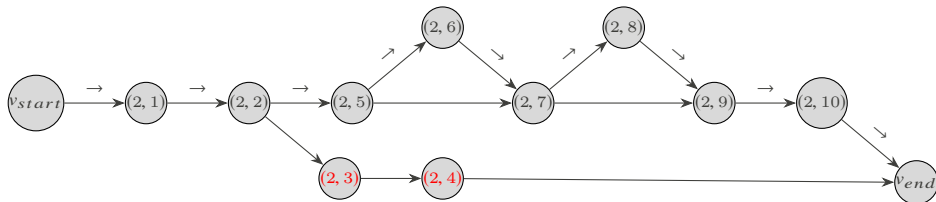
- ▶ Lösche (1, 1) und (1, 2) aus den Knoten im Sukzessionsgraph.

Folge der längsten Kette II



- Lösche (2, 3) und (2, 4).

Folge der längsten Kette II



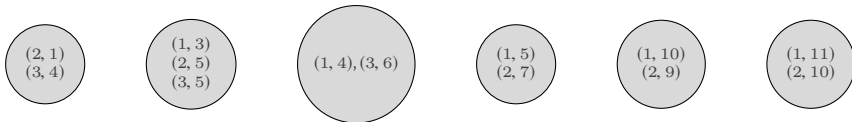
► Lösche (2, 3) und (2, 4).



► Nichts zu löschen.

Das resultierende multiple Alignment

- Damit bleiben nach dem Löschen die folgenden Alignments bestehen:



- Resultierendes Alignment:

1. -----tcGTctgcaCG
2. ---AgtcGtCg---CG
3. aatA---GTcat----

- Schlechteres Alignment als gieriger Ansatz; funktioniert bei größeren Sequenzen hoffentlich besser.

Vereinfachungen, die ich vorgenommen habe

- ▶ Es werden nicht nur die Gewichte der Fragmente selbst betrachtet, sondern zusätzlich die Überschneidungen mit anderen (*Überlappgewichte*).
- ▶ Der Korrekturterm K wurde von mir halbiert. Er funktioniert wohl eher bei langen Sequenzen.
- ▶ Das aktualisieren der Konsistenzgrenzen ist recht kompliziert. Das Verfahren wurde hier nicht vorgestellt. [MDW96]
- ▶ Besserer Ansatz von Abdeddaim. [Abdeddaïm97]
- ▶ Nach dem Min-Cut-Ansatz werden diese Zuweisungsspalten als Ankerpunkte für das normale DIALIGN benutzt. [MPPS06]

1. Motivation

2. Das DIALIGN-Verfahren mit einem Min-Cut-Ansatz für das Konsistenzproblem

3. Umsetzung in der Bachelorarbeit

Gierige Algorithmen und Dynamische Programmierung

- ▶ Gierige Algorithmen:
 - ▶ Treffe lokal die bestmögliche Entscheidung und verwirfe diese danach nicht wieder.
 - ▶ Gute Laufzeit, aber liefern oft nicht das optimale Ergebnis (→ gut geeignet für Heuristiken).
 - ▶ Wird in DIALIGN beim Zusammensetzen des multiplen Alignments benutzt.
 - ▶ Beispiele: Algorithmus von Prim, Fractional Knapsack-Problem

Gierige Algorithmen und Dynamische Programmierung

- ▶ Gierige Algorithmen:
 - ▶ Treffe lokal die bestmögliche Entscheidung und verwurfe diese danach nicht wieder.
 - ▶ Gute Laufzeit, aber liefern oft nicht das optimale Ergebnis (→ gut geeignet für Heuristiken).
 - ▶ Wird in DIALIGN beim Zusammensetzen des multiplen Alignments benutzt.
 - ▶ Beispiele: Algorithmus von Prim, Fractional Knapsack-Problem
- ▶ Dynamische Programmierung:
 - ▶ Berechnung der optimalen Lösung durch Kombination von optimalen Lösungen sich überlappender Teilprobleme, wobei diese gespeichert werden und bei Bedarf abgerufen werden können.
 - ▶ Beispiele: Algorithmus von Needleman-Wunsch, Knapsack-Problem

Flussnetze und der Push-Relabel-Algorithmus

► Flussnetze:

- gerichteter, gewichteter Graph mit ausgewiesener Quelle und Senke; Gewichte der Kante entsprechen maximaler Kapazität an Fluss, die über diese fließen können.
- Ziel: Maximierung des Gesamtflusses eines Flussnetzwerkes von der Quelle zur Senke.
- Min-Cut-Max-Flow-Satz: der Wert eines maximalen Flusses entspricht dem Wert des minimalen Schnitts → benötigen wir zum Auflösen von Inkonsistenzen im Inzidenzgraphen.

Flussnetze und der Push-Relabel-Algorithmus

- ▶ Flussnetze:
 - ▶ gerichteter, gewichteter Graph mit ausgewiesener Quelle und Senke; Gewichte der Kante entsprechen maximaler Kapazität an Fluss, die über diese fließen können.
 - ▶ Ziel: Maximierung des Gesamtflusses eines Flussnetzwerkes von der Quelle zur Senke.
 - ▶ Min-Cut-Max-Flow-Satz: der Wert eines maximalen Flusses entspricht dem Wert des minimalen Schnitts → benötigen wir zum Auflösen von Inkonsistenzen im Inzidenzgraphen.
- ▶ Push-Relabel-Algorithmus:
 - ▶ Zur Berechnung des maximalen Flusses.
 - ▶ besitzt eine bessere Laufzeit, als der in der Veröffentlichung benutzte Edmonds-Karps-Algorithmus.

Algorithmus im Detail

- ▶ Gewichtungsfunktionen für DNA- und Proteinsequenzen (BLOSUM62) mit Überlappungsgewichten.
- ▶ Berechnung von paarweisen Alignments mit Hilfe von dynamischer Programmierung.
- ▶ Aufbau des Inzidenzgraphen und Min-Cut auf den Zusammenhangskomponenten.
- ▶ Aufbau des Sukzessionsgraphen und Löschen von Sites, um Konsistenz herzustellen.
- ▶ Verwendung der so gefundenen alignierten Spalten als Ankerpunkte, um zwischen diesen das klassische DIALIGN laufen zu lassen.
- ▶ Wenn mir am Ende noch langweilig ist: Blick auf DIALIGN TX, um dieses gegebenenfalls zwischen den Ankerpunkten zu verwenden.

Programmierung

- ▶ Programmiersprache C++
- ▶ Boost Graph Library, die bereits ausgeklügelte Graphen, Flussnetzwerke und Algorithmen zur Berechnung des maximalen Flusses bereitstellt
- ▶ variabelster Teil der Bachelorarbeit:
 - ▶ bei wenig Zeit am Ende: Ein- und Ausgabe über simple Textdateien.
 - ▶ bei viel Zeit: Umsetzung mit Formaten wie Clustal und FASTA, sowie Visualisierung über externe Visualisierungssoftware.

Beurteilung der Güte der Ergebnisse

- ▶ Bewertung des Algorithmus ist schwierig, weil er definitiv nicht das perfekte Alignment berechnet, ich aber nicht das Expertenwissen besitze, um die Güte eines Alignments wirklich beurteilen zu können.
- ▶ Test auf BALiBase für global und (D)IRMBASE für lokal verwandte Sequenzen → bereits ausgewertete Sequenzen für die ein richtiges Alignment bekannt ist.
- ▶ Erfordert noch Einarbeitung.

Fragen, Anregungen, Wünsche?

Vielen Dank für Ihre Aufmerksamkeit!

Quellenangaben



Introduction to Algorithms. Cormen et al, MIT Press 1992



Multiple Sequence Alignment with user-defined anchor points. Morgenstern et al, Algorithms for Molecular Biology 2006



Automatic detection of anchor points for multiple sequence alignment. Pitschi et al, BMC Bioinformatics 2010



Segment-based scores for pairwise and multiple sequence alignments. Morgenstern et al, ISMB-98 Proceedings



On Incremental Computation of Transitive Closure and Greedy Alignment. Abdeddaïm, in *Proceedings of the 8th Annual Symposium on Combinatorial Pattern Matching* 1997



A min-cut algorithm for the consistency problem in multiple sequence alignment. Corel et al, Oxford University Press 2010



Multiple DNA and protein sequence alignment based on segment-to-segment comparison. Morgenstern et al, 1996



Per DNA-Verlust zum Menschen?. Reno, Spektrum der Wissenschaft Februar 2018



On the complexity of multiple sequence alignment. Wang et al, Journal of Computational Biology 1994



A general method applicable to the search for similarities in the amino acid sequence of two proteins. Needleman und Wunsch, Journal of Molecular Biology 1970