



Bsp 1) Dann liest  $y$  der Overlap von  $s$  und  $t$   $ov(s,t)$ ,  $ov(s,t) = |ov(s,t)|$

$xyz$  der Menge von  $s$  und  $t$   $\langle s, t \rangle$

$pref(s,t) := x$ ,  $pref(s,t) = |pref(s,t)|$

$suft(s,t) := z$ ,  $suft(s,t) = |suft(s,t)|$

Verallgemeinertes Overlap, wie oben, ohne (ii)  $\overline{ov}(s,t)$

$$\overline{ov}(s,t) = |\overline{ov}(s,t)|$$

## 4. String-Algorithmen

### 4.1. Das String-Matching-Problem

Def. 4.1: Sei  $\Sigma$  ein Alphabet, das (exakte) String-Matching-Problem ist das folgende Berechnungsproblem:

Eingabe: Zwei Strings  $t = t_1 \dots t_n \in \Sigma^n$  (Text) und

$P = P_1 \dots P_m \in \Sigma^m$  (Muster)

Ausgabe: die Menge  $I \subseteq \{1, \dots, n-m+1\}$ , so daß  $i \in I$  gdw

$$t_i \dots t_{i+m-1} = P$$

### Algorithmus 4.1: Naives String-Matching

Eingabe: Muster  $P = P_1 \dots P_m$ , Text  $t = t_1 \dots t_n$

1.  $I = \emptyset$

2. for  $j := 0$  to  $n-m$  do

$i := 1$   
    while  $P_i = t_{j+i}$  and  $i \leq m$  do  
        ~~increase~~  $i := i+1$

    if  $i = m+1$  then  
         $I = I \cup \{j+1\}$

Ausgabe:  $I$

Laufzeit:  $\mathcal{O}(n \cdot (n-m))$

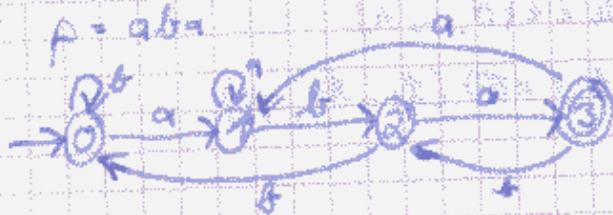
schlechtester Fall:  $P = a^m$ ,  $t = a^n$

Verbesserung der Laufzeit: Preprocessing des Musters

4.2. String-Matching-Automaten

Idee: Konstruiere EA, der für geg. Muster  $P_1 \dots P_m$  alle Texte akzeptiert, die mit  $p$  enden.

Beispiel:  $p = aba$



Def. 4.1: Sei  $P = P_1 \dots P_m \in \Sigma^m$

der String-Matching-Automat für  $p$  ist der EA  $M_P(Q, \Sigma, q_0, \delta, F)$

mit  $Q = \{0, \dots, m\}$

$q_0 = 0$   
 $F = \{m\}$

$\delta(q, a) = \overline{ov}(P_1 \dots P_m, q, a)$  für alle  $q \in Q, a \in \Sigma$

Satz 4.1: Sei  $P = P_1 \dots P_m \in \Sigma^m$ ,  $M_P = (Q, \Sigma, q_0, \delta, F)$  der String-Matching-Automat für  $p$ . Sei  $t = t_1 \dots t_n \in \Sigma^n$

Dann gilt:  $p$  ist Suffix von  $t_1 \dots t_n \iff \delta(q_0, t_1 \dots t_n) \in F$

Lemma 4.1: Sei  $\Sigma$  Alphabet, seien  $n, m \in \mathbb{N}$ , seien  $x = x_1 \dots x_n \in \Sigma^n$ ,

$y = y_1 \dots y_m \in \Sigma^m$  und  $a \in \Sigma$ . Sei  $i = \overline{ov}(x, y)$ . Dann gilt

$$\overline{ov}(x_0, y) = \overline{ov}(y_1 \dots y_i, a, y)$$

Beweis: zeige zunächst  $\overline{ov}(x_0, y) \leq \overline{ov}(x, y) + 1$

zwei Fälle: 1.  $\overline{ov}(x, y) = 0$

2.  $\overline{ov}(x, y) = i > 0 \implies y_1 \dots y_i$  Suffix von  $x_0$

$\implies y_1 \dots y_{i-1}$  Suffix von  $x$

$\implies (4.1)$

Wegen  $i = \overline{ov}(x, y)$  folgt  $x = x' y_1 \dots y_i$  für ein  $x' \in \Sigma^*$

$$\overline{ov}(x' y_1 \dots y_i, a, y) = \overline{ov}(y_1 \dots y_i, a, y)$$

□

Lemma 4.2: Sei  $p = p_1 \dots p_n \in \Sigma^n$ ,  $M_p = (Q, \Sigma, q_0, \delta, F)$  der String-Matching-Automat für  $p$ , sei  $x = x_1 \dots x_n \in \Sigma^n$ .  
 Dann gilt für alle  $i \in \{0, \dots, n\}$

$$\bar{\delta}(q_0, x_1 \dots x_i) = \overline{\text{OV}}(x_1 \dots x_i, p)$$

Beweis: Ind. über  $i$ .

$$i=0 \checkmark$$

$i \rightarrow i+1$ : Sei  $q = \bar{\delta}(q_0, x_1 \dots x_i)$ , dann gilt

$$\bar{\delta}(q_0, x_1 \dots x_{i+1}) = \bar{\delta}(q, x_{i+1})$$

$$\stackrel{\text{Def. 4.1}}{\Rightarrow} = \overline{\text{OV}}(p_1 \dots p_{i+1}, p)$$

Ind. vor:  $q = \overline{\text{OV}}(x_1 \dots x_i, p)$

$$\Rightarrow \bar{\delta}(q_0, x_1 \dots x_{i+1}) = \overline{\text{OV}}(p_1 \dots p_{i+1}, p) = \overline{\text{OV}}(x_1 \dots x_i, p) \quad \square$$

Beweis von Satz 4.1: Lemma 4.2, wenn  $F = \{m\}$   $\square$

Algorithmus 4.2: Konstruktion eines String-Matching-Automaten

Eingabe:  $p = p_1 \dots p_m \in \Sigma^m$

für  $q := 0$  bis  $m$  do

  für each  $a \in \Sigma$  do

$$\text{Berechne } \bar{\delta}(q, a) = \overline{\text{OV}}(p_1 \dots p_q, p)$$

Ausgabe:  $M_p = (Q, \Sigma, q_0, \delta, F)$

Algorithmus 4.3: String-Matching mit endl. Automaten

Eingabe: Text  $x = x_1 \dots x_n \in \Sigma^n$  und Muster  $p = p_1 \dots p_m \in \Sigma^m$

1. Bilde den String-Matching-Automaten  $M_p$  mit Alg. 4.2

2.  $q := q_0$   $I := 0$

  für  $i = 1$  bis  $n$  do

$$q := \bar{\delta}(q, x_i)$$

  if  $q \in F$  then  $I := I \cup \{i-m+1\}$

0. schritt  $\odot$  Laufzeit: Berechnung von  $M_p$ :  
naiv:  $O(|\Sigma| \cdot m^2)$   
optimal:  $O(|\Sigma| \cdot m)$

Schritt 2:  $O(n)$

gesamt:  $O(n + |\Sigma| \cdot m)$

### 4.3. Der Boyer-Moore Algorithmus

Idee: Verwende Sprünge wie im naiven Alg.

zunächst zwei Regeln die Vergleich einsparen

Bad-Character-Regel: Schiebe das Muster bis zum am weitesten rechts stehenden Vorkommen von  $t_j + x_i$  im Muster

[Vergleiche  $p_1 \dots p_m$  mit  $t_{j+1} \dots t_{j+m}$  von hinten nach vorn,  
 $p_i \neq t_{j+i}$  ist erste gefundene nicht übereinstimmende Position]

zur Implementierung der Bad-Character-Regel:

Bestimme Funktion  $\beta: \Sigma \rightarrow \{0, \dots, m\}$ , die jedem Symbol die  
jeweils letzte Position in  $p$  markiert, falls es, 0 and

Lemma 4.3. Gegeben  $p = t_{j+1} \dots t_{j+m}$  und  $p_i \neq t_{j+i} \neq a$

Dann kann das Muster für den nächsten Vergleich um

$i - \beta(a)$  Positionen nach rechts verschoben werden, ohne ein Vorkommen  
von  $p$  in  $t$  zu verpassen.

Beweis: 1)  $a$  kommt nicht in  $p$  vor:  $p$  kann an  $t_{j+1}$  vorbeigleiten, wenn

$\rightarrow$  Verschiebung von  $i - \beta(a) = i$  Pos. möglich

2)  $\beta(a) = k < i$

$p$  von  $i - k = i - \beta(a)$  Positionen verschoben

3)  $\beta(a) = k > i$ : Ignoriere Bad-Character-Regel  $\square$

Falls  $P_{i+1} \dots P_m = x_{j+1+i} \dots x_{j+m}$  und

$P_i \neq x_{j+i}$ , dann verschiebe das Muster um

$m - \sigma(i+1)$  Positionen nach rechts, wobei

$$\sigma(i) = \max \{ 0 \leq k < m \mid P_i \dots P_m \text{ ist Suffix von } P_1 \dots P_k \text{ oder } P_1 \dots P_k \text{ ist Suffix von } P_i \dots P_m \}$$

$$\sigma'(i) = \max \{ 0 \leq k < m \mid P_i \dots P_m \text{ ist Suffix von } P_1 \dots P_k \}$$

$$\sigma''(i) = \max \{ 0 \leq k < m \mid P_1 \dots P_k \text{ ist Suffix von } P_i \dots P_m \}$$

Ziel: Berechne  $\sigma'(i)$  für alle  $i \in \{2, \dots, m\}$

äquivalent: Gegeben  $P_1 \dots P_m$ , bestimme für alle Suffixe  $P_i \dots P_m$  von  $P$  für  $2 \leq i \leq m$  das letzte Vorkommen in  $P_1 \dots P_{m-1}$ .

äquivalent: Gegeben  $P_1 \dots P_m$ , bestimme für alle Präfixe  $P_m \dots P_i$  von  $P^R = P_m \dots P_1$  für  $2 \leq i \leq m$  das erste Vorkommen in  $P_{m-1} \dots P_1$ .

Lemma 4.4: Sei  $M_i = (\{0, \dots, m-i+1\}, \Sigma, 0, \delta_i, \{m-i+1\})$  der String-Matching-Automat für  $P_m \dots P_i$ , sei  $t = t_1 \dots t_n$ . Dann gilt, falls  $t_k = t_{k-m+i}$  das erste Vorkommen von  $P_m \dots P_i$  in  $t$  ist, für alle  $1 \leq j \leq k-m+i$

$$\hat{\delta}_i(q_0, t_1 \dots t_j) = \hat{\delta}_2(q_0, t_1 \dots t_j)$$

Beweis: direkt aus der Konstruktion des String-Matching-Automaten.  $\square$

Ziel: Simuliere die Berechnung aller  $M_i$  für  $2 \leq i \leq m$  durch eine Berechnung von  $M_2$ .

Blatt 7

Speichere die Folge der erreichten Zustände beim Lesen von

$P_{m-1} \dots P_1$ , sei  $q_0: P_{m-1}, \dots, P_1$  diese Folge.

Dann gilt:  $z'(i) = \max \{j \mid q_j = i\}$ , falls der Zustand  $i$  beim Lesen von  $P_{m-1}, \dots, P_1$  erreicht wurde  
sonst 0

Berechnung von  $z''(i) = \max \{0 \leq k < m \mid P_k - P_k \text{ ist Suffix von } P_i - P_m\}$ :

Wenn  $M_2$  für  $P_m \dots P_2$  nach dem Lesen von  $P_{m-1} \dots P_2$  im Zustand  $q_i = j$  endet, dann sind die letzten gelesenen Zeichen  $P_m \dots P_j$  und  $j$  ist minimal mit dieser Eigenschaft:

$\Rightarrow P_m \dots P_j$  ist Suffix von  $P_{m-1} \dots P_1$ , aber für alle  $j' < j$  ist  $P_m \dots P_{j'}$  kein Suffix von  $P_{m-1} \dots P_1$

$\Rightarrow P_j - P_m$  ist Präfix von  $P_1 - P_{m-1}$  und  $j$  ist minimal mit dieser Eigenschaft.

$\Rightarrow z''(i) = m - j + 1$  für alle  $2 \leq i \leq j$

Für alle anderen Werte von  $i$  gilt  $z''(i) = 0$ .

### Algorithmus 4.5. Preprocessing für die Good-Suffix-Regel

Eingabe Muster  $P = P_1 \dots P_m$  über  $\Sigma$

1. Konstruiere den String-Matching-Automaten

$M = (Q, \Sigma, q_0, \delta, f)$  für  $P_m - P_1$ .

2. Bestimme die Zustandsfolge  $q_0, q_{m-1}, \dots, q_1$ , die  $M$  bei der Berechnung auf  $P_{m-1}, \dots, P_1$  durchläuft.

3. for  $i := 2$  to  $m$  do  $z'(i) := 0$

for  $j := 1$  to  $m-1$  do  $z''(m-i+1) := q_j$

Biinf ③ + for  $i = 2$  to  $m$  do

if  $i \leq q$ , then

$$\sigma''(i) := m - q + 1$$

else

$$\sigma''(i) := 0$$

5. for  $i = 2$  to  $m$  do  $\sigma(i) := \max\{\sigma'(i), \sigma''(i)\}$

Ausgabe:  $\sigma$

### Algorithmus 4.6. Boyer - Moore - Algorithmus

Eingabe: Muster  $p = p_1 \dots p_m$  und Text  $t = t_1 \dots t_n$  über  $\Sigma$

1. Berechne  $B$  für die Bad-Character-Regel

2. Berechne  $\sigma$  für die Good-Suffix-Regel

3.  $I := 0$

$j := 0$

while  $j < n - m$  do

$i := m$

while  $p_i = t_{j+i}$  and  $i > 0$  do

$i := i - 1$

if  $i = 0$  then  $I := I \cup \{j\}$

$j := j + \max\{1 - B(t_{j+1}), m - \sigma(i+1)\}$

Ausgabe:  $I$

Laufzeit: worst-case:  $O(n \cdot m)$

aber: schlechterer Fall tritt sehr selten auf.



Def. 4.6. Sei  $t = t_1 \dots t_n \in \Sigma^n$  ein Text.

Ein gerichteter Baum  $T_t = (V, E)$  mit einer Wurzel  $r$  heißt **einfacher Suffix-Baum für  $t$** , falls gilt:

- (1)  $T_t$  hat genau  $n$  Blätter, die mit  $1, \dots, n$  beschriftet sind.
- (2) die Kanten von  $T_t$  sind mit Symbolen aus  $\Sigma$  beschriftet.
- (3) Alle von einem inneren Knoten ausgehenden Kanten zu seinen Kindern sind mit paarweise verschiedenen Symbolen beschriftet.
- (4) der Pfad von der Wurzel  $r$  zu dem Blatt  $i$  trägt die Beschriftung  $t_i \dots t_n$  (als Konkatenation der Kantenbeschriftungen auf dem Pfad)

Frage: Existiert für jeden Text  $t$  ein einfacher Suffix-Baum?  
Nein, existiert nicht, wenn ein Suffix von  $t$  gleichzeitig das Präfix eines anderen Suffixes ist.

Ablhilfe: Hänge an  $t$  ein neues Symbol  $\$ \notin \Sigma$  an und konstruiere den einfachen Suffix-Baum für  $t\$$ .

Algorithmus 4.7: Konstruktion eines einfachen Suffix-Baums:

Eingabe:  $t = t_1 \dots t_n \in \Sigma^n$

1.  $t' := t\$$  für  $\$ \notin \Sigma$ .

2. Initialisiere  $T_{t'}$  mit Wurzel  $r$  und leerer Kantenmenge

3. for  $i=1$  to  $n$  do

(a) Suche von  $r$  ausgehend einen Pfad in  $T_{t'}$ , der mit einem maximalen Präfix  $t_i \dots t_n$  beschriftet ist und in Knoten  $x_i$  endet.

(b) Füge den Baum eines mit  $t_{i+1} \dots t_n \$$  beschrifteten Pfades  $x_i \dots y_{i+1} \dots y_i \dots y_n$  hinzu, wobei



# Blatt 4

Lösung des String-Matching-Problems mit Suffix-Bäumen

Starte in der Wurzel und Suche Pfad in  $T_x$  mit der Beschriftung  $p$ .  
 Falls dieser existiert, ist dieser Pfad eindeutig.

→  $p$  ist Präfix eines Suffixes von  $t$

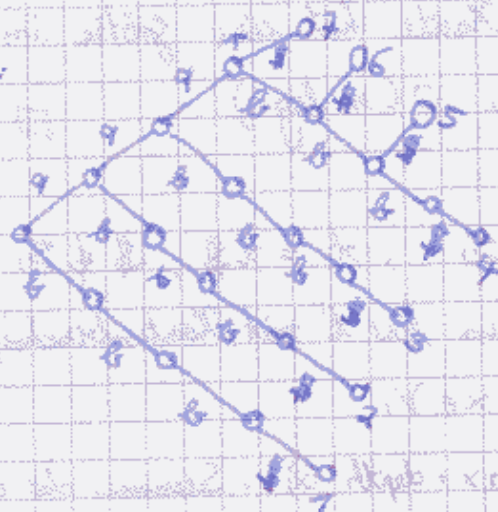
Jedes der Blätter in  $T_x$ , die unter dem Endknoten dieses Pfades hängen, entspricht einem Vorkommen von  $p$  in  $t$ .

Falls in  $T_x$  kein Pfad mit der Beschriftung  $p$  vorkommt, dann ist  $p$  kein Teilstring von  $t$ .

Problem: Das einfache Suffix-Baum für einen Text  $t = t_1 \dots t_n$  kann eine Größe in  $O(|t|^n)$  (Fu!) erreichen.

Beispiel:  $t_n = a^n b^n$

$n=3$ :



Def. 4.7. Sei  $t = t_1 \dots t_n \in \Sigma^n$  ein Text. Ein gerichteter Baum 12.5.2003  
 $T_x = (V, E)$  mit Wurzel  $r$  heißt **kompakter Suffix-Baum** für  $t$ , wenn gilt:

1. Der Baum hat genau  $n$  Blätter, die mit  $1, \dots, n$  beschriftet sind.
2. Jeder innere Knoten von  $T_x$  hat mindestens zwei Kinder.
3. Die Kanten sind mit Teilstrings von  $t$  beschriftet. Dabei wird jeder Teilstring des Stange  $k$  dargestellt durch seine Anfangs- und Endposition in  $t$ , falls  $k \geq 2 \cdot \log_2 (n - |\Sigma|)$

4. Alle Beschriftungen von einem inneren Knoten aus beginnen mit paarweise verschiedenen Symbolen.
5. Der Pfad von  $r$  zu dem Blatt  $i$  trägt die Beschriftung  $t_i \dots t_n$ .

Lemma 4.5. Sei  $t = t_1 \dots t_n \in \Sigma^n$ . Ein kompakter Suffix-Baum für  $t$  hat eine Größe in  $O(n \log n)$ .

Beweis: Jeder Suffix-Baum für  $t$  hat genau  $n$  Blätter.

Da jeder innere Knoten mindestens zwei Kinder hat, hat der Suffix-Baum höchstens  $n-1$  innere Knoten haben.

$\Rightarrow$  insgesamt  $\leq 2n-1$  Knoten

$\Rightarrow \leq 2n-2$  Kanten

Beschriftung jeder Kante hat eine Größe in  $O(\log n)$

$\Rightarrow$  Gesamtgröße  $O(n \log n)$   $\square$

Def. 4.8. Sei  $t = t_1 \dots t_n \in \Sigma^n$ ,  $t' = t \$$ ,  $\$ \notin \Sigma$ . Sei  $T_{t'} = (V, E)$  ein <sup>oder kompakter</sup> einfacher Suffix-Baum für  $t'$  mit Kantenbeschriftung  $\text{label}: E \rightarrow \Sigma^*$ .  
Für jeden Knoten  $x \in V$ :

Stringtiefe von  $x$  -  $\text{depth}(x)$ : Summe der Längen der Kantenbeschriftungen auf dem Pfad von der Wurzel zu  $x$ .

pathlabel  $(x)$ : Beschriftung auf dem Pfad von der Wurzel zu  $x$ .

Falls  $T_{t'}$  einfacher Suffix-Baum:

Postfix: Minimale Beschriftung eines Blattes in dem Teilbaum mit Wurzel  $x$ .

Biolnf ⑬ Für eine Kante  $(v, x)$  gibt  $\text{Pos}(v) + \text{depth}(v)$  die erste Position in  $t$ , an der die Beschriftung von  $(v, x)$  vorkommt.

### Algorithmus 4.8 Konstruktion eines kompakten Suffix-Baums

Eingabe: String  $t = t_1 \dots t_n \in \Sigma^n$

1.  $t' := t\$$  mit  $\$ \notin \Sigma$

Berechne einfachen Suffix-Baum  $T_{t'} = (V, E)$  mit Kantenbeschriftung  
label:  $E \rightarrow \Sigma^*$  für  $t'$ .

2. Eliminiere die Knoten vom Grad 2:

$X := \{v \in V \mid v \text{ hat genau einen Nachfolger}\}$

while  $X \neq \emptyset$  do

Wähle  $x \in X$ , sei  $y$  der Vater von  $x$ , sei  $z$  das Kind

von  $x$ . Ersetze die Kanten  $(y, x)$  und  $(x, z)$  durch

$(y, z)$  mit der Beschriftung  $\text{label}(y, z) = \text{label}(y, x) \text{label}(x, z)$

Lösche  $x$ .

3. Komprimiere lange Kantenbeschriftungen:

for all  $e = (x, y) \in E$  do

if  $|\text{label}(e)| \geq 2 \cdot \log_2(n - |\Sigma|)$  then

$\text{label}'(e) := \lfloor \text{Pos}(y) + \text{depth}(x), \text{Pos}(y) + \text{depth}(x) + |\text{label}(e)| - 1 \rfloor$

else  $\text{label}'(e) := \text{label}(e)$

Ausgabe: Der konstruierte kompakte Suffix-Baum mit Kantenbeschriftung  $\text{label}'$ .

Algorithmus 4.3: String-Matching mit kompaktem Suffix-Baum.

Eingabe: Muster:  $P = P_1 \dots P_m$ , Text  $t = t_1 \dots t_n$  über  $\Sigma$

1. Konstruiere kompakten Suffix-Baum  $T$  für  $t' = t\$$  mit Wurzel  $r$  und Kantenbeschriftung  $label$ .

```

2.  $x := r$ 
    $i := 1$ 
   gefunden := falsch
   möglich := wahr
   while gefunden = falsch und möglich = wahr do
     passende_Kante := falsch;  $U :=$  Menge der Kinder von  $x$ 
     while passende_Kante = falsch und  $U \neq \emptyset$  do
       wähle  $v \in U$ 
       if  $label(x, v) = P_i \alpha$  für  $\alpha \in \Sigma \cup \{\$ \}$  then
         passende_Kante := wahr
         label :=  $label(x, v)$ 
       else if  $label(x, v) = [k..l]$  und  $t_k = P_i$  then
         passende_Kante := wahr
          $l' := \min\{l, k + |P_i| - 1\}$ 
         label :=  $t_k \dots t_{l'}$ 
       else
          $U := U - \{v\}$ 
     if  $(P_i \dots P_m$  kein Präfix von label) und (label kein Präfix von  $P_i \dots P_m)$  then
       möglich := falsch
     else if label ist Präfix von  $P_i \dots P_m$  then
        $x := v$ 
        $i := i + |label|$ 
     else
        $x := v$ 
       gefunden := wahr
   if gefunden = wahr then

```

BiInf 15 Bestimme die Menge I der Blattbeschriftungen im dem suffixbaum mit Wurzel  $x$ .

Ausgabe: I

Satz 4.3: Algorithmus 4.3. löst das String-Matching-Problem für Text  $t = t_1 \dots t_n$  und Muster  $P = P_1 \dots P_m$  über  $\Sigma$  in einer Zeit in  $O(n \log n + m \cdot (|\Sigma| + k))$ , wobei  $k$  die Anzahl der Vorkommen von  $p$  in  $t$  und  $|\Sigma| \leq \frac{n}{\epsilon}$  gelte für eine Konstante  $c \geq 2$ .

Beweis Korrektheit:  $\checkmark$

Zeitkomplexität: Konstruktion des kompakten Suffix-Baums:  $O(n \log n)$

Bestimmung des richtigen Kindes des aktuellen Knotens  $x$ :

Kantenbeschriftung der Form  $[a \dots b]$  kodiert Teilstring der Länge  $O(\log n)$ , da  $|\Sigma| \leq \frac{n}{\epsilon}$  und damit  $2 \cdot \log_{|\Sigma|} (n - |\Sigma|) \in O(\log n)$ .

$\Rightarrow$  Beim Lesen von  $p$  höchstens  $O(\log n)$  solcher komprimierter Kantenbeschriftungen auf einem Pfad von der Wurzel. Lesen einer Kantenbeschriftung:  $O(\log n)$

Jeder aktuelle Knoten hat  $\leq |\Sigma|$  Kinder, die ggf. alle überprüft werden müssen.

Wähle die Kinder von  $x$  in solcher Reihenfolge aus, daß zunächst alle nicht-komprimierten Kanten überprüft werden.

$\Rightarrow$  nur in den Schritten, in denen auch die passende Kante eine komprimierte Beschriftung besitzt, werden komprimierte Beschriftungen gelesen.

$\Rightarrow$  Lesen der komprimierten Kantenbeschriftungen insgesamt:

$$O\left(\frac{m}{\log n} \cdot \log n \cdot |\Sigma|\right) = O(|\Sigma| \cdot m)$$

Lesen der nicht-komprimierten Kantenbeschriftungen insgesamt:

$$O(m \cdot |\Sigma|)$$

Durchsuchen des Teilraums:  $O(k)$ .

□

4.5. Weitere Anwendungen von Suffiz-Bäumen

## 4.5.1. Verallgemeinerte Suffiz-Bäume und das Teilstring-Problem

Def. 4.3: Das Teilstring-Problem ist das folgende Berechnungsproblem:Eingabe: Ein Muster  $P$  und  $N$  Texte  $t_1, \dots, t_N$  über Alphabet  $\Sigma$ Ausgabe: Eine Menge  $I \subseteq \{1, \dots, N\}$ , so daß  $i \in I$  gdw.  $P$  ein Teilstring von  $t_i$ .Ziel: Verwende Suffiz-Baum, der die Suffiz von  $t_1, \dots, t_N$  enthält.Konstruiere Suffiz-Baum für  $t_1 \$_1 t_2 \$_2 \dots t_N \$_N$ , wobei  $\$_i \notin \Sigma$ .

4.5.203

Lemma 4.6: Seien  $t_1, \dots, t_N$  über  $\Sigma$  gegeben. Sei  $T$  der kompakte (nicht komprimierte) Suffiz-Baum für  $t'_1 = t_1 \$_1 t_2 \$_2 \dots t_N \$_N$ , wobei  $\$_1, \dots, \$_N \notin \Sigma$ ,  $\$_1, \dots, \$_N$  paarweise verschieden.Dann treten die Transymbole  $\$_i$ ,  $1 \leq i \leq N$ , in  $T$  nur in den Beschriftungen von Knoten auf, die incident zu dem Blättern sind.Beweis: Annahme:  $\$_i$  liegt auf einer Kante zwischen dem inneren Knoten  $x$  und  $y$ , wobei  $x$  der Vater von  $y$  sei. $\Rightarrow$  Es ex. zwei verschiedene Suffiz  $u \$_i v$  und  $w \$_i z$ , da  $y$  mindestens zwei Kinder hat.Widerspruch, da  $\$_i$  genau einmal in  $t'$  vorkommt. □



Def. 4.10: Seien  $t_1, \dots, t_N$  über  $\Sigma$  gegeben, seien  $\$_1, \dots, \$_N \notin \Sigma$  paarweise verschiedene Trennsymbole. Ein verallgemeinerter Suffix-Baum für  $t_1, \dots, t_N$  entsteht aus einem kompakten, nicht komprimierten Suffix-Baum für  $t' = t_1 \$_1 \dots t_N \$_N$  durch die folgenden Schritte:

1. Ersetze jede Kantenbeschriftung der Form  $u \$_j w$ ,  $w \in (\Sigma \cup \{\$_j \mid 1 \leq j \leq N\})^*$  durch  $u \$_j$ .
2. Beschrifte jedes Blatt mit einem Paar  $(i, j)$  aus dem Index des zugehörigen Textes  $t_i$  und der Startposition  $j$  des zugehörigen Suffixes.
  - i. Index des Trennsymbols auf der inzidenten Kante.
  - j. an der Blattbeschriftung des kompakten Suffix-Baums für  $t$  und den Längen des Textes.
3. Komprimiere lange Kantenbeschriftungen.

Satz 4.1: Seien ein Muster  $p = p_1 \dots p_m$  und  $N$  Texte  $t_1, \dots, t_N$  der Gesamtlänge  $n$  über  $\Sigma$  gegeben.

Dann lässt sich das Teilstring-Problem in linear Zeit in  $O(n \log n + m \cdot (|\Sigma| + k))$  lösen, wobei  $k$  die Anzahl der Vorkommen von  $p$  in  $t_1, \dots, t_N$  sei, und  $|\Sigma| \leq \frac{n}{2}$  für ein  $c \geq 2$ .

Beweis: analog zu Satz 4.3:

Konstruktion des verallgemeinerten Suffix-Baums:  $O(n \log n)$

Suchen des Musters  $p$ :  $O(m \cdot |\Sigma|)$

Durchsuchen des Teilbaums:  $O(k)$

□

## Bioinf 8 4.5.2. Längste gemeinsame Teilstring

Def. 4.11: Das Längst-Gemein-Teilstring-Problem ist das folgende Optimierungsproblem:

Eingabe: Menge  $M = \{t_1, \dots, t_N\}$  von Strings über  $\Sigma$

Zulässige Lösungen: Jeder String  $t$ , der Teilstring von  $t_i$  ist für alle  $1 \leq i \leq N$

Kosten: Für zulässige Lösung  $t$ :  $\text{cost}(t) = |t|$

Optimierungsziel: Maximierung

Algorithmus 4.10: Bestimmung des längsten gemeinsamen Teilstrings

Eingabe: Eine Folge  $(t_1, \dots, t_N)$  von Strings über  $\Sigma$

1. Konstruiere verallgemeinerten Suffix-Baum  $T$  für  $t_1, \dots, t_N$

2. Beschrifte jeden inneren Knoten  $x$  mit einer Menge

$M(x) \subseteq \{1, \dots, N\}$ , so daß  $i \in M(x)$  gdw. in dem Teilbaum mit Wurzel  $x$  ein Blatt mit Beschriftung  $(i, j)$  ex. für bel.  $j$

3. Finde unter allen inneren Knoten von  $T$  mit der Beschriftung  $\{1, \dots, N\}$  einen Knoten  $x_{\max}$  maximaler Stringtiefe

4.  $d_{\max} := \text{pathlabel}(x_{\max})$

Ausgabe:  $d_{\max}$

Satz 4.5: Alg. 4.10. bestimmt den längsten gemeinsamen Teilstring von  $N$  Strings  $t_1, \dots, t_N$  der Gesamtlänge  $n$  in einer Zeit in  $O(n(\log n + N^2 \cdot |\Sigma|))$

Beweis: Korrektheit:  $\checkmark$

Zeitkomplexität: Konstruktion des verallgemeinerten Suffix-Baums  
 $O(n \cdot \log n)$

Bestimmung der Mengen  $M(x)$ : bottom-up, von den Blättern aus:

Punkt 13) Beschriftung für Knoten  $x$  ergibt sich als Vereinigung der Mengen seiner Kinder. Jeder Knoten hat höchstens  $|\Sigma| + N$  Kinder, falls er weder eine Blatte ist,  $|\Sigma|$  Kinder sonst, es gibt  $O(n)$  Knoten  
~~ergibt sich durch Induktion~~  
 $\Rightarrow$  Gesamtaufwand:  $O(n \cdot (|\Sigma| + N) \cdot N)$

Für jeden Knoten  $x$  prüfen, ob  $L(x) \subseteq \{1, \dots, N\}$  oder nicht.

$\rightarrow$  Für jeden Knoten ist die Überprüfung der Bedingung in konstanter Zeit möglich.

$\Rightarrow$  Finden aller Knoten mit Beschriftung  $\{1, \dots, N\}$  in  $O(n \cdot \log n)$

Bestimmung von  $x_{max}$  und  $d_{max}$ :  $O(n \cdot \log n)$

### 4.5.3. Effiziente Bestimmung von Overlaps

Aufgabe: Bestimmung aller paarweisen Overlaps von  $N$  Strings

Laufzeit: naiv (falls alle Strings die Länge  $\ell$  haben):

$$O(N^2 \cdot \ell^2) = O(n^2)$$

Ziel: mit verallgemeinertem Suffix-Baum:  $O(n \cdot (\log n + |\Sigma| + N))$

Algorithmus 4.11: Bestimmung aller paarweisen Overlaps

Eingabe: Eine Folge  $t_1, \dots, t_N$  von  $N$  Strings über  $\Sigma$

1. Konstruiere verallgemeinerten Suffix-Baum  $T$  für  $t_1, \dots, t_N$  mit Wurzel  $r$

2. Beschrifte jeden inneren Knoten  $x$  mit  $L(x) \subseteq \{1, \dots, N\}$ , so daß

$i \in L(x)$  gdw.  $x$  insidert ist zu einer Kante mit Beschriftung  $t_i$

3. for  $j := 1$  to  $N$  do

$x := T$

while  $x$  ist kein Blatt do

for all  $i \in L(x)$  do

if  $\text{depth}(x) < \min\{|t_i|, |t_j|\}$  then

$u(t_i, t_j) := \text{pathlabel}(x)$

$w(t_i, t_j) := \text{depth}(x)$

$x := \text{Kind von } x \text{ auf dem Pfad mit Beschriftung } t_j|_j$

Ausgabe: Die Overlaps  $O_i(t_i, t_j) = u(t_i, t_j)$  und ihren Längen

$w_i(t_i, t_j) = w(t_i, t_j)$  for alle  $1 \leq i, j \leq N$

Satz 4.6: Seien  $N$  Strings  $t_1, \dots, t_N$  der Gesamtlänge  $n$  über  $\Sigma$  gegeben.

Dann bestimmt Alg. 4.11 alle paarweisen Overlaps in einer

Zeit in  $O(n \cdot (\log n + |\Sigma| + N))$

Beweis: Korrektheit: Falls  $i \in L(x)$  für einen Knoten  $x$ , dann

ist  $\text{pathlabel}(x)$  ein Suffix von  $t_i$ .

Gleichzeitig ist  $\text{pathlabel}(x)$  ein Präfix von  $t_j$ , falls  $x$  auf dem Pfad mit Beschriftung  $t_j|_j$  liegt.

$\Rightarrow \text{pathlabel}(x)$  ist Überlappung von  $t_i$  und  $t_j$

Alle möglichen Überlappungen von  $t_i$  und  $t_j$  kommen vor als  $\text{pathlabel}(x)$  für einen Knoten  $x$ .

$\Rightarrow O_i(t_i, t_j) = \text{pathlabel}(x)$  für den tiefsten Knoten  $x$  auf dem Pfad mit Beschriftung  $t_j|_j$ , dessen  $i$ -te Menge den Index  $i$  enthält.

Bedingung  $\text{depth}(x) < \min\{|t_i|, |t_j|\}$  stellt sicher, daß keine Überlappung berücksichtigt wird, bei der  $t_i$  Suffix von  $t_j$  ist

## Lösung 21

oder umgekehrt.

Zeitkomplexität: Konstruktion des Suffix-Baums:  $O(n \cdot \log n)$

Bestimmung von  $L(x)$  bzw.  $O(n \cdot N)$

Gesamtlänge aller untersuchten Pfade:  $O(n)$

Summe der Mächtigkeiten der dabei untersuchten  $L$ -Mengen:  $O(n)$

Berechnung der Pfade insgesamt  $O(n \cdot (|\Sigma| + N))$

$\Rightarrow$  Insgesamt:  $O(n \cdot (\log n + |\Sigma| + N))$   $\square$

18.5.2008

## 5. Alignment-Verfahren

A	T	T	T
C	C	C	C

### 5.1. Alignment von zwei Strings

#### Beispiel 5.1

$s = \text{GACGGATTATG}$

$t = \text{GATCGGAATAG}$

$s' = \text{GA} - \text{CGG} \text{ATTATG}$

$t' = \text{GATCGGA} \text{ATA} - \text{G}$

Def. 5.1: Seien  $s = s_1 \dots s_m$ ,  $t = t_1 \dots t_n$  zwei Strings über  $\Sigma$ . Sei  $- \notin \Sigma$  ein spezielles **Lückensymbol** und  $\Sigma' = \Sigma \cup \{-$ .

$h: (\Sigma')^* \rightarrow \Sigma^*$ ,  $h(a) = a \forall a \in \Sigma$   
 $h(-) = \lambda$

Ein **Alignment** von  $s$  und  $t$  der Länge  $l = \max\{m, n\}$  über  $\Sigma'$  ist ein Paar  $(s', t')$  für die gilt:

i)  $|s'| = |t'| = l \geq \max\{m, n\}$

ii)  $h(s') = s$  und  $h(t') = t$

iii)  $-$  ist keine Position an der sowohl in  $s'$  als auch in  $t'$

Point (a) eine Lücke auftritt

$\forall i \in \{1, \dots, l\} \quad t_i' \neq - \text{ oder } s_i' \neq -$

4. Typen von Spalten in einem Alignment

1.) **Einfügung / Insertion:** Der erste (obere) String besitzt in dieser Spalte eine Lücke.

2.) **Löschung / Deletion:** Der zweite (untere) String besitzt in dieser Spalte eine Lücke.

3.) **Match:** Die Buchstaben in dieser Spalte des Alignment sind identisch

4.) **Mismatch:** Die Buchstaben in dieser Spalte des Alignment sind nicht identisch.

### Bewertung eines Alignment

Def. 5.2: Zwei Strings  $s, t$  über  $\Sigma$

Sei  $p(a, b) \in \mathbb{Q} \quad \forall a, b \in \Sigma$

$g \in \mathbb{Q}$

Definieren: die **Bewertung eines Alignment** von  $s$  und  $t$  sei spaltenweise definiert:

$$\forall x, y \in \Sigma: \delta(x, y) = p(x, y)$$

$$\delta(x, -) = \delta(-, y) = g$$

$s', t'$ ,  $(s', t')$  Alignment von  $s$  und  $t$

$$\delta(s', t') = \sum_{i=1}^l \delta(s_i', t_i')$$

Außerdem sei durch  $\text{goal}_g \in \{\min, \max\}$  ein Optimierungsziel gegeben. Im allgemeinen sollte  $p(a, b) = p(b, a)$

Brotf (23) Def. 5.3: Seien  $s, t \in \Sigma$ ,  $J$  eine Alignment-Bewertung,

goal  $J$  ein Optimierungsziel für  $J$ .

Die Ähnlichkeit (Similarity)  $\text{sim}_J(s, t)$  von  $s$  und  $t$  ist die Bewertung eines optimalen Alignment  $(s', t')$  von  $s$  und  $t$ , d.h.

$$\text{sim}_J(s, t) = \text{goal}_J \{ J(s', t') \mid (s', t') \text{ Alignment von } s \text{ und } t \}$$

### übliche Bewertungsfunktionen

Edit-Distanz / Levenshtein-Distanz

$$P(a, b) = 1 \quad a \neq b$$

$$P(a, a) = 0$$

$$g = 1$$

goal = min

anderes Maß:

$$P(a, a) = 1$$

$$P(a, b) = -1 \quad a \neq b$$

$$g = -2$$

goal = max

### 5.12 Globale Alignment

„Ähnlichkeit der gesamten Strings“

formal: Def. 5.4 Global-Alignment-Problem

Eingabe:  $s, t$  über  $\Sigma$ ,  $J$  eine Bewertungsfunktion,  
goal  $J$  Optimierungsziel für  $J$ .

Zulässige Lösungen: Alle Alignment von  $t$  und  $s$ .

Kosten:  $\forall$  Alignment  $A = (s', t')$

$$\text{cost}(A) = J(A)$$

Ziel: goal  $J$  — im folgenden Maximierung

Point (24)

ATGA - - - - -16 8  
- - - - ATGA

ATGA  
ATGA

4

0

└──┘

└──┘

g = -2

g = 1

Match = 1

Match = 0

Mismatch = -1

Mismatch = 1

Max.

Min.

### dynamische Programmierung:

- Zerlegung eines Problems in Teilprobleme
  - Setzen die Lösungen der Teilprobleme bottom-up zusammen.
  - Lösungen für Teilprobleme werden zwischengespeichert
- (Beispiel: Fibonacci - Zahlen)

### dynamische Programmierung für das Global-Alignment-Problem

$$S = s_1 \dots s_m$$

$$x = x_1 \dots x_n$$

$(m+1) \times (n+1)$  Präfixpaare

→ Matrix  $M$  aus  $(m+1) \times (n+1)$  Einträgen

$M(i, j)$ : ist die Bewertung des optimalen Alignment für  $s_1 \dots s_i$  und  $x_1 \dots x_j$

$M(m, n)$ : Bewertung des optimalen Alignment für  $s, x$

$M$ : Ähnlichkeitsmatrix

$$s_1 \dots s_i \quad \begin{pmatrix} s_1 \dots s_i \\ \text{---} \end{pmatrix} : \quad M(i, 0) = i \cdot g$$

$$x_1 \dots x_j \quad \begin{pmatrix} \text{---} \\ x_1 \dots x_j \end{pmatrix} : \quad M(0, j) = j \cdot g$$



Initialisierung von M:

Initialisiere die erste Spalte und erste Zeile von M

$$\text{mit } M(i, 0) = i \cdot g \quad \forall i \in \{1, \dots, m\}$$

$$M(0, j) = j \cdot g \quad \forall j \in \{1, \dots, n\}$$

Skizze: 

Rekurrenz:

$$M(i, j) = \begin{cases} M(i-1, j) + g & \text{Einfügung} \\ M(i, j-1) + g & \text{Löschen} \\ M(i-1, j-1) + P(s_i, t_j) & \text{Match/Mismatch} \end{cases}$$

Algorithmus 5.1 - Bestimmung der Ähnlichkeit

Eingabe:  $s, t$ ,  $s = s_1 \dots s_m$ ,  $t = t_1 \dots t_n$

1) Initialisierung

for  $i = 0$  to  $m$  do

for  $j = 0$  to  $n$  do

$$M(i, j) := 0$$

2) Initialisierung der Ränder:

for  $i = 0$  to  $m$  do

$$M(i, 0) := i \cdot g$$

for  $j = 0$  to  $n$  do

$$M(0, j) := j \cdot g$$

3.) Auffüllen der Matrix:

for  $i=1$  to  $m$  do

for  $j=1$  to  $n$  do

$$M(i, j) := \max \left\{ \begin{array}{l} M(i-1, j) + g, \\ M(i, j-1) + g, \\ M(i-1, j-1) + p(s_i, t_j) \end{array} \right\}$$

Ausgabe:  $\text{sim}(s, t) = M(m, n)$

Satz 5.1: Seien zwei Strings  $s = s_1 \dots s_m$  und  $t = t_1 \dots t_n$  über  $\Sigma$  gegeben. Dann lässt sich mit dem Algorithmus 5.1 (Berechnung der Ähnlichkeit) und 5.2 (Trace-Back, siehe Folie) ein optimales Alignment von  $s$  und  $t$  in Zeit  $O(n \cdot m)$  bestimmen.

Beweis: Alg. 5.1: Auffüllen der Matrix  $O(n \cdot m)$

Alg. 5.2: Durchlauf durch  $M$ :  $O(n+m)$

Insgesamt:  $O(n \cdot m)$

Berechnung aller möglichen optimalen Alignments

$\Rightarrow$  es können exponentiell viele optimale Alignments existieren.

z.B.  $s = A^{2n}$ ,  $t = A^n$

$$\begin{array}{cccccc} A & A & A & A & A & A \\ - & A & A & A & - & - \end{array} \quad \begin{array}{cccccc} A & A & A & A & A & A \\ A & - & - & - & A & A \end{array}$$

allg.  $\binom{2n}{n}$  exponentiell viele

graphentheoretisches Problem:

Def. 5.5:  $s = s_1 \dots s_m$ ,  $t = t_1 \dots t_n$  Strings über  $\Sigma$   
 $\mathcal{J}$  Alignmentbewertung mit Parametern  $P$  und  $g$ .

Der Edit-Graph von  $s$  und  $t$  bezgl.  $\mathcal{J}$  definiert durch

$$G_{\mathcal{J}}(s, t) = (V, E, c)$$

$$V := \{0, \dots, m\} \times \{0, \dots, n\}$$

$$E := \{ (i, j), (i, j+1) \mid \forall i, j \}$$

$$\cup \{ (i, j), (i+1, j) \mid \forall i, j \}$$

$$\cup \{ (i, j), (i+1, j+1) \mid \forall i, j \}$$

$$c: E \rightarrow \mathbb{Q}$$

$$c((i, j), (i, j+1)) = c((i, j), (i+1, j)) = g$$

$$c((i, j), (i+1, j+1)) = P(s_{i+1}, t_{j+1}) \quad \forall i, j$$

Lösung des global-Alignment-Problem auf Basis des Edit-Graphen  
 - finde den längsten Weg von dem Knoten  $(0, 0)$  zum Knoten  $(m, n)$ .  
 (Polynomielle Laufzeit)

21.5.2023

5.1.3. Lokales und Semiglobales Alignment

Lokales Alignment

Def. 5.6: Seien zwei Strings  $s = s_1 \dots s_m$  und  $t = t_1 \dots t_n$  über  $\Sigma$  gegeben,  
 und sei  $\mathcal{J}$  eine Alignment-Bewertung mit Optimierungsziel  
 Maximierung.

Ein lokales Alignment von  $s$  und  $t$  ist ein (globales) Alignment von  
 Teilstrings  $\bar{s} = s_{i_1} \dots s_{i_2}$  und  $\bar{t} = t_{j_1} \dots t_{j_2}$ .

Ein Alignment  $A = (\bar{s}, \bar{t})$  von  $\bar{s}$  und  $\bar{t}$  ist ein optimales lokales Alignment falls gilt:

$$J(A) = \max \{ \text{sim}(S, T) \mid S \text{ ist Teilstring von } s, T \text{ ist Teilstring von } t \}$$

Def. 5.7: das Local-Alignment-Problem ist das folgende Optimierungsproblem:

Eingabe: Zwei Strings  $s$  und  $t$  über  $\Sigma$  und Alignment-Bewertung  $f$  mit Optimierungsziel Maximierung.

Zulässige Lösungen: Alle lokalen Alignments von  $s$  und  $t$ , d.h. alle globalen Alignments für alle möglichen Teilstrings  $S$  von  $s$  und  $T$  von  $t$ .

Kosten: Für ein lokales Alignment  $A = (S', T')$  gilt  $\text{cost}(A) = J(A)$

Optimierungsziel: Maximierung

Beispiel:  
 $s = \text{AAAAA CTCTCTCT}$   
 $t = \text{GCGCGCGCAAAAA}$

opt. lokales Alignment:

	$s$	
	AAAAA	(CTCTCTCT)
(GCGCGCGC)		AAAAA
		$t$

Berechnung eines opt. lokalen Alignments: dynamische Programmierung  
 Berechne  $(n+1) \times (m+1)$ -Matrix  $M$

$M(i, j) =$  höchste Bewertung eines Alignments zwischen einem Suffix von  $s_1 \dots s_i$  und einem Suffix von  $t_1 \dots t_j$ .

$$M(i, j) = \max \begin{cases} M(i-1, j) + g \\ M(i, j-1) + g \\ M(i-1, j-1) + f(s_i, t_j) \\ 0 \end{cases}$$

Initialisierung: Zeile 0 und Spalte 0 mit Nullen initialisieren

Bewertung eines opt. lokalen Alignments: höchster in  $M$  vorkommender Wert

Bestimmung des opt. lokalen Alignments: Traceback von höchstem Wert in  $M$

# Semiglobale Alignments

Beispiel 5.6:

$s = \text{ACTTATGCCTGCT}$

$t = \text{ACAGGCT}$

$d$	Match	+1
	Mismatch	-1
	Deletion	-2

opt. globales Alignment:

ACTTATGCCTGCT  
AC--A-G---GCT

Problem:  $t$  wird stark verstückelt

opt. lokales Alignment

ACTTATGCCTGCT  
(ACAG)GCT

„besseres Alignment“:

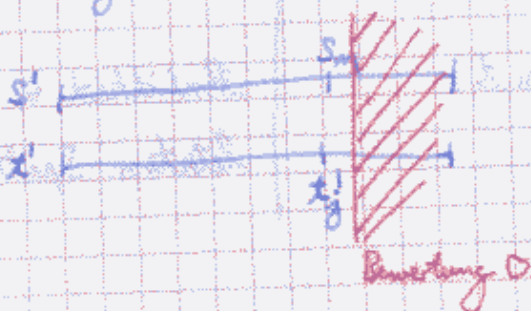
ACTTAT- GCCTGCT  
-----ACAGGCT---

Lücken vor  $t_1$  und hinter  $t_n$  ignorieren: Bewertung 0 (Statt -13 bei glob. Alignment)

1. Kostenlose Lücken am Ende von  $s$ :

Sei  $(s', t')$  ein Alignment von  $s$  und  $t$  der Länge  $l$ , bei dem

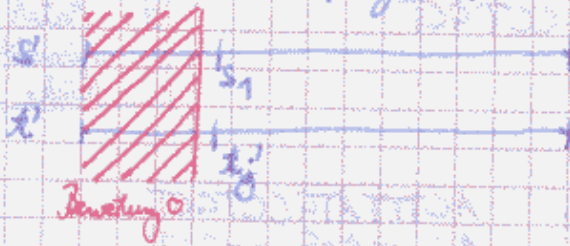
$t_{j'} = s_m$  gilt für ein  $j' \leq l$ .



⇒ Bewertung eines globalen Alignments von  $s$  mit dem Präfix  $t_1 \dots t_{j'}$  von  $t$ .

⇒ In der Alignment-Matrix werden in der letzten Zeile alle Bewertungen von Alignments von  $s$  mit allen Präfixen von  $t$  berechnet → Suche Maximum der letzten Zeile.

2. Kostenlose Lücken am Anfang von S:



- ⇒ Berechne opt. globales Alignment von S mit Suffix von T
- ⇒ Initialisieren die erste Zeile der Matrix mit Nullen

Kostenlose Lückensymbole Änderung des Algorithmus 5.1

am Anfang von S Initialisierung der ersten Zeile von M mit Nullen

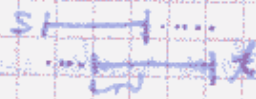
am Ende von S Ähnlichkeit  $\pm$  Maximum der letzten Zeile

am Anfang von T Initialisierung der ersten Spalte von M mit Nullen

am Ende von T Ähnlichkeit  $\pm$  Maximum der letzten Spalte

Beispiel: Bestimmung von Näherungswerten overlaps

→ mache Lücken am Ende von S und am Anfang von T kostenlos:



Bewertung von Lücken:

Def. 5.8: Seien  $s$  und  $t$  zwei Strings, sei  $(s', t')$  ein Alignment von  $S$  nach  $T$ .

Einen Teilstring  $s'_i, \dots, s'_{i+k} = -^k$  mit  $s'_i, s'_{i+k} \neq -$

bzw. einen Teilstring  $t'_j, \dots, t'_{j+k} = -^k$  mit  $t'_j, t'_{j+k} \neq -$

bezeichnen wir als Lücke der Länge  $k$ .

Beispiel: Bewertung einer Lücke der Länge  $k$ :  $k \cdot g$

offene Lückenbewertung: Bewertung einer Lücke der Länge  $k$  mit

$$-(g + \sigma \cdot k), \quad g, \sigma > 0$$

$\Rightarrow$  Anteil  $\sigma \cdot k$  proportional zur Länge der Lücke  
 zusätzlich Kosten für das Öffnen einer Lücke

Berechnung eines opt. Alignments mit offener Lückenbewertung: siehe Übung

Bewertungsmatrizen:

Motivation: Vergleich von Protein-Sequenzen.

Substitution zwischen bestimmten Aminosäuren sind wahrscheinlicher als zwischen anderen.

Non-triviale  $20 \times 20$ -Matrix von Werten  $p(a, b)$ : Bewertungsmatrix

PAM - Matrizen:

PAM = „percent of accepted mutations“

Def. 5.9: Eine akzeptierte Mutation ist eine Mutation, die die Funktionsweise des Proteins nicht oder nur so wenig ändert, daß sie wieder revertiert werden kann.

Zwei Protein-Sequenzen  $s$  und  $t$  sind eine PAM-Einheit voneinander entfernt, wenn  $s$  in  $t$  überführt wurde durch eine Folge von akzeptierten Punkt-Mutationen (Substitutionen einzelner Aminosäuren, keine Einfügungen oder Lösungen), so daß durchschnittlich eine akzeptierte Punkt-Mutation pro hundert Aminosäuren auftritt.

Beachte: Zwei Protein-Sequenzen, die  $k$  PAM-Einheiten entfernt sind, unterscheiden sich nicht notwendig an  $k$  Positionen der Stellen, da mehrere Mutationen an der gleichen Stelle auftreten können.

Idee: Konstruiere  $k$ -PAM-Matrix für den Vergleich von Protein-Sequenzen, die  $k$  PAM-Einheiten voneinander entfernt sind.

Bestimmung einer  $k$ -PAM-Matrix im idealen Fall:

Annahme 1- Wir kennen viele Paare von homologen Protein-Sequenzen von denen wir wissen, daß sie  $k$  PAM-Einheiten voneinander entfernt sind.

2- Wir kennen für jedes Paar das opt. Alignment, also die Positionen der Lücken.

Sei  $A$  die Menge aller Alignments der geg. Sequenz-Paare, sei  $Sp(A)$  die Multimenge aller Spalten in  $A$ .



$$\text{freq}(a_i, a_j) = \frac{\text{Anzahl der Vorkommen von } (a_i, a_j) \text{ und } (a_j, a_i) \text{ in SP(A)}}{2 \cdot |S(A)|}$$

$$\text{freq}(a_i) = \frac{\text{Anzahl der Vorkommen von } a_i \text{ in allen Sequenzen}}{\text{Gesamtlänge aller Sequenzen}}$$

$$\text{PAM}_k(i, j) = \log \frac{\text{freq}(a_i, a_j)}{\text{freq}(a_i) \cdot \text{freq}(a_j)}$$

Anschauliche Bedeutung:  $\text{PAM}_k(i, j)$  beschreibt das Verhältnis zwischen der W'keit, sich die  $a_i$  durch akzeptierte Mutationen in  $a_j$  umgewandelt wird, und der W'keit, dass dieses Paar zufällig in einem Alignment auftritt.

Praktischer Ansatz:

Wähle Menge sehr ähnlicher Sequenzen, die von gemeinsamen Vorfahren abstammen und nur eine PAM-Einheit voneinander entfernt sind.

→ diese Sequenzen haben ungefähr dieselbe Länge

→ Anordnung der Buchst. ist einfach zu bestimmen

⇒ Bestimme 1-PAM-Matrix

2. Verwende 1-PAM-Matrix zur Berechnung von  $k$ -PAM-Matrizen:

Sei  $F$  eine  $20 \times 20$ -Matrix, so daß  $F(i, j)$  die W'keit angibt,

daß  $a_i$  in einer PAM-Einheit zu  $a_j$  mutiert (unabhängig von der

Kontext der Buchst. von  $a_i$ ).

⇒  $F^k(i, j)$  gibt die W'keit an, daß  $a_i$  in  $k$ -Einheiten zu  $a_j$  mutiert.

$$\Rightarrow \text{PAM}_k(i, j) = \log \frac{\text{freq}(a_i, a_j) \cdot F^k(i, j)}{\text{freq}(a_i) \cdot \text{freq}(a_j)} = \log \frac{F^k(i, j)}{\text{freq}(a_i)}$$

In der Praxis: Beim Vergleich von zwei Protein-Sequenzen ist  $k$  nicht bekannt.

⇒ Standardwerte:  $k=40, k=100, k=250$

## 5.3.1. Definition und Bewertung von multiplen Alignments

Def. 5.10: Seien  $k$  Strings  $s_1 = s_{11} \dots s_{1m_1}, \dots, s_k = s_{k1} \dots s_{km_k}$  über  $\Sigma$  gegeben. Sei  $- \in \Sigma$  Lückensymbol,  $\Sigma' = \Sigma \cup \{-\}$ .

Sei  $h: (\Sigma')^* \rightarrow \Sigma^*$  ein Homomorphismus, definiert durch  $h(a) = a$  für alle  $a \in \Sigma$ ,  $h(-) = \epsilon$ .

Ein **multipler Alignment** von  $s_1, \dots, s_k$  ist ein Tupel  $(s_1', \dots, s_k')$  von Strings der Länge  $l \geq \max\{m_i \mid 1 \leq i \leq k\}$  über  $\Sigma'$ , so daß gilt:

$$(a) |s_1'| = |s_2'| = \dots = |s_k'| = l$$

$$(b) h(s_i') = s_i \text{ für alle } 1 \leq i \leq k$$

(c) es gibt keine Position, an der in allen  $s_i'$  ein Lücke vorkommt.

Def. 5.11: Gegeben sei ein multipler Alignment  $(s_1', \dots, s_k')$  der Länge  $l$ . Ein String  $c = c_1 \dots c_l \in \Sigma^l$  heißt **Consensus** für  $(s_1', \dots, s_k')$ , falls gilt

$$c_j = \operatorname{argmax}_{a \in \Sigma} |\{s_{ij}' = a \mid 1 \leq i \leq k\}| \text{ für alle } 1 \leq j \leq l$$

Der **Abstand** eines Alignment  $(s_1', \dots, s_k')$  von einem Consensus  $c$  ist definiert als

$$\operatorname{dist}(c, (s_1', \dots, s_k')) = \sum_{j=1}^l |\{s_{ij}' \mid s_{ij}' \neq c_j, 1 \leq i \leq k\}|$$

Lemma 5.1: Sei  $(s_1, \dots, s_k)$  ein multipler Alignment und seien  $c$  und  $\bar{c}$  zwei Consensus-Strings für  $(s_1, \dots, s_k)$ . Dann gilt

$$\text{dist}(c, (s_1, \dots, s_k)) = \text{dist}(\bar{c}, (s_1, \dots, s_k))$$

Beweis: Betrachte bel. Spalte  $j$  des Alignments. Falls  $c_j \neq \bar{c}_j$ , müssen in Spalte  $j$  die Symbole  $c_j$  und  $\bar{c}_j$  gleich häufig vorkommen.

$$\Rightarrow |\{s_{ij} \mid s_{ij} \neq c_j, 1 \leq i \leq k\}| = |\{s_{ij} \mid s_{ij} \neq \bar{c}_j, 1 \leq i \leq k\}| \quad \square$$

$\Rightarrow$  Abstand zum Consensus  $\hat{=}$  Abstand zu beliebigem Consensus

Beispiel 5.3:

$$\begin{array}{r} s_1 = \text{AATGCT} \\ s_2 = \text{A-TTC-} \\ s_3 = \text{- - -TCC} \\ \hline c = \text{AATTCT} \end{array}$$

Abstand  $1+2+1+1+0+2 = 7$

Def. 5.1: Das Multi-Consensus-Align-Problem ist das folgende Optimierungsproblem:

Eingabe: Eine Menge  $S = \{s_1, \dots, s_k\}$  von Strings über einem Alphabet  $\Sigma$ .

zulässige Lösungen: alle multiplen Alignments von  $S$ .

Kosten: Die Kosten eines multiplen Alignments  $(s_1, \dots, s_k)$  mit Consensus  $c$  sind:

$$\text{cost}(s_1, \dots, s_k) = \text{dist}(c, (s_1, \dots, s_k))$$

Optimierungsziel: Minimierung

Bioinf (36) Def. 5.13: Sei  $\Sigma$  ein Alphabet,  $- \notin \Sigma$  Lückensymbol und  $d$  Bewertungsfunktion für das Alignment von zwei Strings mit Optimierungsziel Minimierung, erweitert um geeignete Definition für  $d(-, -)$ .

Die Bewertung  $J_{SP}$  eines multiplen Alignment  $(s_1, \dots, s_k)$  der Länge  $l$  als Summe der Paare (SP-Bewertung) für eine Spalte:

$$J_{SP}(s_{1j}, \dots, s_{kj}) = \sum_{i=1}^k \sum_{r=i+1}^k d(s_{ij}, s_{rj})$$

$$J_{SP}(s_1, \dots, s_k) = \sum_{j=1}^l J_{SP}(s_{1j}, \dots, s_{kj})$$

Beispiel 5.10:  $s_1 = AATGCT$

$s_2 = A-TTC-$

$s_3 = - - -TCC$

$d$ : Edit-Distanz, d.h.  $d(a, b) = 0$ , falls  $a = b$ ,  $d(a, b) = 1$  sonst

$$\begin{aligned} J_{SP}(s_1, s_2, s_3) &= \sum_{j=1}^6 \sum_{i=1}^3 \sum_{r=i+1}^3 d(s_{ij}, s_{rj}) \\ &= \sum_{j=1}^6 (d(s_{1j}, s_{2j}) + d(s_{1j}, s_{3j}) + d(s_{2j}, s_{3j})) \\ &= 2 + 2 + 2 + 2 + 0 + 3 = 11. \end{aligned}$$

Def. 5.14: Das Mult-SP-Align.-Problem ist das folgende Optimierungsproblem:

Eingabe: Eine Menge  $S = \{s_1, \dots, s_k\}$  von Strings über  $\Sigma$  und eine Alignment-Bewertung  $d$  mit Ziel Minimierung.

Zulässige Lösungen: Alle multiplen Alignment von  $S$

Kosten: Die Kosten eines multiplen Alignment  $(s_1, \dots, s_k)$  sind:

$$\text{cost}(s_1, \dots, s_k) = d_{sp}(s_1, \dots, s_k)$$

Optimierungsziel: Minimierung

### 5.3.2. Exakte Bestimmung multipler Alignments

Ansatz: dynamische Programmierung:

Für  $k$  Strings verwende  $k$ -dimensionales Array  $M$ , dessen Einträge  $M(i_1, \dots, i_k)$  die Bewertung eines opt. Alignments der Präfixe  $s_{11} \dots s_{1i_1}, s_{21} \dots s_{2i_2}, \dots, s_{k1} \dots s_{ki_k}$  enthält.

- Probleme:
- Aufwändige Datenstruktur, falls  $k$  nicht vorher bekannt.
  - Laufzeit und Speicherbedarf hängen exponentiell von  $k$  ab.  
 ein Eintrag des Arrays: Minimumbildung über  $2^k - 1$  Werte  
 (jede mögliche Kombination von Lücken in der aktuell betrachteten Spalte)
  - insgesamt Zeitkomplexität  $O(k^2 \cdot 2^k \cdot n^k)$ .

Def. 5.15: Die Entscheidungsvariante des Mult-SP-Align-Problem ist wie folgt definiert: Das Mult-SP-Align-Problem (DMSPAP)

Eingabe: nat. Zahl  $k$  und Menge  $S = \{s_1, \dots, s_k\}$  über  $\Sigma$ , eine Bewertungsfunktion  $d: (\Sigma \cup \{-\})^2 \rightarrow \mathbb{Q}$  und ein  $d \in \mathbb{N}$

Ausgabe: JA, falls multiples Alignment ex., das SP-Bewertung begl.  $d$  hat, die sd ist.  
 NEIN, sonst.

Def. 5.16: Die Entscheidungsvariante des Problems der kürzesten gemeinsamen Supersequenz über  $\Sigma = \{0, 1\}$  (D-(0,1)-SSP) ist wie folgt definiert:

Eingabe:  $k \in \mathbb{N}$ ,  $S = \{s_1, \dots, s_k\}$  von Strings über  $\Sigma = \{0, 1\}$ ,  $m \in \mathbb{N}$

Ausgabe: JA, falls gemeinsame Supersequenz  $t$  der Strings aus  $S$  ex. mit  $|t| \leq m$ .

NEIN, sonst.

Lemma 5.2: D-(0,1)-SSP ist NP-vollständig. □

Satz 5.3: D-MSPAP ist NP-vollständig.

Beweis: D-MSPAP  $\in$  NP ✓

Polynomialzeit-Reduktion vom D-(0,1)-SSP auf

D-MSPAP:

Seien  $S = \{s_1, \dots, s_k\}$  Menge von Strings über  $\{0, 1\}$  und  $m \in \mathbb{N}$  als Eingabe für D-(0,1)-SSP gegeben.

O.B.d.A.  $\max\{|s_i| \mid 1 \leq i \leq k\} \leq m$

Konstruiere  $m+1$  verschiedene Eingabeinstanzen für D-MSPAP

und zeige, daß es eine Supersequenz der Länge  $\leq m$  für

$S$  gibt gdw. eine der konstruierten Eingabeinstanzen

für D-MSPAP ein multiples Minimum mit SP-Bewertung

unter der geg. Schranke hat.

- Für alle  $i, j \in \mathbb{N}$  mit  $i+j = m$  sei  $X_{i,j} = S_{i,j}(a, b)$ ,  $a, b \in \Sigma$

-  $d = (k-1) \cdot \|S\| + (2k+1) \cdot m$ , wobei  $\|S\| = \sum_{s \in S} |s|$

-  $\delta = (\{0, 1, a, b, -\})^2 \rightarrow \mathbb{N}$ .

	0	1	a	b	-
0	2	2	1	2	1
1	2	2	2	1	1
a	1	2	0	k	1
b	2	1	k	0	1
-	1	1	1	1	0

Zeige: Die Menge  $S$  hat eine Supersequenz  $x$  mit  $|x|=m$  genau dann, wenn eine der Menge  $X_{i,j}$  ein multiples Alignment mit SP-Bewertung  $\leq d$  bzgl.  $\delta$  hat.

1. Sei für ein  $X_{i,j}$  ein multiples Alignment  $A = (s_1, \dots, s_k, \alpha, \beta)$  mit SP-Bewertung  $\leq d$  gegeben. Dabei sei  $\alpha$  die zu  $a^i$  gehörige Zeile, und  $\beta$  die zu  $b^j$  gehörige Zeile.

Betrachte zunächst die Einschränkung  $A'$  von  $A$  auf die Zeilen  $s_1, \dots, s_k$ .

Zeige: Bewertung von  $A'$  ist  $(k-1) \cdot \|S\|$ .

Untersuche hierfür  $A'$  spaltenweise:

Spalte mit  $l$  Lücken hat Bewertung von

$$\underbrace{l \cdot (k-l)}_{\substack{\text{Vergleich Lücken} \\ \leftrightarrow \text{Mißlichkeiten}}} + 2 \cdot \underbrace{\frac{(k-l)(k-l-1)}{2}}_{0,1 \leftrightarrow 0,1} = (k-1) \cdot (k-l)$$

Sei  $x$  die Anzahl der Spalten von  $A'$  und  $y$  die Anzahl aller Lückensymbole in  $A'$ . Dann gilt  $\|S\| = k \cdot x - y$ .

Sei  $l_p$  die Anzahl der Lücken in Spalte  $p$  von  $A'$  für alle  $1 \leq p \leq x$ . Dann gilt:

$$\begin{aligned} \delta_{sp}(A) &= \sum_{p=1}^x (k-1) \cdot (k-l_p) = (k-1) \sum_{p=1}^x (k-l_p) \\ &= (k-1) \cdot (k \cdot x - y) = (k-1) \cdot \|S\|. \end{aligned}$$

Buch: ISBN 3-519-00395-8 / Teubner

Algorithmische Grundlagen der Bioinformatik

Böckenhauer / Bongartz

Zeige, daß es ein Alignment von  $X_{ij}$  mit SP-Bewertung  $d$  nur dann geben kann, wenn dieses Alignment die Länge  $m$  hat und in keine Spalte  $a$  und  $1$  oder  $b$  und  $0$  gemeinsam vorkommen.

Vergleich von  $d$  mit  $\beta$ . Kosten von  $m$ , falls in keine Spalte  $a$  und  $b$  gleichzeitig vorkommen.

Zusätzliche Kosten  $k$  für jede Spalte, in der  $a$  und  $b$  vorkommen.

Falls die Länge  $x$  des Alignment  $< m$  ist, haben wir Kosten von  $\geq m + (m-x) \cdot k$

Vergleich von  $d$  mit  $s_1', \dots, s_k'$  bzw. von  $\beta$  mit  $s_1', \dots, s_k'$ :

Kosten:  $\geq 2 \cdot k \cdot x - z$ , wobei  $z$  die Anzahl der Paare von Lückensymbolen

Da  $d$  nur  $x-i$  Lücken und  $\beta$  nur  $x-j$  Lücken enthält, gilt:  $z \leq k \cdot (x-m)$ , da  $i+j = m$

$$\Rightarrow \text{Kosten: } \geq 2 \cdot k \cdot x - k \cdot (x-m) = k \cdot (x+m)$$

Zusätzliche Kosten für jede Spalte, in der  $a$  und  $1$  oder  $b$  und  $0$  auftreten.

$$\Rightarrow d_{SP}(A) \geq (k-1) \cdot \|S\| + \max\{m, m+k(m-1)\} + k \cdot (x+m)$$

Minimum wird erreicht für  $x=m$ :

$$\begin{aligned} (k-1) \cdot \|S\| + m + k \cdot (x+m) &= (k-1) \cdot \|S\| + m + 2 \cdot k \cdot m \\ &= d \end{aligned}$$



→ Falls Alignment mit SP-Bewertung  $\leq d$  ex., dann gilt  $x = m$  und in keiner Spalte treten  $a$  und  $b$  oder  $a$  und  $1$  oder  $b$  und  $0$  gemeinsam auf.

⇒ In keiner Spalte treten  $0$  und  $1$  gemeinsam auf

⇒ Bestimme Supersequenz  $t$  von  $S_1, \dots, S_k$  der Länge  $m$  wie folgt:

$$t_k = \begin{cases} 0, & \text{falls } \alpha_k = a \\ 1, & \text{falls } \beta_k = b \end{cases}$$

2. Sei eine gemeinsame Supersequenz  $t$  für  $S$  der Länge  $m$  gegeben.

Sei  $i$  die Anzahl der Nullen in  $t$ ,  $j$  die Anzahl der Einsen in  $t$ .

Dann besitzt  $X_{i,j}$  ein multiples Alignment mit SP-Bewertung  $\leq d$ : für jedes  $p \in S$  ex. ein Alignment von  $p$  und  $t$  ohne Mismatches der Länge  $m$ .

⇒ Fasse diese Alignments zusammen, ordne die  $a$ -Symbole an  $a^i$  den  $0$ -Spalten und die  $b$ -Symbole an  $b^j$  den  $1$ -Spalten zu.

Analog Rechnung wie in 1. zeigt, daß dieses multiples Alignment eine Bewertung von  $d$  hat.

□

5.3.3. Zusammenfügen paarweiser Alignments

Def. 5.17: Sei  $S = \{s_1, \dots, s_k\}$  eine Menge von Strings,  
 sei  $T = \{s_{i_1}, \dots, s_{i_m}\} \subseteq S$ .

Sei  $A' = (s_{i_1}', \dots, s_{i_m}')$  ein multiples Alignment von  $S$ .

$A'' = (s_{i_1}'', \dots, s_{i_m}'')$  ein multiples Alignment von  $T$ .

$A'$  heißt **kompatibel** zu  $A''$ , falls die Einschränkung von  $A'$  auf die Zeilen  $i_1, \dots, i_m$ , bei der alle Spalten eliminiert werden, die nur aus Lücken bestehen, gleich  $A''$  ist.

Beispiel 5.11

$$\begin{array}{l} A - C G G \\ A - - T G \\ A T C G G \end{array} \quad \begin{array}{l} A - C G G \\ A - - T G \end{array} \rightarrow \begin{array}{l} A C G G \\ A - T G \end{array}$$

$$\begin{array}{l} A - - T G \\ A T C G G \end{array} \quad \neq \quad \begin{array}{l} A T - G - \\ A T C G G \end{array}$$

Def. 5.18: Sei  $S = \{s_1, \dots, s_k\}$  eine Menge von Strings.

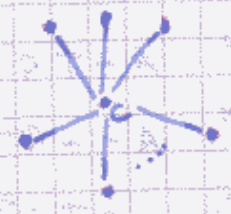
Ein Baum  $T = (V, E)$  mit  $V = \{s_1, \dots, s_k\}$ , bei dem jede Kante  $\{s_i, s_j\} \in E$  mit einem optimalen Alignment  $(s_i', s_j')$  beschriftet ist, heißt **Alignment-Baum** für  $S$ .

Satz 5.3. Sei  $S = \{s_1, \dots, s_k\}$  eine Menge von Strings, sei  $T = (V, E)$  ein Alignment-Baum für  $S$ .

Dann läßt sich ein multiples Alignment  $(s_1'', \dots, s_k'')$  für  $S$ , das mit den optimalen paarweisen Alignments  $(s_i', s_j')$  kompatibel ist für alle  $\{s_i, s_j\} \in E$ , effizient bestimmen.

□

Im folgenden: Spezialfall: T ist ein Stern  
⇒ Star-Alignment



Algorithmus 5.3 - Star-Alignment

Eingabe: Menge  $S = \{s_1, \dots, s_k\}$  von Strings

1. Berechne das Zentrum c des Sterns:

a) Bestimme optimales paarweises Alignment von  $s_i$  und  $s_j$   
für alle  $1 \leq i, j \leq k$  und berechne dessen Bewertung  
 $\text{sim}(s_i, s_j)$

b) Bestimme c als denjenigen String  $t$ , der  $\sum_{s \in S} \text{sim}(t, s)$   
minimiert.

c) Setze T als Stern mit Zentrum c und Blättern  $S - \{c\}$

2. Bestimme kompatibles multiples Alignment:

for  $i := 2$  to  $k$  do

Bestimme ein zu T kompatibles multiples Alignment von  
c und  $s_1, \dots, s_i$  aus den bereits bestimmten  
kompatiblen multiplen Alignment c und  $s_1, \dots, s_{i-1}$   
und dem optimalen paarweisen Alignment von c und  $s_i$   
nach dem Prinzipi „Einmal eine Lücke - immer eine Lücke“

Ausgabe: Das zu T kompatible multiple Alignment von S.

Beispiel:  
c: ATG-CATT      c': A-TGC-ATT  
s<sub>1</sub>: A-GTCAAT      s<sub>2</sub>: A-CTGTAAAT  
s<sub>3</sub>: -TCTGA--

$C''$ : A-T-C-A-T-T  
 $C'''$ : A--GTC-AAT  
 $C''''$ : --TCTC-A--  
 $C'''''$ : ACTC-TAATT

Def. 5.13: Eine Bewertungsfunktion  $\delta: (\Sigma \cup \{-\})^2 \rightarrow \mathbb{Q}$  heißt gut, wenn gilt:

- (i)  $\delta(a, a) = 0$  für alle  $a \in \Sigma \cup \{-\}$
- (ii)  $\delta(a, c) \leq \delta(a, b) + \delta(b, c)$  für alle  $a, b, c \in \Sigma \cup \{-\}$   
(Dreiecksungleichung)

Lemma 5.3: Sei  $\delta$  eine gute Bewertungsfunktion. Dann gilt  $\delta(a, b) \geq 0$  für alle  $a, b \in \Sigma \cup \{-\}$

Beweis: Es gilt  $0 = \delta(a, a) \leq \delta(a, b) + \delta(b, a) = 2 \cdot \delta(a, b)$  für alle  $a, b \in \Sigma \cup \{-\}$ . □

Lemma 5.4: Sei  $\delta: (\Sigma \cup \{-\})^2 \rightarrow \mathbb{Q}$  eine gute Bewertungsfunktion.

Sei  $S = \{c, s_1, \dots, s_k\}$  eine Menge von Strings über  $\Sigma$ .

Sei  $T = (S, E)$  ein Stern mit Zentrum  $c$  und sei

$(c', s_1', \dots, s_k')$  ein zu  $T$  kompatibles multiples Alignment

von  $S$ . Dann gilt für alle  $i, j \in \{1, \dots, k\}$

$$\delta(s_i', s_j') \stackrel{(1)}{\leq} \delta(s_i', c') + \delta(c', s_j') \stackrel{(2)}{=} \text{sim}(s_i, c) + \text{sim}(c, s_j)$$

Beweis: (1) Dreiecksungleichung

(2) folgt, da die von  $(c', s_1', \dots, s_k')$  induzierten paarweisen Alignments zwischen  $s_i$  und  $c$  sowie zwischen  $c$  und  $s_j$  optimal sind. □

Satz 5.4: Sei  $\delta$  eine gute Bewertungsfunktion, sei  $\delta_{sp}$  die von  $\delta$  induzierte SP-Bewertungsfunktion. Sei  $S = \{s_1, \dots, s_k\}$  eine Menge von Strings, sei  $\text{sim}_{sp}(S)$  die SP-Bewertung eines optimalen multiplen Alignment von  $S$ .

Dann gilt für das von Abg. 5.3 beschriebene multiple Alignment  $(s_1', \dots, s_k')$ :

$$\delta_{sp}(s_1', \dots, s_k') \leq \left(2 - \frac{2}{k}\right) \cdot \text{sim}_{sp}(s_1, \dots, s_k).$$

Beweis: Sei  $(s_1'', \dots, s_k'')$  ein optimales multiples Alignment von  $S$ , d.h.

$$\delta_{sp}(s_1'', \dots, s_k'') = \text{sim}_{sp}(s_1, \dots, s_k)$$

definiere  $v(s_1', \dots, s_k') = \sum_{i=1}^k \sum_{j=1}^k \delta(s_i', s_j') = 2 \cdot \delta_{sp}(s_1', \dots, s_k')$

und  $v(s_1'', \dots, s_k'') = \sum_{i=1}^k \sum_{j=1}^k \delta(s_i'', s_j'') = 2 \cdot \delta_{sp}(s_1'', \dots, s_k'')$   
 $= 2 \cdot \text{sim}_{sp}(s_1, \dots, s_k)$

Es reicht zu zeigen:  $\frac{v(s_1', \dots, s_k')}{v(s_1'', \dots, s_k'')} \leq 2 - \frac{2}{k}$

Sei  $M := \min_{t \in S} \sum_{s \in S} \text{sim}(s, t) = \sum_{s \in S} \text{sim}(c, s) = \sum_{s \in S - \{c\}} \text{sim}(c, s)$

a.B.d.A.:  $c = s_k$

Nach Lemma 5.4 gilt:

$$v(s_1', \dots, s_k') = \sum_{i=1}^k \sum_{j=1}^k \delta(s_i', s_j') \leq \sum_{i=1}^k \sum_{j=1}^k (\text{sim}(s_i, c) + \text{sim}(s_j, c))$$

Bolnt. 46

$$\begin{aligned} &= \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} (\text{sim}(s_i, c) + \text{sim}(s_j, c)) \\ &= \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} \text{sim}(s_i, c) + \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} \text{sim}(s_j, c) \\ &= 2 \cdot (k-1) \cdot \sum_{i=1}^{k-1} \text{sim}(s_i, c) = 2 \cdot (k-1) \cdot M \end{aligned}$$

$$\begin{aligned} v(s_1'', \dots, s_k'') &= \sum_{i=1}^k \sum_{j=1}^k \delta(s_i'', s_j'') \geq \sum_{i=1}^k \sum_{j=1}^k \text{sim}(s_i, s_j) \\ &\geq k \cdot \sum_{j=1}^k \text{sim}(c, s_j) = k \cdot M \end{aligned}$$

$$\frac{v(s_1', \dots, s_k')}{v(s_1'', \dots, s_k'')} \leq \frac{2 \cdot (k-1) \cdot M}{k \cdot M} = 2 - \frac{2}{k} \quad \square$$

26.2.2013

## 5.2. Algorithmen zur Datenbanksuche

### 5.2.1. Das FASTA-Verfahren

Ziel: Vergleiche zu suchendes Muster (Datenbankanfrage) nacheinander mit allen gespeicherten Sequenzen (Datenbank-Strings).

#### Prinzipielle Vorgehensweise:

1. Wähle Parameter  $k$  und suche alle exakten Matches der Länge  $k$  (k-Mer-Spots)

übliche Werte:  $k=6$  für DNA  
 $k=2$  für Proteine

2. Fasse mehrere Dot-Spots zusammen:

Betrachte Matrix  $M$  für Muster  $p$  und Text  $t$

$$M(i, j) = \begin{cases} 1 & \text{falls } p_i = t_j \\ 0 & \text{sonst} \end{cases}$$

$\Rightarrow$  Dot-Spot  $\hat{=}$  Abschnitt einer Diagonalen

Suche 10 besten diagonalen Stücke, die mit Dot-Spot beginnen und enden.

Bewertung: Anzahl der Dot-Spots: positiv  
Länge des Stückes dazwischen: negativ

Bestimm. für jeden Lauf <sup>diagonalen</sup> opt. lokales Alignment

3. Setze die so berechneten Teilalignments zu längerem Alignment zusammen.

4. Berechne Alternativlösung durch lokales Alignment, beschränkt auf einen Streifen konstanter Breite um das beste lokale Alignment aus Schritt 2 herum.

Anschließend statistische Bewertung der berechneten Lösungen.

## 5.22. Das BLAST-Verfahren

Suchalgorithmus:

1. Suche ähnliche Teilstrings, sogenannte **kits**, gegebenes Länge  $w$ .

übliche Werte:  $w=11$  für DNA  
 $w=3$  für Proteine

2. Suche alle Paare von Oligos, die Abstand  $\leq d$  haben.
3. Erweitere die Paare von Oligos an beiden Enden, bis sich die Alignment-Bewertung nicht mehr erhöht.  
 Falls Bewertung über einen Schwellenwert  $S$  liegt:  
High-Scoring-Pair

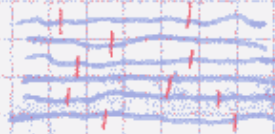
Anschließend statistische Bewertung der Lösungen.

## Teil II - DNA-Sequenzierung

Menschliches Genom: 35 Gbp

dichte Sequenzierung:  $< 1000$  bp

- erzeugen Kopien von  $k$ :

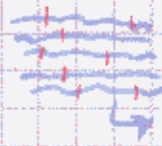


50-300 kbp

Ordnung der Fragmente geht verloren!

→ Physikalische Kartierung

- nehme ein Fragment aus einer sogenannten physikalischen Karte  
 Sequenziere dieses Fragment



~ 1000 bp

} DNA-Sequenzierung

Human-Genom-Projekt

Calera Genomics

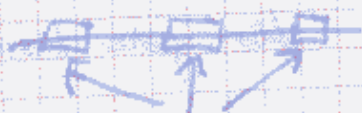


7. Physikalische Kartierung:

Def. 7.1: Sei  $D$  eine DNA-Sequenz.

Eine **Physikalische Karte** von  $D$  besteht aus einer Menge  $M$  von **Markern** und einer Funktion  $p: M \rightarrow \text{Pot}(\mathbb{N})$ , die für jeden Marker dessen Position in  $D$  angibt.

$m \in M$

7.1 Restriktionsstellen - Kartierung

Verwenden: Restriktionsenzyme mit spezifischen Restriktionsstellen zum **erschneiden (Verdau)** der DNA.

→ Restriktionsstellen dienen als Marker

→ Aufgabe: Bestimme die Reihenfolge der durch den Verdau entstandenen Fragmente

7.1.1. Das Double-Digest-Verfahren

Beachte:

- $\Delta(A)$ ,  $\Delta(B)$ ,  $\Delta(AB)$  sind Multimengen
- (idealerweise) Vollverdau, bzw. full-digest, d.h. falls eine Restriktionsstelle vorhanden ist, so wird das Molekül auch dort geschnitten.

Def. 7.2: Sei  $X = \{x_1, \dots, x_n\}$  eine Multimenge mit Elementen aus  $\mathbb{N} \setminus \{0\}$ . Sei  $\pi$  eine Anordnung dieser Elemente,  $\pi = (x_{i_1}, \dots, x_{i_n})$ . Dann

$$\text{Pos}(\pi) = \left\{ 0, x_{i_1}, x_{i_1} + x_{i_2}, \dots, \sum_{j=1}^n x_{i_j} \right\}$$

Positionsmenge der Anordnung  $\pi$ .

## Biotuf 50

Umgekehrt sei  $X$  eine Menge von Elementen aus  $\mathbb{N}$  und sei  $Y = \{y_1, \dots, y_p\}$ ,  $y_1 < y_2 < \dots < y_p$ .

Dann  $\text{Dist}(Y) = \{ |y_i - y_j| \mid i < j \in \{1, \dots, p\} \}$   
Distanzmenge der Menge  $Y$

Es gilt:  $\text{Dist}(P_{\text{as}}(X)) = X$  für jede Anordnung  $\pi$  einer Multimenge  $X$ .

Def. 7.3: Seien  $A, B, C$  Multimengen, mit  $|A| = n$ ,  $|B| = m$ .

Sei  $\pi$  eine Anordnung der Elemente aus  $A$  und  $\rho$  eine Anordnung der Elemente aus  $B$ .

Das Paar  $(\pi, \rho)$  heißt zulässige Lösung für  $A, B, C$ , wenn gilt:

$$\text{Dist}(P_{\text{as}}(\pi) \cup P_{\text{as}}(\rho)) = C$$

Def. 7.4: Das double-digest-Problem (DDP) ist das folgende Berechnungsproblem:

Eingabe: Multimengen  $A, B, C$

Ausgabe: Ein Element aus der Menge

$$U = \{ (\pi, \rho) \mid (\pi, \rho) \text{ zulässige Lösung von } A, B, C \} \\ \text{oder den Wert } 0, \text{ wenn } U = \emptyset.$$

→ die Menge  $U$  heißt auch die Menge der zulässigen Lösungen für das DDP.

$$U = \{ (\pi, \rho) \mid (\pi, \rho) \text{ zulässige Lösung von } A, B, C \}$$

naive Ansatz: Teste alle möglichen Anordnungen  $\pi$  von  $\phi$ .  
 $\Rightarrow |A|! \cdot |B|!$  viele Möglichkeiten testen ... "unreal"

Def. 7.5: die Entscheidungsvariante des Double-Subset-Partition-Problems, kurz Dec-DDP, lautet wie folgt:

Eingabe: Multimengen  $A, B, C$

Ausgabe: Ja, falls  $|W| \geq 1$   
 Nein, sonst

Def. 7.6: Das Set-Partition-Problem:

Eingabe:  $X = \{x_1, \dots, x_n\}$  mit Elementen aus  $\mathbb{N} \setminus \{0\}$

Ausgabe: Ja, wenn eine Zerlegung von  $X$  in 2 disjunkte Mengen  $Y$  und  $Z$  existiert, so dass

$$\sum_{y \in Y} y = \sum_{z \in Z} z$$

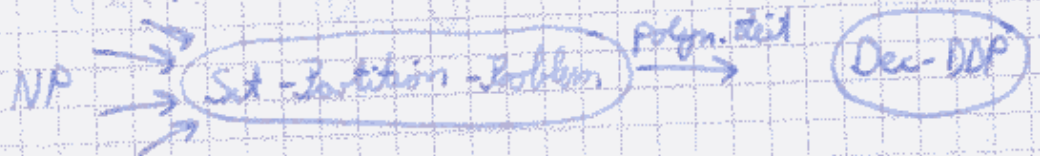
Nein, sonst

Es ist bekannt, dass das Set-Partition-Problem NP-vollständig ist.

Satz 7.1: Das Dec-DDP ist NP-vollständig!

Beweis: (i) Dec-DDP  $\in$  NP ✓

(ii) Alle Probleme aus NP in polynomiellem Zeit lassen sich auf Dec-DDP zurückführen



Sei  $X$  eine Eingabe für das Set-Partition-Problem.

O.B.d.A. sei die Summe aller Werte in  $X$  gerade.

$\rightarrow$  konstruiere Eingabe für Dec-DDP  $(A, B, C)$

- $A := X = \{x_1, \dots, x_n\}$
- $B := \left\{ \frac{\alpha}{2}, \frac{\alpha}{2} \right\}$ , wobei  $\alpha = \sum_{x \in X} x$
- $C := X$

zz. Es ex. eine Lösung für das St-Partition-Problem mit Eingabe  $X$  genau dann wenn eine Lösung für das Dec-DDP mit Eingabe  $A, B, C$  (wie oben) existiert.

$\Rightarrow$  Sei  $x$  teilbar in 2 Mengen  $Y$  und  $Z$ , mit  $Y$  und  $Z$  sind disjunkt

$$\sum_{y \in Y} y = \sum_{z \in Z} z$$

Dann ist  $\pi = (P(Y), P'(Z))$  und  $\rho = \left( \frac{\alpha}{2}, \frac{\alpha}{2} \right)$  mit  $P, P'$  beliebige Ordnungen, eine Lösung für das Dec-DDP.

$\Leftarrow$  Sei  $(\pi, \rho)$  eine zulässige Lösung für das Dec-DDP. Sei  $\pi = (x_{j_1}, \dots, x_{j_n})$  dann existiert ein Index  $i_0$  mit

$$\sum_{i=1}^{i_0} x_{j_i} = \sum_{i=i_0+1}^n x_{j_i}$$

da diese Aufteilung wegen  $B = \left\{ \frac{\alpha}{2}, \frac{\alpha}{2} \right\}$  notwendig ist.

### 7.1.2. Das Partial-Digest-Verfahren

- unvollständiger Verdau (partial-digest)
- Multimenge  $\Delta_p(A)$  im Gegensatz DDF

Im Folgenden gehen wir von idealen Daten aus.

Def. 7.7 die ermittelte Datenmenge  $\Delta_p(A)$  heißt **ideal**, wenn sie die Länge jedes durch Restriktionstellen bzw. Enden des Moleküls begrenzten Fragments enthält.

d.h. wir haben Schnittstellen (bzw. Endpunkte)

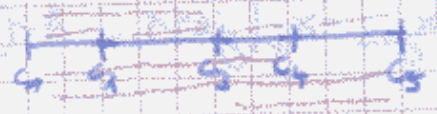
$$c_1 < c_2 < \dots < c_k$$

dann gilt:

$$\Delta_p(A) = \{c_j - c_i \mid 1 \leq i < j \leq k\}$$

→ eine ideale Multimenge für das PD enthält

$\binom{k}{2}$  Elemente (k ist die Anzahl der Restriktionstellen + Endpunkte)



Def. 7.8: Sei  $A$  eine Multimenge mit  $\binom{k}{2}$  Elementen aus  $\mathbb{N} - \{0\}$ , sei  $P = \{x_1, \dots, x_k\}$  eine Menge von natürlichen Zahlen mit  $x_1 = 0$  und  $x_1 < x_2 < \dots < x_k$ .  
 Dann heißt  $P$  auch **Punktmenge**.

Zu  $P$  lässt sich eine Multimenge aller paarweisen Distanzen definieren:

$$\text{Dist}(P) = \{x_j - x_i \mid 1 \leq i < j \leq k\}$$

Eine Restriktionsmenge  $P$  heißt zulässige Lösung für  $A$  falls

$$\text{Dist}_p(P) = A$$

Def. 7.9: Partial-digest-Problem (PDP)

Eingabe: Multimenge  $A$  mit  $\binom{k}{a}$  Elementen aus  $\mathbb{N} - \{0\}$

Abgabe: Ein Element aus

$$\mathcal{N} = \{P, P \text{ ist zulässige Lösung für } A\}$$

oder  $\emptyset$ , falls  $\mathcal{N} = \emptyset$ .

Beachte: Im Gegensatz zum DDP ist hier "keine Anordnung" von den Elementen in  $A$  gesucht.

Def. 7.10: Sei  $A$  eine Multimenge mit  $\binom{k}{a}$  Elementen aus  $\mathbb{N} - \{0\}$ . Sei  $P = \{x_1, \dots, x_k\}$  eine zulässige Lösung für  $A$ , mit  $x_i \geq 0$  und  $x_1 \leq x_2 \leq \dots \leq x_k$ . Dann bezeichne

$$\text{level}_p(i) = \{x_j \mid j \in \{1, \dots, k-1\}\} \subseteq A$$

die Multimenge von Distanzen deren Endpunkte den Abstand  $i$  in der Lösung  $P$  besitzen.

Bemerkung 7.1:  $A$  Multimenge und  $P$  zulässige Lösung (wie oben)

- $|\text{level}_p(i)| = k - i$

- $\text{level}_p(k)$ , Distanzen zwischen benachbarten Restriktionsstellen  $\rightarrow$  full digest

$\rightarrow$  atomare Distanzen

- $level_p(k-1) \cong$  maximale Distanz in  $A$ , Länge des betrachteten Moleküls.
- Level bilden "disjunkte" Zerlegung von  $A$ .

$$level_p(1) \cup level_p(2) \cup \dots \cup level_p(k-1) = A$$

Naiver Ansatz:

- wähle  $k-1$  atomare Distanzen aus  $A$   $\binom{k}{2}$
- überprüfe alle Anordnungen:  $(k-1)!$
- Laufzeit:  $O((k-1)! \binom{k}{2})$  exponentiell
- Ermittle atomare Distanzen durch full-digest-Experiment bestimmen  $\rightarrow O((k-1)!)$

Backtracking:

- sukzessive Spezifikation von Teillösungen
- falls eine Spezifikation nicht "abgeschlossen"  $\rightarrow$  Backtracking Schritte, nehme letzte Spezifikation zurück

Line: Teillösungen  $\cong$  festgelegte Positionen in der Punktmenge.

Notation:  $y \in \mathbb{N}$ , Multimenge  $X = \{x_1, \dots, x_n\}$

$$\delta(y, X) = \{ |x-y| \mid x \in X \}$$

Schematische Darstellung des Backtracking-Algorithmus: (Idee)

- 1) plane längste Distanz  $\rightarrow$  Intervall, in dem alle weiteren Punkte liegen, festgelegt.

2) Für die jeweils längste verbleibende Distanz:

- überprüfe, ob Platzierung am linken Rand möglich  
 → falls ja, platziere links
- ansonsten, überprüfe, ob Platzierung am rechten Rand möglich  
 → falls ja, platziere rechts
- ansonsten: Backtracking-Schritt

3) Gebe Lösung aus, wenn alle Distanzen erfolgreich platziert wurden.

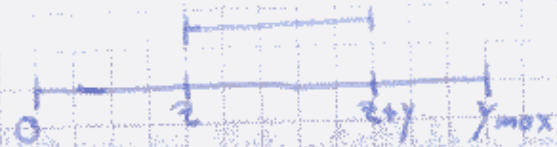
Satz 7.2: Sei  $A$  eine Eingabe für das PDP.

Falls eine zulässige Lösung für das PDP mit Eingabe  $A$  existiert, dann berechnet der Algorithmus 7.1 eine solche.

Beweis: Der Algorithmus 7.1 untersucht alle Lösungsmöglichkeiten, bei denen die jeweils längste verbleibende Distanz  $y$  entweder am linken oder rechten Rand platziert wird.  
 → zeige: damit werden alle möglichen Lösungen durchlaufen.

Sei  $X = \{0, x_1, x_2, \dots, x_{max}\}$  die bisher konstruierte Teillösung,  $A$  Multimenge mit verbleibenden Distanzen,  $y$  momentan größter Distanz.





Angenommen die Distanz  $y$  könnte "in der Mitte" platziert werden, also in einem Intervall  $[z, z+y]$  ( $z \neq 0, z+y \neq y_{\max}$ )

- wenn  $z+y \neq X$ , so muss  $z+y \in A$   
 $\rightarrow$  aber  $z+y > y \Rightarrow y$  ist nicht größtes Element in  $A$  ✓
- wenn  $z \neq X \Rightarrow y_{\max} - z \in A$   
 $\rightarrow$  es gilt  $y_{\max} - z > y \Rightarrow y$  ist nicht größtes Element in  $A$  ✓
- wenn  $z \in X$  und  $z+y \in X$   
 $\Rightarrow$  notwendige Distanz der Länge 0 in  $A$  ✓

□

Satz 7.3: Das Algorithmus 7.1 hat im schlechtesten Fall eine Laufzeit von  $O(2^k \cdot k \log k)$  für eine Eingabe mit  $\binom{k}{2}$  Elementen.

- Beweis:
- Initialisierung:  $O(1)$
  - Sortierung von  $A$ :  $O(k^2 \cdot \log k^2) = O(k^2 \cdot \log k)$
  - Platzierung von  $k-1$  Distanzen, wobei Platzierung am rechten- und linken Rand möglich sind.  
 $\rightarrow 2^{k-1}$  Platzierungsmöglichkeiten
  - Funktion  $f$  liefert  $O(k)$  zu überprüfende Distanzen  
 Aufwand ob Platzierung möglich / Platzierung durchführen  
 $O(k \cdot \log k^2) = O(k \cdot \log k)$   
 Aufwand für das Backtracking (analog zur Platzierung)  
 $O(k \cdot \log k)$
- $\Rightarrow O(k^2 \cdot \log k + 2^{k-1} \cdot k \cdot \log k) = O(2^k \cdot k \log k)$

Algorithmus 7.1 hat im schlechtesten Fall exponentiellen Aufwand  
(es existieren entsprechende Eingaben)

Satz 7.4: Sei die Anzahl der Backtracking-Schritte, die  
Algorithmus 7.1 benötigt, unabhängig von  $k$  für eine  
Eingabe  $A$  der Größe  $\binom{k}{2}$ , so kann die Laufzeit durch  
 $O(k^2 \log k)$  abgeschätzt werden.

Gemessen an der Eingabegröße  $\binom{k}{2}$  ergibt sich dann die  
Laufzeit  $O(n \cdot \log n)$  ( $n = \binom{k}{2}$ ).

Beweis: Plaziere  $O(k)$  Elemente (nicht wie in Satz 7.3  $2^{k-1}$ )  
(analog zu Beweis von Satz 7.3)

Zusammenfassung:

DDP: einfacher Exprimiert  $\leftrightarrow$  "schweres" Problem

PDP: "komplexes" Exprimiert  $\leftrightarrow$  "leichtes" Problem  
(guter Algorithmus)

Konzept des Fingerprinting

- Ordnet jedem Fragment spezifische Agarose zu  
↓  
Fingerprints / Fingerabdrücke
- Fingerprints sind leicht zu ermitteln
- Zwei Fragmente überlappen einander  $\Leftrightarrow$  ähnliche Fingerprints

Kandidaten für Fingerprints

- Restriktionsstellen - Kartierung
- Fragmentgrößen nach Verdau durch Restriktionsenzyme
- Hybridisierungsdaten

Eckkurs: Biologische Grundlagen: Hybridisierung - Klonierung -  
 DNA-Chips

Hybridisierung: Aneinanderlagerung (komplementärer) Nucleinsäure-  
 stränge

A	C	T	A
T	G	A	T

$\Rightarrow$  Test, ob ein bestimmter Teilstrang in einem Nucleinsäurestrang auftritt.

Ziel: Viele Hybridisierungsexperimente gleichzeitig durchführen

dazu: Vermehrung von DNA

1. Klonierung:

- Einsetzen der zu kopierenden DNA
- Onkret in einen Vektororganismus (Plast)
- Replikation des Plast repliziert auch den insert
- Extraktion des insert aus dem Plast

Länge des Inrets: 15-50 kbp für Bakterien / Viren  
 mehrere Millionen für tierische Chromosomen  
Problem: z.B. Übermengen mit der Host-DNA.

2) Polymase-Kettenreaktion (PCR)

Schritt 1: Gabe in Reagenzglas

- zu kopierende DNA d
- Primer  $P_1$  und  $P_2$  (DNA-Stränge (Einkettstränge), die komplementär zum Anfang, bzw. zum Ende von d sind).
- alle Nucleotide in ausreichender Menge
- DNA-Polymerase (Enzym, das aus einem Primer entsprechend einer Folge sukzessive einen kompletten Molekülstrang aufbaut).

Schritt 2: (wiederholte beliebig oft)

- denaturiere DNA (durch Erhitzen werden die beiden Einzelstränge voneinander getrennt).
- abkühlen  $\rightarrow$  Primer lagern sich an die Einzelstränge an  
 $\rightarrow$  DNA-Polymerase verlängert die Primer zu einem vollständigen komplementären Strang.

n Generationen  $\rightarrow$  bis zu  $2^n$  Kopien

Voraussetzung: Kenntnis der Primer

Problem: Fehler in frühen Generationen können sich exponentiell fortplanzen.

DNA-chips:

-  $s, t$  DNA-Sequenzen (Einzeltsträngig)

Aufgabe: Teste, ob  $t$  in  $s$  als Teilstrang vorkommt

→ Führe Hybridisierungsexperiment von  $s$  und  $t'$  durch, wobei  $t'$  die zu  $t$  komplementäre Sequenz ist.

-  $s, t_1, \dots, t_n$  DNA-Sequenzen (Einzeltsträngig)

Teste „parallel“, ob  $t_1, \dots, t_n$  als Teilstrang in  $s$  auftreten

→ 1.) Markiere  $t_1, \dots, t_n$  auf einer Oberfläche (Positionen von  $t_1$  bis  $t_n$  sind bekannt).

$t_1, \dots, t_n \hat{=} \text{Probes / Sonden}$

Oberfläche mit Probes  $\hat{=} \text{DNA-Chip}$

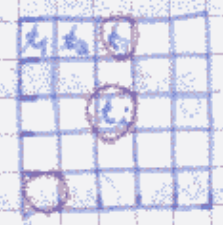
2.) Gebe markierte (z.B. fluoreszierende) Kopien der zu untersuchenden DNA auf den Chip

- Kopien der DNA  $\hat{=} \text{Clones}$

- Clones hybridisieren mit entsprechenden Probes

3.) Entferne nicht-hybridisierten Clones („abwaschen“)

4.) Bestimme (anhand der Markierungen), wo Hybridisierung stattfanden.



Fehlertypen:

- falsch positiv: Experiment meldet Hybridisierung von  $s$  und  $t_i$ , obwohl  $t_i$  kein Teilstrang von  $s$  ist.

- falsch negativ: Experiment meldet keine Hybridisierung von  $s$  und  $t_i$ , obwohl  $t_i$  ein Teilstrang von  $s$  ist.

Beachte: unbekannt, wie oft  $t_i$  als Teilstring in  $s$  auftritt.

### Methode 7.3 Kartierung durch Hybridisierung

Gegeben: zu untersuchende DNA  $d$

1) Gruppe mit hoher Wahrscheinlichkeit überlappende Fragmente von  $d$ .

2) Gruppe Kopien der Fragmente  $\rightarrow$  Clones

3)  $C = \{c_1, \dots, c_n\}$  Clones (clone-library), volle Menge von Proben  $P = \{P_{1,m}, \dots, P_m\}$

4) Führe alle Hybridisierungsexperimente  $(c_i, P_j)$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq m$

Ausgabe:  $(n \times m)$ -Hybridisierungsmatrix  $H$

$$H(i, j) = \begin{cases} 1, & \text{falls } c_i \text{ mit } P_j \text{ hybridisiert} \\ 0, & \text{sonst} \end{cases}$$

Ziel: über  $H$  die ursprüngliche Anordnung der Clones ableiten

Def. 7.11: Das Problem der Kartierung durch Hybridisierung (KdH)

Eingabe:  $(n \times m)$ -Hybridisierungsmatrix

Ausgabe: Anordnung der Clones/Proben, die durch Hybridisierungsmatrix möglichst gut erklärt.

- Kartierung mit eindeutigen Proben (unique Proben)  $\leftarrow$
- Kartierung mit mehdeutigen Proben (non-unique Proben)

Keine Eindeutigkeit  $\rightarrow$  STS-Proben

Aufgabe: Suche eine Permutation der Proben, die der realen Anordnung entspricht.

Def. 7.13: Sei  $A$  eine  $(m \times n)$ -Matrix mit Einträgen aus  $\{0, 1\}$

$A$  hat die Eigenschaft der aufeinanderfolgenden Einsen / Consecutive Ones Property (COP), falls eine Permutation  $\pi$  der Spalten von  $A$  existiert, so daß in jeder Zeile keine Null zwischen zwei Einsen steht.

Falls  $A$  schon diese Form besitzt, dann hat  $A$  die Consecutive Ones Form (COF)

Aufgabe: Besitzt eine Hybridisierungsmatrix die COP, wenn ja berechne eine entsprechende Permutation.

naiv:  $\rightarrow$  teste alle Permutationen  $\rightarrow O(n!)$  hilft nicht!

PQ-Bäumen:

Def. 7.13: Sei  $U = \{u_1, \dots, u_n\}$  eine endliche Menge von Elementen. Ein PQ-Baum über  $U$  eine Struktur

$T = (V, E, r, B, \text{label}, \text{type})$  mit:

- (i)  $(V, E)$  ist ein geordneter Baum
- (ii)  $r \in V$  ist der Wurzel von  $(V, E)$
- (iii)  $B \subseteq V$  ist die Menge der Blattknoten
- (iv)  $\text{label}: B \rightarrow U$  ist eine bijektive Abbildung der Blätter auf  $U$
- (v)  $\text{type}: V \rightarrow \{P, Q\}$  ist eine Zuordnung der inneren Knoten zu einer

Def. 7.14: Sei  $U$  eine Menge,  $T$  ein PQ-Baum über  $U$  und  $(v_1, \dots, v_n)$  die Blätter von  $T$  entsprechend ihrer Anordnung von links nach rechts.

Front von  $T$ :  $\text{Front}(T) = (\text{label}(v_1), \dots, \text{label}(v_n))$

Def. 7.15: Sei  $T$  ein PQ-Baum über einer Menge  $U = \{u_1, \dots, u_n\}$

legale Operationen:

- (i) die Anordnung der Kinder eines P-Knotens darf beliebig verändert werden.
- (ii) Die Anordnung der Kinder eines Q-Knotens darf invariant werden.

$$v_1 v_2 \dots v_n \rightarrow v_2 v_{n-1} \dots v_1 v_n$$

konsistente Permutationen

$\text{Consist}(T) = \{ \pi \mid \pi \text{ ist eine Permutation über } U, \pi \text{ ergibt sich als Front von } T \text{ durch eine Abfolge legaler Operationen} \}$

universellen PQ-Baum:  $U = \{u_1, \dots, u_n\}$



→ stellt alle möglichen Permutationen der Elemente in  $U$  dar.

Eingabe: Eine Menge  $U = \{u_1, \dots, u_n\}$  paarweise verschiedene Elemente (Proben) und eine Menge von Restriktionen  $R \in \text{Pot}(U)$  (Enten in den Ziffern der Hybridisierungsmatrix)



Aussage: Alle Permutationen  $\pi$  der Elemente aus  $\mathcal{K}$ , so daß für alle Restriktionen  $R$  gilt:  
Die Elemente aus  $R$  folgen aufeinander.

Lösung: - starte mit universellem PQ-Baum über  $\mathcal{K}$   
- Für jede Restriktion transformiere den aktuellen PQ-Baum, so daß die Menge der konsistenten Permutationen die Restriktion  $R$  erfüllt.

Satz 7.5: Sei  $M$  eine  $(n \times m)$ -Matrix mit Einträgen  $\{0,1\}$  und sei  $k$  die Anzahl der Einsen in  $M$ .  
Dann existiert ein Algorithmus (basierend auf PQ-Bäumen), der das  $\{1\}$ -Problem in Zeit  $O(n+m+k)$  löst.

15.6.2003

7.2.2 Kartierung mit eindeutigen Proben und Fehlern

Angenommen, die Anordnung der Proben in dem untersuchten Polyploid ist bekannt.

Was folgt dann für bestimmte Fehlertypen?

- falsch negativ: 

```
00001111111011111100
                    ↑
```
- falsch positiv: 

```
00010000111111110000
                    ↑
```
- Chimären: 

```
00001111100001111000
                    ↑↑↑
```

Lücken  $\hat{=}$  Block von 0, der durch Einsen begrenzt wird.

Ziel: Suche Permutation der Spalten in der Hybridisierungsmatrix, so daß die Anzahl der Lücken minimiert wird.

Def. 7.16: Das **Problem der minimalen Lückenzahl**, MinLM:

Eingabe:  $(n \times m)$ -Matrix über  $\{0,1\}$

zulässige Lösungen: Für alle Eingaben  $A$ ,

$$U(A) = \{ (i_1, \dots, i_m) \mid \text{Permutationen der Spalten in } A \}$$

Kosten: für eine zulässige Lösung  $\pi \in U(A)$

$\text{cost}(\pi, A) =$  Anzahl der Lücken in der Matrix  $A_{\pi}$ ,  
die durch Permutation des Spalten entsprechend  $\pi$  resultiert.

Optimierungsziel: Minimierung

Annahme: In der betrachteten Matrix existiert keine Zeile,  
die nur aus Nullen besteht.

Def. 7.17: Sei  $A$  eine binäre  $(n \times m)$ -Matrix. Dann ist der  
**Spaltenmatrizen-Graph**  $G_A = (V, E, c)$  definiert durch

-  $V = \{1, \dots, m\}$

-  $E = \{ (i, j) \mid 1 \leq i, j \leq m, i \neq j \}$

-  $c : E \rightarrow \mathbb{N}, c(i, j) = \left| \{ k \mid (1 \leq k \leq n, A(k, i) + A(k, j)) \} \right|$

Hamming-Distanz

→ Wenn eine Lücke in einer Zeile auftritt

→ entlang des entsprechenden Pfades in  $G_A$

zwei Bitwechsel vor, Kosten erhöhen sich um 2.

Umformulierung MinLM auf ein Problem über  $G_0$ :

"Finde einen Pfad in  $G_0$ , der jeden Knoten genau einmal besucht und minimale Kosten besitzt"

Problem: - Falls die Permutation mit einem Einserblock beginnt und endet, so tragen nur die Lücken zu den Kosten des Pfades bei.

- Falls die Permutation mit einem Nullenblock beginnt und endet, so tragen die Lücken und der 0-Wechsel am Anfang und der 1-0-Wechsel am Ende zu den Kosten bei.

→ vereinfachen! ~~...~~

Ziel: Füge eine spezielle Spalte, die nur aus Nullen besteht hinzu und betrachte statt Pfaden Kreise!

Satz 7.6: -  $A$  binäre  $(n \times n)$  Matrix

-  $A'$  binäre  $(n \times (n+1))$  Matrix ( $A$  ergänzt um Nullspalte  $p'$ )

-  $G_0$  ist der zu  $A'$  gehörende Spaltenabstamlograph

-  $\pi$  Permutation der Spalten in  $A$

-  $A_{\pi}$  die aus Anwendung der Permutation resultierende Matrix

-  $K_{\pi}$  sei nun der Kreis  $p'(\pi_1 \dots \pi_n)p'$  in  $G_0$

$$\Rightarrow \text{cost}(K_{\pi}) = 2 \cdot l + 2 \cdot n$$

wobei  $l$  die Anzahl der Lücken in  $A_{\pi}$  ist.

Beweis: folgt aus der Annahme der Spalte  $p'$ .

$\Rightarrow$  MinLM auf den Spaltenabstamlographen entspricht TSP.

Frage: Besitzt  $G_\Delta$  eine bestimmte Eigenschaft?

→ die Kosten auf den Kanten von  $G_\Delta$  erfüllen die Dreiecksungleichung.

$v_1, v_2, v_3 \in G_\Delta$ :

$$c(\{v_1, v_3\}) \leq c(\{v_1, v_2\}) + c(\{v_2, v_3\})$$

TSP auf Graphen mit  $\Delta$ -Ungleichung:

$\frac{3}{2}$ -Approximation (Christofides)

→ 2-Approximation

Spannbaum-Algorithmus für  $\Delta$ -TSP

Satz 3.1: Der Spannbaum-Algorithmus ist ein 2-Approximationsalgorithmus für das  $\Delta$ -TSP.

Beweis: 1) Laufzeit:  $\left. \begin{array}{l} \text{- Berechnung des min. Spannbaums} \\ \mathcal{O}(E \cdot \log(V)) \\ \text{- Tiefensuche } \mathcal{O}(V) \end{array} \right\} \text{polynom.}$

2) Approximationsgüte

Sie Hopf eine optimale Lösung für eine Eingabe

$G = (V, E, d)$

-  $\text{cost}(T) \leq \text{cost}(\text{Hopf})$ , denn die Löschung einer

Kante aus einem beliebigen Hamilton-Kreis

ergibt einen Spannbaum. Alle Kantenkosten

sind positiv.

- Sei  $w$  das Pfad, das bei der Tiefensuche durchlaufen

wird  $\text{cost}(w) = 2 \cdot \text{cost}(T)$ , da jede Kante des

Spannbaums einmal 2-mal durchlaufen wird.

Boinf (2)

- H entspricht W, wobei Knoten nicht mehr mehrfach besucht werden, sondern durch eine Kante überbrückt.  
wegen Dreiecksungleichung  
 $\text{cost}(H) \leq \text{cost}(W)$

- Insgesamt:  $\text{cost}(H) \leq \text{cost}(W) = 2 \cdot \text{cost}(T) \leq 2 \cdot \text{cost}(\text{Hopt})$

## 8 Bestimmung der Basensequenz

### Gelelektrophorese und Kettensbruch-Methode

#### Gelelektrophorese:

- Trennung von DNA-Molekülen entsprechend ihrer Länge
- DNA-Moleküle sind negativ geladen

→ in einem elektrischen Feld wandern sie in Richtung des positiven Pols.

- Ziele: Gelbe Gemisch zu trennende DNA auf Gel-artigen Träger

- lege elektrisches Feld an

→ große Moleküle → langsame Wanderung im Gel

→ kleine Moleküle → schnelle Wanderung im Gel

→ Wanderungsgeschwindigkeit ist antiproportional zur Größe des Moleküls.

#### Kettensbruchmethode:

- Verfahren zur (direkten) Sequenzierung von DNA mittels

Gelelektrophorese

Gegeben: 4 Reagenzgläser, beschriftet mit A, C, G, T

- 1) Erzeuge viele Kopien der zu sequenzierenden DNA (Einzelschnüre)
- 2) Verteile die Kopien auf Reagenzgläser.
- 3) Gebe in jedes Reagenzglas  $I \in \{A, C, G, T\}$  alle Nucleotide außer I hinzu (also gebe C, G, T in A)
- 4) Gebe in jedes Reagenzglas  $I \in \{A, C, G, T\}$  in einem bestimmten Verhältnis:

- Nucleotid I
- chemisch veränderte Form von I, an das der Restteil eines komplementären Strangs durch DNA-Polymerase anbindet.

5) Gebe in jedes Reagenzglas

- DNA-Polymerase und Primer

- es werden komplementäre Stränge synthetisiert
- mit hoher Wahrscheinlichkeit enthält I alle Einzelschnüre, die auf dem Nucleotid I enden (wegen dem Einbau der chem. veränderten Form)

6) Gebe die Inhalte der Reagenzgläser nebeneinander auf ein Gel und starte die Gel-Elektrophorese

- Trenne Stränge entsprechend ihrer Länge
- lese die Sequenz ab (Read)

→ Damit können Moleküle bis zu 1000 bp lang sequenziert werden.

8.1. Shotgun-SequenzierungMethode 8A: Shotgun-Sequenzierung:Eingabe: Ein DNA-Molekül  $D$ 

- 1) Erzeuge Kopien  $C = \{D_1, \dots, D_n\}$  von  $D$
- 2) Zerlege jede Kopie in kleine Fragmente (zufällig)  
 $\rightarrow$  Menge einander überlappende DNA-Fragmente
- 3) Ermittle die Sequenz der Fragmente (Kettenabbruch-Methode) (oder Anfangsstücke der Fragmente)  
 $\rightarrow$  Eine Menge von Strings über dem Alphabet

$$\Sigma_{DNA} = \{A, C, G, T\}$$

Ausgabe: Die Menge der ermittelten DNA-Fragment-Sequenzen  
 $S = \{s_1, \dots, s_n\}$

Definition: Fragment-Assembly-ProblemEingabe: Menge von Strings  $S = \{s_1, \dots, s_n\}$ Ausgabe: Anordnung der Strings, die der ursprünglichen Anordnung in dem Molekül entspricht.Beispiel: DNA der Länge  $L = 100$  kbpFragmentanzahl  $n = 1500$ durchschnittliche Länge der Fragmente  $f = 500$  bpDaher:  $n \cdot f = 750$  kbpdurchschnittliche Überdeckung / Coverage:  $7,5 = \frac{n \cdot f}{L}$

Schem zur Lösung des Fragment-Assembly-Problem:

- 1) Overlap-Bestimmung:
  - Überlappungen zwischen den einzelnen Strings ermitteln
  - Überlappungen müssen nicht notwendig Suffiz-Präfix-Paar sein  $\rightarrow$  Alignment

2) Layout: Anordnung der Strings  $\rightarrow$  semiglobales multiples Alignment

3) Consensus: Bestimmung der Sequenz aus dem Layout.

$\rightarrow$  Ziel: Modellbildung für die Layout-Phase.

8.1.1. Fehlerquellen und Probleme beim Fragment-Assembly

- Sequenzierungsfehler:

- Einfügung / Insertion: Einfügen einer Base wo diese nicht vorhanden ist

AGTAXCA

- Löschung / Deletion: Löschen einer Base

AGTATTGCA

- Ersetzung / Substitution: Ersetzen einer Base durch eine andere

ACCTGAAC



- Chimer: Fragmente, die aus unterschiedlichen Bereichen der Ursprungs-DNA stammen und sich zusammenschließen.





- unvollständige Überdeckung des Ursprungs-DNA

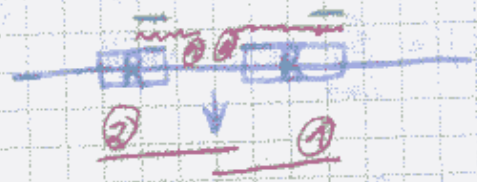


z.B. das Fragment toxische Wirkung auf einen Gastorganismus hat.

- Unbekannte Orientierung der Fragmente

- Repeats: Teilstrings der DNA, die an mehreren verschiedenen Stellen (identisch / fast identisch) auftreten. Große Variation über die Länge, Häufigkeit, Ähnlichkeit der Repeats.

Problematik: - Überlappung wegen Repeats



- Anordnung der Repeats unklar

- Länge der Ausgangs-DNA ins Modell einbinden

⇒ ideale Daten ≙ keine der vorangegangenen Zelltypen folgt auf.

8.1.2. Shortest-Common-Superstring-Problem

Def. 8.3: Das Shortest-Common-Superstring-Problem (SCS) ist das folgende Optimierungsproblem.

Eingabe: Eine Menge  $S = \{s_1, \dots, s_n\}$  von Strings über einem Alphabet  $\Sigma$ .

Zulässige Lösungen: Jeder Superstring  $w$  von  $S$ , d.h.  $w$  enthält alle  $s_i$  ( $i \in \{1, \dots, n\}$ ) als Teilstring

Kosten: Länge von  $w$ ,  $\text{cost}(w) = |w|$

Optimierungsziel: Minimierung

Def. 8.3: Eine Menge  $S$  von Strings heißt Teilstringfrei, wenn kein Paar  $s, t \in S$  existiert mit  $s \neq t$  und  $s$  ist Teilstring von  $t$ .

Def. 8.4: Der triviale Superstring  $w_T$  einer Menge  $S = \{s_1, \dots, s_n\}$  entspricht der Konkatenation aller Strings in  $S$ , also

$$w_T = s_1 \cdot s_2 \cdot \dots \cdot s_n$$

Die Länge des trivialen Superstrings

$$|w_T| = \sum_{i=1}^n |s_i| = |S|$$

Def. 8.5: Sei  $w$  ein Superstring einer Menge  $S = \{s_1, \dots, s_n\}$ . Die Kompression von  $w$  ist definiert als

$$\text{comp}(w, S) = |S| - |w|$$

Def. 8.6: Das Maximum-Kompression-Längen-Superstring-Problem (MKSS)

Eingabe:  $S = \{s_1, \dots, s_n\}$  von Strings

Zulässige Lösungen: Jeder Superstring  $w$  von  $S$ .

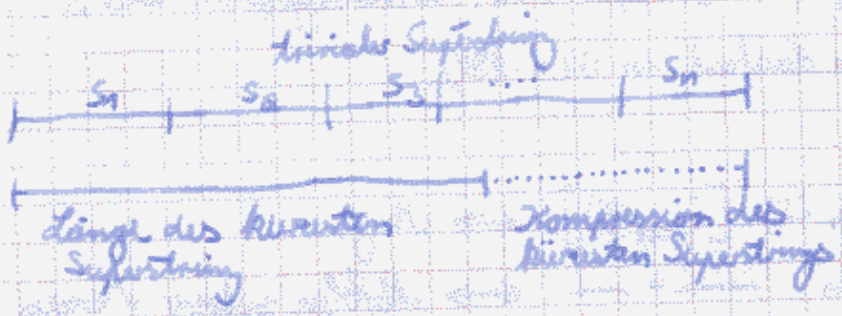
Kosten: Die Kompression des Superstrings

$$\text{cost}(w) = \text{comp}(w, S)$$

Optimierungsziel: Maximierung

SCS: Shortest-Common-Superstring-Problem mit Längenmaß

MCCS: Shortest-Common-Superstring-Problem mit Kompressionsmaß



Zur Erinnerung: Seien  $s, t$  Strings

- die Menge  $\langle s, t \rangle$  von  $s$  und  $t$  die Verschmelzung zweier Strings mit maximalem Überlappungsblock,  $s = uv$ ,  $t = vw$ ,  $v$  maximal  $\langle s, t \rangle = uvw$

- Overlap,  $ov(s, t) = v$ ,

$ov(s, t) = |v|$

- Prefix,  $pref(s, t) = u$ ,

$pref(s, t) = |u|$  (Distanz)

(- Suffix,  $suff(s, t) = w$ ,

$suff(s, t) = |w|$ )

$\rightarrow pref(s, t) = |s| - ov(s, t)$

Def. 8.7:  $S = \{s_1, \dots, s_n\}$  von Strings

Overlap-Graph  $G_{ov}(S)$ , gerichteter, gewichteter, vollständiger Graph  $(V, E, c)$

$V := S, E := V^2$

$c: E \rightarrow \mathbb{N}$  mit  $c(s_i, s_j) = ov(s_i, s_j) \quad \forall s_i, s_j \in V$

analog: Distanz-Graph  $G_{pref}(S)$

$V := S, E := V^2$

$c: E \rightarrow \mathbb{N}$  mit  $c(s_i, s_j) = pref(s_i, s_j)$

Satz 8.1: Sei  $S = \{s_1, \dots, s_n\}$  Menge von Strings über  $\Sigma$ .

Die Erstellung eines Overlap- bzw. Distanzgraphen aus  $S$  benötigt eine Laufzeit in  $O(n \cdot |S|)$ .

- Kante in  $G_{ov}(S)$  [ $G_{pref}(S)$ ]  $\leftrightarrow$  Merge
- Kantenfolge in  $\xrightarrow{\cdot} \xrightarrow{\cdot} \xrightarrow{\cdot} \leftrightarrow$  Abfolge mehrerer Merges  
 $(s_1, s_2, s_3, \dots, s_n)$  gerichteter Pfad in  $G_{ov}(S)$  [ $G_{pref}(S)$ ]

$$\begin{matrix} \updownarrow \\ \langle s_1, s_2, s_3, \dots, s_n \rangle = \text{Pref}(s_1, s_2) \text{Pref}(s_2, s_3) \dots \text{Pref}(s_{n-1}, s_n) s_n \end{matrix}$$

Def 8.3: Sei  $S = \{s_1, \dots, s_n\}$  eine Menge von Strings

$\Pi$  eine Anordnung der Strings

dann ist die durch  $\Pi$  induzierte Superstring

$$\text{LSP} = \langle \Delta_{i_1}, \Delta_{i_2}, \dots, \Delta_{i_n} \rangle$$

$$(\Pi = (i_1, i_2, \dots, i_n))$$

$\rightarrow$  Zurückführung des KCS auf TSP

$$\text{Opt}_{TSP} \leq \text{Opt}_{SC} - \text{ov}_{max} \leq \text{Opt}_{SC}$$

Satz 8.2: Die Entscheidungsvariante für das KCS ist NP-vollständig!

$\rightarrow$  Reduktion auf das HK-Problem.

## Biolnf (77) Algorithmus 8.1: Greedy-Superstring

Eingabe: Eine Menge von Strings  $S = \{s_1, \dots, s_n\}$

while  $|S| > 1$

1) Bestimme  $s_i, s_j \in S$ ,  $s_i \neq s_j$  mit dem maximalen Overlap aller Strings in  $S$ .

2) Sei  $s' = \langle s_i, s_j \rangle$  der Menge von  $s_i, s_j$

3) Lösche  $s_i, s_j$  aus  $S$  und füge  $s'$  hinzu

Ausgabe: den einzigen verbleibenden String  $s_{\text{greedy}} \in S$ .

30.6.03

### Greedy-Algorithmus bzgl. Overlap-Graphen

- wähle immer die Kante  $(s_i, s_j)$  mit dem größten Gewicht, die keine Schleife ist.

- verschmelze die Knoten  $s_i$  und  $s_j$  zu einem einzigen Knoten  $\langle s_i, s_j \rangle$

- lösche alle ausgehenden Kanten von  $s_i$  und alle eingehenden Kanten von  $s_j$

⇒ bis nur noch ein Knoten verbleibt.

... noch abstrakter:

⇒ Hamiltonischer Pfad im Overlap-Graphen

### Beispiel 8.3:

$S = \{ababaa, cabo, aabdd, aabca, aacab\}$

$\langle cabo, ababaa \rangle = cababaa$

$\{aabdd, aabca, aacab, cababaa\}$

$\langle aacab, cababaa \rangle = aacababaa$

$\{aabdd, aabca, aacababaa\}$

$\langle aacababaa, aadabd \rangle = aacababaaadabd$

$\{aa, bca, aaca, ba, baaadabd\}$

$\langle aa, bca, aaca, ba, baaadabd \rangle = aabcaacaabaaadabd$

Länge 16, Kompression um 9

Satz 8.3: Der Algorithmus Greedy-Superstring benötigt für eine Eingabe  $S = \{s_1, \dots, s_n\}$  eine Laufzeit in  $O(n \cdot |S|)$ .

Beweis: (verwende die Anweisung der Berechnung eines Hamiltonischen Pfades im Overlap-Graphen)

$O(n \cdot |S|)$  {

- Erstellung des Overlap-Graphen ( $\rightarrow$  Satz 8.1)  
 $\Rightarrow O(n \cdot |S|)$
- Sortierung der Kanten entsprechend ihres Gewichts (Overlap)  
 ~~$O(n \log n)$~~   $\Rightarrow$  mit Bin-Sort  $\underbrace{O(n^2 + |S|)}_{\leq O(n \cdot |S|)}$
- $\Rightarrow$  Adjazenzmatrix und geordnete Liste der Kanten  
 $O(n^2 + |S|) \Rightarrow O(n \cdot |S| + |S|) \Rightarrow O(n \cdot |S|)$

zu Schritt 1: Zugriff auf die Kante mit max. Gewicht  $\Rightarrow O(1)$

zu Schritt 2 und 3: gewähltes Paar  $(s_i, s_j)$

$O(n)$  {

- Verbot von ausgehenden Kanten von  $s_i$   
 $\rightarrow$  markiere die  $i$ -te Zeile
- Verbot von eingehenden Kanten von  $s_j$   
 $\rightarrow$  markiere die  $j$ -te Spalte
- Verbot eines Zykels:  
 $\rightarrow$  markiere die Position  $(j, i)$  in der Adjazenzmatrix

- lösche markierte Elemente vom Beginn der geordneten Liste

$\Rightarrow O(n)$

- (Protokolliere den durchgeführten Merge  $O(1)$ )

$\Rightarrow$  die Schritte 1 bis 3 werden  $O(n)$ -mal durchgeführt

$\Rightarrow$  Insgesamt:  $O(n \cdot |S|) + O(n) \cdot O(n)$   
 $= O(n \cdot |S|)$

□

Beispiel 8.4:

$S = \{c(ab)^m, (ba)^m, (ab)^m c\}$

Greedy:  $\langle c(ab)^m, (ab)^m c \rangle = c(ab)^m c$

$\langle (ba)^m, c(ab)^m c \rangle = (ba)^m c (ab)^m c$

$2m+2+2m = 4m+2$

Optimal:

$cababab \dots ab$   
 $babab \dots aba$   
 $abab \dots ababc$

$\Rightarrow ca(ba)^m bc$  Länge  $2m+4$

Approximationsgüte =  $\frac{4m+2}{2m+4} \xrightarrow{m \rightarrow \infty} 2$

Satz 8.4: Der Algorithmus Greedy-Superstring ist ein 4-Approximationsalgorithmus für das SCS. □

Satz 8.5: Der Algorithmus Greedy-Superstring ist ein 2-Approximationsalgorithmus für das MGS. □

Satz 8.6: Der Algorithmus Greedy-Superstring ist ein 3-Approximationsalgorithmus für das MGS. □

Beweis:  $w_{opt} = \langle S_{i_1}, \dots, S_{i_n} \rangle$ ,  $Comp(w_{opt}) \triangleq \sum$  der Anordn. der Merges

$$w_{greedy} = \langle S_{j_1}, \dots, S_{j_n} \rangle$$

• Ein Merge von Greedy kann maximal 3 Merges der optimalen Lösung verhindern!

$$m = \langle S_{j_k}, S_{j_{k+1}} \rangle$$

1)  $m$  tritt auch in der optimalen Lösung auf.  
 $\Rightarrow$  kein Merge verhindert

$$2) m = \langle S_{i_k}, S_{i_{k+1}} \rangle$$

$\Rightarrow$  zwei Merges werden verhindert, nämlich

$$\langle S_{i_k}, S_{i_{k+1}} \rangle \text{ und } \langle S_{i_{k+1}}, S_{i_{k+2}} \rangle$$

$$3) m = \langle S_{i_{k-1}}, S_{i_k} \rangle$$

$\Rightarrow$  drei Merges werden verhindert, nämlich

$$\langle S_{i_{k-1}}, S_{i_k} \rangle, \langle S_{i_k}, S_{i_{k+1}} \rangle, \text{ ein}$$

$$\text{Merge der Form } \langle S_{i_k}, S_{i_{k+1}} \rangle \dots \langle S_{i_{k+n-1}}, S_{i_{k+n}} \rangle$$

(Zykel)

$\rightarrow$  Da Greedy die Knoten mit max. Anordn. wählt, ergibt sich ein Superstring, der mindestens  $\frac{2}{3}$  des Anordn. der optimalen Lösung leistet.  $\square$

Def. 8.10: Sei  $G = (V, E, c)$  ein vollständiger, gerichteter, gewichteter Graph mit einer Gewichtsfunktion  $c: E \rightarrow \mathbb{N}$ .

Ein **Cycle-Cover**  $\mathcal{C}$  von  $G$  besteht aus einer Menge gerichteter Kreise in  $G$ ,  $\mathcal{C} = \{C_1, \dots, C_k\}$ , so dass jeder Knoten in  $G$  in genau einem Kreis  $C_i$  vorkommt.



Die Kosten  $\text{cost}(C)$  eines Cycle-Covers  $C$  mit  $C = \{C_1, \dots, C_k\}$  entsprechen der Summe des Kantengewichts in den einzelnen Kreisen

$$\text{cost}(C) := \sum_{i=1}^k \sum_{e \in C_i} c(e)$$

Ein minimaler  $C_{\min}$  ist ein Cycle-Cover mit minimalen Kosten.

→ Berechnung min. Cycle-Cover ist in polynomieller Zeit möglich.

Algorithmus 8.2 Cycle-Cover-Superstring-Algorithmus

Eingabe:  $S = \{s_1, \dots, s_n\}$ ,  $G_{\text{pref}}(S)$  der zugehörigen Präferenzgraphen.

- 1) Berechne den min. Cycle-Cover  $C$  von  $G_{\text{pref}}(S)$ .
- 2) Für jedes Kreis  $C_i \in C$  wähle (beliebig) einen Repräsentanten  $\tau_i$ .  
Die Menge der Repräsentanten  $R = \{\tau_i \mid C_i \in C\}$ .
- 3) Bestimme den durch  $R$  induzierten Teilgraphen  $G'$  von  $G_{\text{pref}}(S)$ .
- 4) Berechne min. Cycle-Cover  $C'$  auf  $G'$ .
- 5) In jedem Kreis  $C' = (c'_1, \dots, c'_k) \in C'$  lösche die Kante, die ein Overlap-Graphen das minimale Gewicht in  $C'$  besitzt.

Sei o.B.d.A.  $(c'_i, c'_j)$  diese Kante.

Beschreiben wir die Strings in dem entstehenden Pfad

$$\forall C' \in C' : \tau_i = \langle c'_1, \dots, c'_k \rangle$$

6) Kombiniere alle diese Strings  $\tau_i$  für alle  $C' \in C'$  und beschrifte den resultierenden String mit  $w$ .

7) Sei  $C = (C_1, \dots, C_k) \in C$ . o.B.d.A. sei  $\tau_i = c_1$ .  
Früher jeden Repräsentanten  $\tau_i$  in  $w$  durch die Kombination aller Präfixe im Kreis  $C_i$ , also durch

$$\text{Pref}(\tau_i, c_2) \text{Pref}(c_2, c_3) \dots \text{Pref}(c_{k-1}, c_k) \text{Pref}(c_k, \tau_i) \tau_i$$

Brotf (82) Berechne den resultierenden String mit  $w$ .

Ausgabe: Die Superstring  $v$  von  $(S)$  kann (D),  $\text{best} = \text{min}$

Satz 8.7: Algorithmus Cycle-Cover-Superstring ist ein 3-Approximationsalgorithmus für SS.

27.2003

Beweis: Skizze:

- (i) Kosten min Cycle-Covers auf dem Distanzgraphen  $\leq$  Länge des kürzesten Superstrings
- (ii) Abschätzen der Länge des Strings  $w$  (Schritt 5 Alg.)
- (iii) Abschätzen der Länge von  $w$  bzgl. Cycle-Cover  $C'$  auf  $G'$  und dem Overlap der gelöschten Kanten
- (iv) Abschätzen der Länge von  $w$  bzgl. Cycle-Cover  $C'$ , Cycle-Cover  $C$  und dem Overlap der gelöschten Kanten
- (v) Abschätzen des Overlaps durch den Cycle-Cover  $C$

Berechnungen:

- $\text{Opt}_{SS}(S) \hat{=}$  Länge des kürzesten Superstrings für  $S$
- $\text{Opt}_{CC}(G) \hat{=}$  Kosten eines optimalen Cycle-Covers für  $G$
- $C$ : Kreise im Cycle-Cover  $C$  von  $G$
- $C'$ : Kreise im Cycle-Cover  $C'$  von  $G'$
- $\text{min-ov}_C \hat{=}$  Kosten des minimalen Overlap in  $C \in C$
- $\text{min-ov}_{C'} \hat{=}$  Summe der  $\text{min-ov}_C \forall C \in C'$
- $\text{sum-ov}_C \hat{=}$  Summe aller Overlaps in einem Kreis  $C \in C$
- $\text{sum-ov}_{C'} \hat{=}$  Summe der  $\text{sum-ov}_C \forall C \in C'$

Böhrf 83

Zu (i): Kreis  $(v_{i_1}, \dots, v_{i_k}) \rightarrow$  String  $\langle v_{i_1}, \dots, v_{i_k} \rangle$

$$\text{cost}((v_{i_1}, \dots, v_{i_k})) = \text{pref}(v_{i_1}, v_{i_2}) + \dots + \text{pref}(v_{i_{k-1}}, v_{i_k}) + \text{pref}(v_{i_k}, v_{i_1})$$

$$\text{cost}(\langle v_{i_1}, \dots, v_{i_k} \rangle) = \text{pref}(v_{i_1}, v_{i_2}) + \dots + \text{pref}(v_{i_{k-1}}, v_{i_k}) + |v_{i_k}|$$

da  $\text{pref}(v_{i_k}, v_{i_1}) \leq |v_{i_k}|$

$$\Rightarrow \text{cost}(C) = \text{Opt}_{\text{cc}}(G) \leq \text{Opt}_{\text{scs}}(S) \quad (1)$$

$$\Rightarrow \text{cost}(C) = \text{Opt}_{\text{cc}}(G) \leq \text{Opt}_{\text{scs}}(R) \quad (2)$$

da  $G'$  Teilgraph von  $G$  ist:

$$\Rightarrow \text{cost}(C') \leq \text{cost}(C) \quad (3)$$

Zu (ii): Auflösen eines Kreises  $c' = (c'_1, \dots, c'_k) \in C$  an einer Kante  $(c'_k, c'_1)$  mit min. Overlap.

$$\begin{aligned} |u_{c'}| &= \text{pref}(c'_1, c'_2) + \dots + \text{pref}(c'_{k-1}, c'_k) + |c'_k| \\ &= \underbrace{\text{pref}(c'_1, c'_2) + \dots + \text{pref}(c'_{k-1}, c'_k) + \text{pref}(c'_k, c'_1)}_{= \text{cost}(c')} - \text{pref}(c'_k, c'_1) \\ &= \text{cost}(c') + \underbrace{|c'_k| - \text{pref}(c'_k, c'_1)}_{= \text{ov}(c'_k, c'_1)} \\ &= \text{cost}(c') + \text{min-ov}_{c'} \end{aligned}$$

$$\Rightarrow |u_{c'}| = \text{cost}(c') + \text{min-ov}_{c'} \quad (4)$$

zu (iii):  $w'$  ergibt sich als Kombination aller  $u_i'$

$$\Rightarrow |w'| \leq \underbrace{\sum_{c \in C} \text{cost}(c)}_{\text{cost}(C)} + \underbrace{\sum_{c \in C} \min_{e \in c} \text{cost}(e)}_{\min_{e \in E} \text{cost}(e)} \quad (5)$$

zu (iv): Ersetzen  $\tau_c$  in  $w'$  durch

$$\tau_c \rightarrow \text{Pref}(v_c, c_1) \dots \text{Pref}(c_{i-1}, c_i) \text{Pref}(c_i, \tau_c) \tau_c$$

diese Länge entspricht dem Kosten des  
Kreises  $\text{cost}(c)$

$\Rightarrow$  Verlängerung um  $\text{cost}(c) \quad \forall c \in C$

$$\Rightarrow |w'| \leq |w'| + \text{cost}(C)$$

$$\stackrel{(5)}{\leq} \text{cost}(C) + \min_{e \in E} \text{cost}(e) + \text{cost}(C)$$

$$\stackrel{(6)}{\leq} 2 \cdot \text{cost}(C) + \min_{e \in E} \text{cost}(e) \quad (6)$$

zu (v): Jeder Kreis besteht aus mindestens 2 Kanten

$$\Rightarrow \min_{e \in E} \text{cost}(e) \leq \frac{1}{2} \cdot \sum_{e \in E} \text{cost}(e) \quad (7)$$

$$\Rightarrow |w'| \leq 2 \cdot \text{cost}(C) + \frac{1}{2} \cdot \sum_{e \in E} \text{cost}(e)$$

Hilfssatz:

Lemma 5.1: Seien  $c_1$  und  $c_2$  Kreise in einem minimalen Cycle-Cover und  $s_1 \in c_1$  und  $s_2 \in c_2$  zwei Strings in diesen Kreisen. Dann gilt:

$$\text{cost}(c_1, s_1, s_2) \leq \text{cost}(c_1) + \text{cost}(c_2)$$

(ohne Beweis)

□

Da alle Knoten in einem Kreis  $C \in C'$  zu unterschiedlichen Kreisen in  $C$  gehören, folgt:

$$\sum_{e \in C} \text{cost}(e) \leq 2 \cdot \text{cost}(C) \quad (5)$$

mit (5) und (1) folgt:

$$\underline{|W|} \leq 2 \cdot \text{cost}(C) + \frac{1}{2} \cdot \sum_{e \in C} \text{cost}(e)$$

$$\stackrel{(8)}{\leq} 3 \cdot \text{cost}(C)$$

$$\stackrel{(1)}{\leq} \underline{3 \cdot \text{Opt}_{SS}(C)}$$

□

### 8.13 Das Reconstruction-Modell

Ziel: Einberingung von Sequenzierungsfehlern und der unbekanntem Orientierung.

Def. 8.11: Seien  $s, t$  Strings über  $\Sigma$ . Dann ist die Teilstring-Edit-Distanz  $\text{eds}(s, t)$  definiert als

$$\text{eds}(s, t) = \min_{x \in \text{Substr}(t)} \text{ed}(s, x)$$

wobei  $\text{Substr}(t)$  die Menge aller Teilstrings von  $t$  und  $\text{ed}$  die übliche Edit-Distanz (siehe Abjunkt)

Def. 8.12: Das Reconstruction-Problem

Eingabe:  $S = \{s_1, \dots, s_n\}$ ,  $\epsilon \in [0, 1]_{\mathbb{R}}$  (Fehlertoleranzwert)

zulässige Lsg: Jeder String  $w$ , so dass für alle  $i=1, \dots, n$

$$\min\{\text{eds}(s_i, w), \text{eds}(w, s_i)\} \leq \epsilon \cdot |s_i|$$

wobei  $\bar{s}_i$  das rechte Komplement von  $s_i$  ist.

Kosten: Die Länge von  $w$ .

Ziel: Minimierung

Satz 8.5 Das Reconstruction-Problem ist NP-schwer. □

## 8.2 Sequenzierung durch Hybridisierung

Methode 8.2: Sequenzierung durch Hybridisierung (SBH)

Eingabe: Die zu sequenzierende DNA und  $l \in \mathbb{N}$

- 1) Erzeuge einen DNA-Chip mit allen verschiedenen Proben der Länge  $l$ .
- 2) Erzeuge Kopien der DNA.
- 3) Führe Hybridisierungsexperiment durch

Ausgabe: Menge  $S = \{s_1, \dots, s_n\} \subseteq \Sigma_{DNA}^l$  der Länge  $l$ ,  
 (die als Teilstrings in der DNA auftreten.  
 → Spektrum (der DNA)

Def 8.16: Sei  $w$  ein String und  $S = \{s_1, \dots, s_n\}$  ein Spektrum mit Strings der Länge  $l$ .

$w$  kompatibel zu  $S$ , falls  $w$  jeden String aus  $S$  als Teilstring enthält und keinen anderen String der Länge  $l$  enthält.

$w$  einfach-kompatibel zu  $S$ , falls  $w$  kompatibel zu  $S$  und  $w$  kein String in  $S$  mehrfach enthält.

↳ Ziel: Finde einen kompatiblen String  $w$ .

Def. 8.18:

Sei  $S$  ein Spektrum mit Strings der Länge  $l$ .

Spektrum-Graph  $G_{\text{Spektrum}}(S) = (V, E, \text{label})$

- $V := \Sigma^{l-1}$  (alle Strings der Länge  $l-1$ )

- $E = \{(x, y) \mid x, y \in V, \text{ es ex. ein } s \in S \text{ mit } \langle x, y \rangle = s\}$

- $\text{label}(x, y) = y_{l-1}$

Sei  $x_1, x_2, \dots, x_k$  ein Pfad in  $G_{\text{Spektrum}}(S)$

$\text{pathlabel}(x_1, x_2, \dots, x_k) = x_1 \text{label}(x_1, x_2) \dots \text{label}(x_{k-1}, x_k)$

Eulerischer Pfad := Ein Pfad in einem Graphen, der jede Kante genau einmal durchläuft (Pfad nicht notwendig knotendisjunkt)

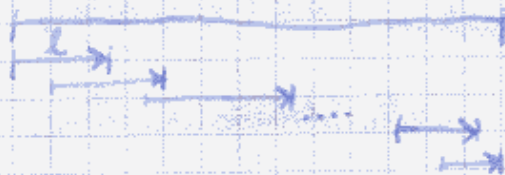
Eulerischer Kreis := Ein Kreis in einem Graphen, der jede Kante genau einmal durchläuft.

Graph  $G$ , der einen Eulerischen Kreis enthält, heißt eulerisch.

Satz 8.3: Sei  $S$  ein Spektrum mit Strings der Länge  $l$ . Ein String  $w$  ist genau dann einfach-kompatibel zu  $S$ , wenn er der Beschriftung eines Eulerischen Pfades im Spektrum-Graphen entspricht.

Beweis:  $t \in S$ ,  $t = x \text{ label}(i, y, l)$   $x, y \in V$

$\Rightarrow w$  ein einfach-kompatibler String.



„Jeder Block der Länge  $l$  in  $w$  entspricht einer Kante in  $\mathcal{L}_{\text{Spektrum}}(S)$ “

$\Leftarrow$  Das Pathlabel eines Euler- Pfades entspricht einem einfach-kompatiblen String.  $\square$

$\Rightarrow$  Eulerkreise / Eulerpfade sind effizient berechenbar.

$\Rightarrow$  es ist effizient möglich einfach-kompatible Strings zu einem Spektrum zu finden.

- Fehlerbehandlung?

- große Anzahl verschiedener Eulerpfade!

7.7.2003

## 9. Bestimmung von Signalen in DNA-Sequenzen

Ziel: Finde „interessante“ Regionen in der DNA, z.B.

- Restriktionsstellen

- Gene

- Bindungsstellen von Proteinen

- ...

Ansätze: meist statistische Methoden



### 3.1. Gleiche und ähnliche Teilstrings

Ziel: Bestimmung von Bindungsstellen, d.h. von Regionen der DNA, an denen sich ein bestimmtes Protein anlagern kann, um z.B. die Transkription eines Gens zu steuern.

Biologisches Experiment: liefert DNA-Fragment der Länge  $m$ , in dem sich (mit hoher W'keit) eine Bindungsstelle der Länge  $l \ll m$  befindet.

Problem: Finde in dem DNA-Fragmenten mehrere solche Experimente einen möglichst langen gemeinsamen Teilstring (magisches Wort).

Exakte Lösung: mit Suffix-Bäumen (Kgl. 4.10)

Probleme:

- komplizierter Muster: Restriktionsstelle des Restriktions-Enzyms XcmI: CCA \*\*\*\*\* \*\*\* TGG
- Mismatches: Finde magisches Wort, das in allen DNA-Fragmenten ungefähr vorkommt

⇒ Spezialfall eines lokalen multiplen Alignment ohne Lücken.

Def. 3.1: Seien  $S = s_1, \dots, s_m$  und  $t = t_1, \dots, t_m$  zwei Strings der Länge  $m$ . Der Hamming-Abstand  $d_H(s, t)$  von  $s$  und  $t$  ist definiert als die Anzahl der Positionen  $1 \leq i \leq m$  mit  $s_i \neq t_i$ .

Def. 3.2: Das Consensus-String-Problem:

Eingabe: Menge von  $n$  Strings  $\{s_1, \dots, s_n\} \in \Sigma^m$ ,  $l \in \mathbb{N}$

zulässige Lösungen: Alle  $(n+1)$ -Tupel  $(t, t_1, \dots, t_n)$ , wobei  $t_i$  Teilstrings der Länge  $l$  von  $s_i$  für  $1 \leq i \leq n$ . der String  $t \in \Sigma^l$  wird Median-String genannt.

Kosten:  $\text{cost}(t, t_1, \dots, t_n) = \sum_{i=1}^n d_H(t, t_i)$

Optimierungsziel: Minimierung

Satz 3.1: Das Consensus-String-Problem ist NP-schwer. □

⇒ Approximations-Algorithmus:

Algorithmus 3.1 Consensus-String-Approximation

Eingabe: Menge  $S = \{s_1, \dots, s_n\} \in \Sigma^m$ ,  $l, r \in \mathbb{N}$

1. Initialisierung:  $c' := \infty$

$u' := \emptyset$

for  $i := 1$  to  $n$  do  $v_i' := \emptyset$

2. for all  $(u_1, \dots, u_r)$ , wobei  $u_i \in \Sigma^l$  für alle  $i$  ein Teilstring eines Strings in  $S$  ist do

Berechne  $u$  als den Consensus von  $u_1, \dots, u_r$  (Def. 5.11)

for  $i := 1$  to  $n$  do

Berechne  $v_i$  als den Teilstring von  $s_i$  mit minimalen Hamming-Abstand zu  $u$

$$c := \sum_{i=1}^n d_H(u, v_i)$$

if  $c < c'$  then

$c' := c$

$u' := u$

for  $i := 1$  to  $n$  do  $v_i' := v_i$

Ausgabe:  $(u', v_1', \dots, v_n')$  mit Kosten  $c'$

Satz 3.2: der Alg. 3.1 ist für jedes  $\tau \geq 3$  ein polynomiales Approximationsalg. für das Consensus-String-Problem mit einer Approximationsgüte von

$$1 + O\left(\sqrt{\frac{\log \tau}{\tau}}\right)$$

und einer Laufzeit in  $O\left((m-l+1)^{\tau+1} \cdot n^{\tau+1} \cdot l\right)$

Beweisidee: Laufzeit: Es gibt  $n$  Strings in  $S$  und  $m-l+1$  mögliche Anfangspositionen in jedem der Strings aus  $S$  für jedes  $u_i$

$\Rightarrow (n \cdot (m-l+1))^\tau$  Möglichkeiten,  $(u_1, \dots, u_\tau)$  zu wählen.

Approximationsgüte: Sei  $S = \{s_1, \dots, s_n\} \in \Sigma^m$ ,  $l, \tau \in \mathbb{N}$ ,  $\tau \geq 3$

Sei  $(s, t_1, \dots, t_n)$  eine opt. Lösung mit Kosten

$$c_{\text{opt}} = \sum_{i=1}^n d_H(s, t_i)$$

Für alle  $(i_1, \dots, i_\tau) \in \{1, \dots, n\}^\tau$  sei  $s_{i_1, \dots, i_\tau}$  ein Consensus von  $t_{i_1}, \dots, t_{i_\tau}$  und es sei  $c_{i_1, \dots, i_\tau} = \sum_{i=1}^n d_H(s_{i_1, \dots, i_\tau}, t_i)$

Idee: Approximiere den opt. Consensus  $s$  durch ein  $s_{i_1, \dots, i_\tau}$  für ein  $(i_1, \dots, i_\tau)$ .

Man kann zeigen, dass sich für  $\tau$ -unabhängig gleichverteilt zufällig gewählte Stellen  $i_1, \dots, i_\tau$  der Erwartungswert von  $c_{i_1, \dots, i_\tau}$  wie folgt abschätzen lässt:

$$E[c_{i_1, \dots, i_\tau}] \leq \left(1 + O\left(\sqrt{\frac{\log \tau}{\tau}}\right)\right) \cdot c_{\text{opt}}$$

Baum (3)

⇒ es ist  $t_1, \dots, t_r$  mit

$$C_{i_1, \dots, i_r} \leq \left(1 + O\left(\frac{1}{\log^2}\right)\right) \cdot \text{const}$$

Der Alg. hat alle möglichen  $r$ -Tupel untersucht ⇒ Behauptung.  $\square$

### 9.3. Häufige und seltene Teilstrings

Ziel: Häufigkeit der vorkommenden Teilstrings analysieren

Idee: Teilstrings, die signifikant häufiger oder seltener auftreten, weisen auf interessante Regionen hin.

Beispiel: DNA von Bakterienlagen:

4-6 Basenpaare lange Teilstrings, die den Bindungsstellen von Restriktionsenzymen entsprechen, treten signifikant seltener auf.

Problem: Für jeden Teilstring  $t$  der Länge  $l$  in einem gegebenen String  $s$  der Länge  $n$  vergleiche die Anzahl der Vorkommen von  $t$  mit der erwarteten Anzahl von Vorkommen in einem zufälligen String der Länge  $n$ .

⇒ Berechne Erwartungswert und Varianz

Alg. 9.5: Häufigkeits-Analyse für Teilstrings

Eingabe: String  $s$  der Länge  $n$ ,  $l \in \mathbb{N}$

für alle  $t \in \Sigma^l$  do

- Bestimme die Anzahl  $h(t)$  der Vorkommen von  $t$  in  $s$ .
- Bestimme die erwartete Anzahl der Vorkommen von  $t$  in einem zufälligen String der Länge  $n$  und deren Varianz.

Ausgabe: Alle Strings  $t$ , deren tatsächliche Häufigkeit signifikant von der erwarteten Häufigkeit abweicht.

Es gilt:  $\text{Var}(X) = E[(X - E[X])^2]$

Problem: Varianz hängt von dem konstanten String selbst ab, nicht nur von seiner Länge.

Def. 9.7: Sei  $t = t_1 \dots t_\ell$  ein String. Die **Autokorrelation** von  $t$  ist ein Binärstring  $c(t) = c_0^{(t)} \dots c_{\ell-1}^{(t)}$ , wobei

$$c_i^{(t)} = \begin{cases} 1 & \text{falls } t_1 \dots t_{\ell-1-i} = t_{i+1} \dots t_\ell \\ 0 & \text{sonst} \end{cases}$$

das **Autokorrelationspolynom** von  $t$  ist definiert als:

$$\text{corr}_t(x) = \sum_{i=0}^{\ell-1} c_i^{(t)} \cdot x^i$$

Beachte:  $c_0^{(t)} = 1$  für alle  $t$

Vereinfachung: der String  $s$  sei **zyklisch**, d.h. wir zählen auch Vorkommen  $t_1 \dots t_\ell$

Satz 9.6: Sei  $\Sigma$  ein Alphabet der Größe  $k$ , sei  $s = s_1 \dots s_m$  ein **zyklischer Permuti-String** der Länge  $m$ , d.h. jedes  $s_i$  sei **zufällig gleichverteilt** und **unabhängig** aus  $\Sigma$  gewählt.

Sei  $t \in \Sigma^\ell$  ein in  $s$  **auftretendes Muster**.

Sei für  $1 \leq i \leq m$  die **zufallsvariable**  $X_i$  def. durch

$$X_i = \begin{cases} 1 & \text{falls } t \text{ an Position } i \text{ von } s \text{ beginnt} \\ 0 & \text{sonst} \end{cases}$$

Dann ist die **Anzahl** der Vorkommen von  $t$  in  $s$  durch

Bolnt (34)

die Zufallsvariable  $X = \sum_{i=1}^m X_i$  gegeben.

Sei  $p = E[X_i] = \frac{1}{\lambda}$ . Dann gilt

$$(a) E[X] = m \cdot p$$

$$(b) \text{Var}[X] = p \cdot m \cdot \left( 2 \cdot \text{corr}_X \left( \frac{1}{\lambda} \right) - (2\lambda - 1)p - 1 \right)$$

Beweis: (a) Linearität des Erwartungswerts.  $\checkmark$

$$(b) \text{Var}[X] = E[X^2] - E[X]^2$$

$$= \sum_{1 \leq i, j \leq m} (E[X_i \cdot X_j] - E[X_i] \cdot E[X_j])$$

Sei  $d(i, j)$  der kürzeste Abstand zwischen den Positionen  $i$  und  $j$  in dem zyklischen String  $s$ .

Dann gilt:

$$\text{Var}[X] = \sum_{\substack{1 \leq i, j \leq m \\ d(i, j) \geq 2}} (E[X_i \cdot X_j] - E[X_i] \cdot E[X_j]) \quad \} S_1$$

$$+ \sum_{\substack{1 \leq i, j \leq m \\ d(i, j) = 1}} (E[X_i \cdot X_j] - E[X_i] \cdot E[X_j]) \quad \} S_2$$

$$+ \sum_{\substack{1 \leq i, j \leq m \\ d(i, j) = 0}} (E[X_i \cdot X_j] - E[X_i] \cdot E[X_j]) \quad \} S_3$$

$S_1 = 0$ , da  $X_i$  und  $X_j$  unabhängig sind.

$$S_2 = \sum_{1 \leq i \leq m} (E[X_i \cdot X_i] - E[X_i] \cdot E[X_i]) = m \cdot (p - p^2)$$

$\underbrace{\hspace{10em}}_{E[X_i] = p} \quad \underbrace{\hspace{10em}}_{p^2}$

$$S_3 = \sum_{i=1}^m \sum_{r=1}^{l-1} \sum_{d(i,r,j)=r} (E[X_i \cdot X_j] - E[X_i] \cdot E[X_j])$$

Wenn  $C_r^{(k)} = 0$ , dann ist  $X_i \cdot X_{i+r} = 0$  für alle  $i$

Wenn  $C_r^{(k)} = 1$ , dann ist es möglich, dass  $k$  in 2 anderen Pos.  $i$  und  $i+r$  beginnt:

$E[X_i \cdot X_{i+r}] =$  Produkt aus der W'keit des Auftretens von  $k$  an der Pos.  $i+r$  und der W'keit, dass  $\Delta_i \dots \Delta_{i+r-1} = k_{1 \dots r}$

$$\Rightarrow E[X_i \cdot X_{i+r}] = C_r^{(k)} \cdot p \cdot \frac{1}{k^r}$$

Für jedes  $i$  gibt es genau zwei Positionen  $j$  mit  $d(i,r,j)=r$ .

$$\begin{aligned} \Rightarrow \sum_{i=1}^m \sum_{r=1}^{l-1} \sum_{d(i,r,j)=r} E[X_i \cdot X_j] &= \sum_{i=1}^m \sum_{r=1}^{l-1} 2 \cdot p \cdot C_r^{(k)} \cdot \frac{1}{k^r} \\ &= \sum_{i=1}^m 2p \cdot \left( \text{corr}_k \left( \frac{1}{k} \right) - 1 \right) \cdot \frac{1}{k^0} \end{aligned}$$

$$\Rightarrow \sum_{i=1}^m \sum_{r=1}^{l-1} \sum_{d(i,r,j)=r} (E[X_i \cdot X_j] - E[X_i] \cdot E[X_j])$$

$$= \sum_{i=1}^m 2p \cdot \left( \text{corr}_k \left( \frac{1}{k} \right) - 1 \right) - 2(l-1) \cdot p^2$$

$$= p \cdot m \cdot \left( 2 \text{corr}_k \left( \frac{1}{k} \right) - 2 - 2(l-1) \cdot p \right)$$

$$\Rightarrow \text{Var}[X] = p \cdot m \cdot \left( 2 \cdot \text{corr}_k \left( \frac{1}{k} \right) - 2 - 2(l-1) \cdot p \right) + m(p-p^2)$$

$$= p \cdot m \cdot \left( 2 \cdot \text{corr}_k \left( \frac{1}{k} \right) - (2(l-1) \cdot p - 1) \right)$$

□

Ziel: Auffinden von sog. CG-Inseln, d.h. Bereichen in einer DNA-Sequenz, in denen CG wesentlich häufiger vorkommt als im Rest der DNA-Sequenz.

Def. 3.8: Ein Hidden-Markov-Modell (HMM) ist ein Quintupel  $\mathcal{M} = (\Sigma, Q, q_0, \delta, \eta)$ , wobei

- $\Sigma$  Alphabet
- $Q$  endl. Zustandsmenge
- $q_0 \in Q$  Anfangszustand
- $\delta$  eine  $(|Q| \times |Q|)$ -Matrix von Transitionswahrscheinlichkeiten
- $\eta$  eine  $(|Q| - 1) \times |\Sigma|$ -Matrix von Emissionswahrscheinlichkeiten

$\delta(p, q) \triangleq$  Wkt. des Übergangs von Zustand  $p$  nach Zustand  $q$

$$\delta(p, p) = 0 \text{ und } \sum_{q \in Q} \delta(p, q) = 1 \text{ für alle } p \in Q$$

$\eta(q, a) \triangleq$  Wkt., daß im Zustand  $q$  das Symbol  $a$  ausgegeben wird für  $q \in Q - \{q_0\}$ ,  $a \in \Sigma$

$$\sum_{a \in \Sigma} \eta(q, a) = 1$$

Ein Path im  $\mathcal{M}$  ist eine Folge  $\pi = q_0, q_1, \dots, q_n$  von Zuständen.



Beispiel 3.2: Würfel-Experiment

Für eine Reihe von Würfeln wird teilweise ein fairer Würfel verwendet, der jede Zahl  $\in \{1, \dots, 6\}$  mit  $1/6$  Wkt  $\%$  wirft, und teilweise ein unfairer Würfel, der 6 mit  $1/2$  Wkt  $\%$  und jede Zahl  $\in \{1, \dots, 5\}$  mit  $1/10$  Wkt  $\%$  wirft. Der Würfel wird mit  $1/2$  Wkt  $\%$  gerade oder mit  $1/2$  Wkt  $\%$  ungerade gewürfelt. Am Anfang wird mit  $1/2$  Wkt  $\%$  einer der Würfeln gewürfelt.

Lemma 3.3: Sei  $\mathcal{M} = (\Sigma, \mathcal{Q}, q_0, \delta, \eta)$  ein HMM, sei  $\pi = q_0, q_1, \dots, q_n$  ein Pfad in  $\mathcal{M}$  und sei  $x = x_1 \dots x_n \in \Sigma^n$ . dann gilt:

$$\text{Prob}[x \sim \pi] = \prod_{i=1}^n (\delta(q_{i-1}, q_i) \cdot \eta(q_i, x_i))$$

Beweis:  $\text{Prob}[\pi] = \prod_{i=1}^n \delta(q_{i-1}, q_i)$

$$\text{Prob}[x|\pi] = \prod_{i=1}^n \eta(q_i, x_i)$$

$$\text{Prob}[x \sim \pi] = \text{Prob}[\pi] \cdot \text{Prob}[x|\pi] \quad \square$$

Beispiel 3.3: Würfelpfad aus Bsp. 3.2

$$\pi = q_0, U, U, F \quad x = 666$$

$$\text{Prob}[x \sim \pi] = \prod_{i=1}^3 (\delta(q_{i-1}, q_i) \cdot \eta(q_i, x_i))$$

$$= \delta(q_0, U) \cdot \eta(U, 6) \cdot \delta(U, U) \cdot \eta(U, 6) \cdot \delta(U, F) \cdot \eta(F, 6)$$

$$= \frac{1}{2} \cdot \frac{1}{10} \cdot \frac{1}{2} \cdot \frac{1}{10} \cdot \frac{1}{2} \cdot \frac{1}{6}$$

$$= \frac{1}{1200} \approx 0.00083$$

Def. 3.3: Das HMM-Decodier-Problem:

Eingabe: HMM  $\mathcal{M} = (\Sigma, Q, q_0, \delta, \pi)$ ,  $x = x_1 \dots x_n \in \Sigma^n$

zulässige Lösungen: Alle Pfade  $\pi = q_0, q_1, \dots, q_n$  der Länge  $n$  in  $\mathcal{M}$ .

Kosten: Für zul. Lösung  $\pi$ :  $\text{cost}(\pi) = \text{Prob}[x|\pi]$

Optimierungsziel: Maximierung

Lösungsansatz: Dynamische Programmierung:

Bestimme für alle Zustände und alle Präfixe von  $x$  den wahrscheinlichsten Pfad, der in diesem Zustand endet und diesen Präfix ermittelt

Lemma 3.4: Sei  $\mathcal{M} = (\Sigma, Q, q_0, \delta, \pi)$  ein HMM und sei  $x = x_1 \dots x_n \in \Sigma^n$ .

Sei  $\sigma_q(i)$  die  $i$ -te Komponente des wahrscheinlichsten Pfades für  $x_1 \dots x_i$ , der in  $q$  endet, für alle  $q \in Q$ , für alle  $0 \leq i \leq n$ .

Dann gilt:  $\sigma_{q_0}(0) = 1$

$$\sigma_q(0) = 0 \text{ für } q \in Q - \{q_0\} \quad (3.2)$$

und für alle  $q \in Q - \{q_0\}$  und für alle  $1 \leq i \leq n$  gilt:

$$\sigma_q(i) = \pi(q, x_i) \cdot \max_{p \in Q} (\sigma_p(i-1) \cdot \delta(p, q)) \quad (3.3)$$

Beweis: (9.2) ✓

(9.3): Sei  $\pi = q_0, q_1, \dots, q_i$  der Pfad mit der höchsten Emissions-  
w'keit für  $x_1, \dots, x_i$ , der in  $q_i$  endet.

$$\Rightarrow \sigma_{q_i}(i) = \sigma_{q_{i-1}}(i-1) \cdot \delta(q_{i-1}, q_i) \cdot \eta(q_i, x_i)$$

Da  $\pi$  der Pfad mit der höchsten Emissionsw'keit für  $x_1, \dots, x_i$  ist, der  
in  $q_i$  endet, kann kein Pfad mit anderem vorletztem Zustand  
eine höhere Emissionsw'keit aufweisen.

$$\Rightarrow \sigma_{q_i}(i) = \eta(q_i, x_i) \cdot \max_{p \in Q} (\sigma_p(i-1) \cdot \delta(p, q_i)) \quad \square$$

$\Rightarrow$  Viterbi-Algorithmus

Satz 9.7: Der Viterbi-Algorithmus löst das HMM-Dekodier-Problem  
für ein HMM mit  $k$  Zuständen und einem String der  
Länge  $n$  in Zeit  $O(n \cdot k^3)$ .

Beweis: Korrektheit: Lemma 3.4.

Laufzeit Initialisierung:  $k$

Schritt 2: n.k.  $k = n \cdot k^3$

Schritt 3:  $O(n)$  □

In der Praxis: Verwendung Logarithmus der Wahrscheinlichkeiten

$\Rightarrow$  - Zahlen weniger nahe bei 0  $\rightarrow$  weniger Rundungsfehler

- Additionen statt Multiplikationen  $\rightarrow$  schneller

$$\Rightarrow \sigma_q(i) = \log \eta(q_i, x_i) + \max_{p \in Q} (\sigma_p(i-1) + \log \delta(p, q_i))$$

Modellierung des CG-Island-Problems als HMM:

$$M_{CG} = (\Sigma_{DNA}, Q, q_0, \delta, \eta)$$

$$\Sigma_{DNA} = \{A, C, G, T\}$$

$$Q = \{q_0, A^+, C^+, G^+, T^+, A^-, C^-, G^-, T^-\}$$

- Transitionsw'keiten aus Isolatoren, aber durch Ausschließen von DNA-Sequenzen in denen man die CG-Islands schon kennt.

$$w_{hid} + \rightarrow - \text{klein als } - \rightarrow +$$

$$+ \rightarrow + \text{klein als } - \rightarrow -$$

$$\delta(C^+, G^+) \gg \delta(C^-, G^-)$$

$$\eta(A^+, A) = \eta(A^-, A) = 1$$

$$\eta(C^+, C) = \eta(C^-, C) = 1$$

$$\eta(G^+, G) = \eta(G^-, G) = 1$$

$$\eta(T^+, T) = \eta(T^-, T) = 1$$

## 11. Phylogenetische Bäume

Ziel: Rekonstruktion von Verwandtschaftsbeziehungen, von z.B.

biologischen Arten oder Genen (Taxa)

→ Konstruktion eines phylogenetischen Baumes (Phylogenie)

Blätter  $\hat{=}$  Taxa

innere Knoten  $\hat{=}$  Verzweigung

Abstand im Baum  $\hat{=}$  Grad der Verwandtschaft

Meist Binärbaum

Kern des Baumes  $\hat{=}$  gemeinsame Vorfahr aller Taxa

14.7.2003

Verschiedene Modelle phylogenetischer Bäume:

- Verzweigungsgrad: meist Binärbäume
- gerichtete oder ungerichtete Bäume
- Kantenlängen spezifiziert oder nur Topologie des Baums

Verschiedene zur Verfügung stehende Daten:

- Distanzmaß, das je zwei Taxa einen Abstand zuordnet.  
z.B. Alignment-Bewertung.
- Menge von Merkmalen  
z.B. phänotypische Merkmale  
oder: multiples Alignment von Gen-Sequenzen, jede  
Position definiert ein Merkmal.

11.1. Ultrametrische Distanzen

Def. 11.1: Sei  $A$  eine Menge von Taxa, sei  $d: A \times A \rightarrow \mathbb{Q}^{\geq 0}$ .  
Dann ist  $d$  eine **Metrik auf  $A$** , wenn gilt:

- (i)  $d(a, b) = 0 \iff a = b$  für alle  $a, b \in A$
- (ii)  $d(a, b) = d(b, a)$  für alle  $a, b \in A$  (Symmetrie)
- (iii)  $d(a, b) \leq d(a, c) + d(c, b)$  für alle  $a, b, c \in A$  (Dreiecksungl.)

Def. 11.2: Sei  $A$  eine Menge von Taxa, sei  $d: A \times A \rightarrow \mathbb{Q}^{\geq 0}$  eine  
Metrik auf  $A$ . Dann ist  $d$  eine **Ultrametrik auf  $A$** , wenn  
zusätzlich die folgende **drei-Punkt-Bedingung** gilt:

Für alle  $a, b, c \in A$  sind zwei der Distanzen  $d(a, b)$ ,  $d(a, c)$ ,  $d(b, c)$  gleich und nicht kleiner als die dritte.

$$\begin{aligned} d(a, b) &\leq d(a, c) = d(b, c) && \text{oder} \\ d(a, c) &\leq d(a, b) = d(b, c) && \text{oder} \\ d(b, c) &\leq d(a, b) = d(a, c). \end{aligned}$$

Ziel: Bestimme phylogenetischen Baum mit einer Wurzel, bei dem auch die Kantenlängen bekannt sind, so daß die Pfadlängen von der Wurzel zu einem beliebigen Blatt gleich sind.

→ ultrametrischer Baum

ideales Evolutionsmodell, in dem die Evolutionsgeschwindigkeit in jedem Ast des Baumes genau gleich ist.

Def. 11.3: Sei  $A = \{a_1, \dots, a_n\}$  eine Menge von Taxa, ein gerichteter kantengewichteter Baum  $T = (V, E, d)$  mit einer Wurzel  $r$  und Kantenbewertung.

$d: E \rightarrow \mathbb{Q}^{\geq 0}$  ist ein ultrametrischer Baum für

$A$ , wenn gilt:

- (i)  $T$  ist binärer Baum.
- (ii)  $T$  hat genau  $n$  Blätter, die mit dem Taxa beschriftet sind.
- (iii) die Summe der Kantenbeschriftungen auf jedem Pfad von der Wurzel zu einem beliebigen Blatt ist gleich.

die Distanz zwischen  $x, y \in V$  ist die Summe der Kantenbeschriftungen auf dem Pfad von  $x$  nach  $y$  in  $T$ .

Bezeichnung:  $\text{dist}_T(x, y)$

Beispiel 11.1:



Lemma 11.1: Sei  $A$  eine Menge von Taxa, sei  $T = (V, E, d)$  ein ultrametrischer Baum für  $A$ . Dann definieren die Distanzen  $\text{dist}_T$  zwischen den Taxa eine Ultrametrik auf  $A$ .

Beweis: Seien  $a, b, c$  Blätter von  $T$ .

zeige, daß  $\text{dist}_T(a, b), \text{dist}_T(a, c), \text{dist}_T(b, c)$  die Dreieck-  
Bedingung erfüllen.

1. Fall:  $a = b$

$$\Rightarrow 0 = \text{dist}_T(a, b) \leq \text{dist}_T(a, c) = \text{dist}_T(b, c) \quad \checkmark$$

2. Fall:  $a, b, c$  paarweise verschieden.

$$\text{dist}_T(r, a) = \text{dist}_T(r, b) = \text{dist}_T(r, c)$$

(folgt aus (11.1))



$$\text{dist}_T(v, a) = \text{dist}_T(v, b)$$

$$\Rightarrow \text{dist}_T(a, c) = \text{dist}_T(a, v) + \text{dist}_T(v, u) = \text{dist}_T(b, v) + \text{dist}_T(v, u) = \text{dist}_T(b, c)$$

noch zu zeigen:  $\text{dist}_T(a, b) \leq \text{dist}_T(a, c)$

aus (11.1) folgt:  $\text{dist}_T(v, a) = \text{dist}_T(v, c)$

$$\begin{aligned} \Rightarrow \text{dist}_T(a, c) &= 2 \cdot \text{dist}_T(v, a) \\ &= 2 \cdot (\text{dist}_T(v, a) + \text{dist}_T(v, u)) \\ &> 2 \cdot \text{dist}_T(v, a) \\ &= \text{dist}_T(a, b) \end{aligned} \quad \square$$

andere Richtung: Gegeben eine Menge  $A$  von Taxa mit ultrametric Distanzfunktion  $d$ , finde ultrametric Baum:

- Idee: Konstruiere Baum, so daß Knoten  $\hat{=}$  Teilungen von  $A$
- Starte mit allen einelementigen Teilungen von  $A$  ( $\rightarrow$  Blätter)
  - Berechne paarweisen Abstand
  - Solange möglich, wähle zwei Menge  $X, Y$  ohne Verfahren mit minimalen Abstand, füge  $X \cup Y$  als neuen Knoten hinzu, verbinde ihn mit  $X$  und  $Y$ , und berechne Abstand zu allen anderen Knoten.
- $\rightarrow$  nach  $n-1$  Schritten entsteht Baum mit Wurzel  $A$
- $\rightarrow$  **UPGMA-Algorithmus** (Alg. 11.1)

Satz 11.1: Der Algorithmus 11.1 berechnet für eine gegebene Menge  $A = \{a_1, \dots, a_n\}$  von Taxa und eine ultrametric Distanzfunktion  $d$  auf  $A$  einen ultrametric Baum für  $A$  in einer Zeit  $O(n^3)$ .

Beweis: Korrektheit: offenbar konstruiert der Alg. einen Baum, in dem der Abstand von der Wurzel zu einem beliebigen Blatt immer gleich ist.

Zeige, daß die in 2.f) betrachteten Kontingenzwerte nicht negativ werden. Es reicht aus, für neu konstruierten Knoten  $D$  und seine Kinder  $C_1$  und  $C_2$  zu zeigen, daß  $\text{height}(D) \geq \text{height}(C_1)$  und  $\text{height}(D) \geq \text{height}(C_2)$



Sei  $D_i$  der in der  $i$ -ten Iteration von Schritt 2 neu hinzugefügter Knoten, seien  $C_{1,i}, C_{2,i}$  dessen Kinder,  $V_i$  die Menge der Knoten nach Iteration  $i$ ,  $T_i$  die Menge  $T$  nach Iteration  $i$ .

Zeige: Für alle  $1 \leq i \leq n-1$  und für alle  $X, Y \in V_i$ :  
 $\text{height}(D_i) \geq \text{height}(X)$

Vollst. Ind. über  $i$ :  $i=0, \forall$

Ind. Schritt: zeige  $\text{height}(D_{i+1}) \geq \text{height}(D_i)$

$$\Rightarrow \text{zu zeigen: } \text{dist}(C_{1,i+1}, C_{2,i+1}) \geq \text{dist}(C_{1,i}, C_{2,i})$$

1. Fall:  $D_i$  kein Kind von  $D_{i+1}$ :

$$\Rightarrow C_{1,i+1}, C_{2,i+1}, C_{1,i}, C_{2,i} \in T_{i+1}$$

Da in der  $i$ -ten Iteration in Schritt 2a)  $C_{1,i}$  und  $C_{2,i}$  als diejenigen mit minimalem Abstand gewählt wurden, folgt (M.2)

2. Fall:  $D_i$  ist Kind von  $D_{i+1}$ : o.B.d.A.  $D_i = C_{1,i+1}$

$$\Rightarrow C_{1,i+1}, C_{2,i}, C_{2,i+1} \in T_{i+1}$$

$$\Rightarrow \text{dist}(C_{2,i}, C_{2,i+1}) \leq \text{dist}(C_{1,i}, C_{2,i+1}) \text{ und}$$

$$\text{dist}(C_{1,i}, C_{2,i}) \leq \text{dist}(C_{2,i}, C_{2,i+1})$$

$$\stackrel{2a)}{\Rightarrow} \text{dist}(C_{1,i+1}, C_{2,i+1}) \leq \text{dist}(D_i, C_{2,i+1}) = \text{dist}(C_{1,i}, C_{2,i+1}) + \text{dist}(C_{2,i}, C_{2,i+1})$$

$$\geq \frac{\text{dist}(C_{1,i}, C_{2,i}) + \text{dist}(C_{2,i}, C_{2,i+1})}{2}$$

$$= \text{dist}(C_{1,i}, C_{2,i})$$

□

Laufzeit: Schritt 1:  $O(n^3)$

Schritt 2:  $n-1$  Durchläufe

pro Durchlauf: a) :  $O(n^2)$

b) :  $O(1)$

c) :  $O(n)$

d) :  $O(1)$

e) :  $O(1)$

f) :  $O(1)$

→ gesamt:  
 $O(n^3)$

$O(n^3)$

## 11.2. Additive Bäume

Problem: Reale Daten sind häufig nicht ultrametrisch

→ schwächer Voraussetzung für das Distanzmaß,  
anderes Modell: Tacke an beliebigen Knoten des Baums.

Def. 11.4: Sei  $A$  eine Menge von  $n$  Tacke, sei  $\delta: A \times A \rightarrow \mathbb{Q}^{\geq 0}$  eine Metrik auf  $A$ ,  $T = (V, E, d)$  ein kantengewichteter Baum mit  $A \subseteq V$ .

Für alle  $a, b \in A$  sei  $\text{dist}(a, b)$  die Summe der Kantengewichte auf dem Pfad von  $a$  nach  $b$  in  $T$ .

$T$  heißt **additiver Baum** für  $A$  und  $\delta$ , wenn  $\text{dist}(a, b) = \delta(a, b)$  gilt, für alle  $a, b \in A$ .

Problem: Gegeben  $A$  und  $\delta$ , ex. ein additiver Baum für  $A$  und  $\delta$ ?  
Falls ja, wie kann man ihn berechnen?

Def. 11.5: Sei  $A$  eine Menge von Java,  $\delta$  Metrik auf  $A$ .

Ein additiver Baum  $T = (V, E, d)$  für  $A$  und  $\delta$  heißt **kompakter additiver Baum** für  $A$  und  $\delta$ , falls  $V = A$ .

Def. 11.6: Das **compact-add-tree-problem** ist das folgende Berechnungsproblem:

Eingabe: Menge  $A$  von Java, Metrik  $\delta: A \times A \rightarrow \mathbb{Q}^{\geq 0}$ .

Ausgabe: Kompakter additiver Baum für  $A$  und  $\delta$ , falls er, Selbstmeldung sonst.

Def. 11.7: Sei  $A$  eine Menge von Java,  $\delta$  metrisches Distanzmaß auf  $A$ . Der **Distanz-Graph** für  $A$  und  $\delta$  ist der vollständige, kantengewichtete Graph  $G(A, \delta) = (V, E, d)$  mit  $V = A$  und  $d(a, b) = \delta(a, b)$  für alle  $a, b \in A$ .

16.7.2003

Satz 11.2: Sei  $A$  eine Menge von Java,  $\delta$  Metrik auf  $A$ . Falls ein kompakter additiver Baum  $T$  für  $A$  und  $\delta$  existiert, dann ist  $T$  der eindeutig bestimmte minimale Spannbau von  $G(A, \delta)$ .

Beweis: Sei  $T$  ein kompakter additiver Baum für  $A$  und  $\delta$ .

Zeige: keine Kante, die nicht in  $T$  enthalten ist, kann in einem minimalen Spannbau enthalten sein.

Sei  $e = \{x, y\}$  eine nicht in  $T$  enthaltene Kante.

- Pfad von  $x$  nach  $y$  in  $T$  hat Gesamtgewicht von  $\delta(x, y)$
- Alle Kantengewichte in  $T$  sind  $\geq 0$ .

$\Rightarrow$  Jede Kante auf dem Pfad von  $x$  nach  $y$  in  $T$  hat Gewicht  $< d(x, y)$ .

Annahme:  $e$  ist im minimalen Spannbaum  $T' = (A, E')$  enthalten.  
 Sei  $G' = (A, E' - \{e\})$ , seien  $S$  und  $S'$  die Zusammenhangskomponenten von  $G'$ , o.B.d.A.  $x \in S, y \in S'$ .

Sei  $e'$  die Kante auf dem Pfad  $P$  von  $x$  nach  $y$  in  $T$ , die von  $S$  nach  $S'$  verläuft.

$\Rightarrow e' \neq e, e' \notin T'$

Setze  $T'' = (A, (E' - \{e\}) \cup \{e'\})$ . Dann ist  $T''$  Spannbaum von  $G(A, \delta)$ .

$\Rightarrow T''$  hat echt geringere Kosten als  $T'$ .  $\checkmark$

□

Lemma 11.2: Sei  $G = (V, E, d)$  ein vollständiger kostenbewerteter Graph,  $u \neq v, \{u, v\}$  Kante minimaler Kosten mit  $u \in U, v \in V \setminus U$ .

Dann gibt es einen min. Spannbaum von  $G$ , der  $\{u, v\}$  enthält.

Beweis: Annahme: kein min. Spannbaum enthält  $\{u, v\}$

Sei  $T$  ein min. Spannbaum. Dann enthält  $T \cup \{u, v\}$  einen Kreis.

$\Rightarrow$  Da  $u$  und  $v$  auch in  $T$  verbunden sind, es. andere Kante  $\{u', v'\}$  mit  $u' \in U$  und  $v' \in V \setminus U$ .

$(T \cup \{u, v\}) \setminus \{u', v'\}$  ist auch ein Spannbaum von  $G$ , der nicht teurer ist als  $T$ .  $\checkmark$

□

Satz 11.3: Der Algorithmus 11.2 löst das Compact-Add-Tree-Problem für eine Menge  $A$  von  $n$  Taxa und eine Metrik  $d$  auf  $A$  in einer Zeit in  $O(n^3 \cdot \log n)$ .

Beweis: Korrektheit: In jedem Durchlauf durch die While-Schleife gilt die Aussage von Lemma 11.2, also berechnet der Alg. einen minimalen Spanbaum, falls dieser eindeutig ist.

Laufzeit: Bestimmung des Distanzgraphen:  $O(n^3)$

Sortieren des Kantengewichts:  $O(n^2 \cdot \log n)$

While-Schleife:  $O(n^3)$  □

#### 11.4. Das Parsimony-Prinzip und die Quartett-Methode

Ziel: Phylogenie bestimmen mit Hilfe der DNA-Sequenzen homologer Gene als Merkmale.

Taxa: Menge gleichlanger Strings

$\hat{=}$  DNA-Sequenzen einer Menge homologer Gene.

Methode: Spalten eines mult. Alignment dieser Strings

Ziel: Bestimme Topologie eines binären phyl. Baums ohne Wurzel, dessen Blätter den Taxa entsprechen.

Zwei Schritte: 1. Bestimme Kostenmaß für gegebene Topologie für gegebene Menge von Taxa und Merkmalen.

2. Bestimme Topologie mit geringsten Kosten.

Def. 11.12: Sei  $S = \{s_1, \dots, s_n\}$  eine Menge von Taxa. Ein ungerichteter phylogenetischer Baum für  $S$  ist ein ungerichteter Binärbaum ohne Wurzel, der genau  $n$  Blätter hat, die (bipolar) mit den Taxa aus  $S$  beschriftet sind.

Def. 11.13: Das Parsimony-Problem:

Eingabe: Menge  $S = \{s_1, \dots, s_n\}$  von Strings der Länge  $k$  über  $\Sigma$  und ein ungerichteter phylogenetischer Baum  $T = (V, E)$  für  $S$ .

Zulässige Lösungen: für eine Eingabe-Instanz ist jede Funktion  $\beta: V \rightarrow \Sigma^k$  eine zulässige Lösung, wobei die Blätter auf die in der Eingabe gegebenen Strings abgebildet werden.

Kosten: für zulässige Lösungen  $\beta$ :

$$\text{cost}(\beta) = \sum_{\{x, y\} \in E} \text{dist}_H(\beta(x), \beta(y))$$

( $\text{dist}_H$  = Hamming-Abstand)

Optimierung: Minimierung

Idee: zur Lösung des Parsimony-Problems:

- Füge an beliebiger Stelle eine Wurzel in den Baum ein.
- Durchlaufe den Baum von den Blättern her, speichere für jeden Knoten Menge möglicher Beschriftungen.
- Durchlaufe den Baum von der Wurzel aus und wähle für jeden Knoten eine der möglichen Beschriftungen aus.

- Entferne die Wurzel

⇒ Fitch-Algorithmus

Satz 11.7: Der Alg. 11.5. löst das Parsimony-Problem in linearer Zeit in  $O(n \cdot k)$ .

Beweis: Korrektheit: klar

Laufzeit: Sowohl bei dem Bottom-up-Durchlauf, wie auch bei dem Top-down-Durchlauf wird jeder Knoten für jedes Merkmal genau einmal betrachtet  $\rightarrow O(n \cdot k)$ .  $\square$

Ziel: Finde ungerichteten phylogenetischen Baum, der die Parsimony-Bewertung minimiert.

Def. 11.11: Min-Par-Top-Problem

Eingabe: Menge  $S$  von  $n$  Strings der Länge  $k$

Zulässige Lösungen: Jeder ungerichtete Phylogenetische Baum für  $S$ .

Kosten: optimale Parsimony-Bewertung

Optimierungsziel: Minimierung

Naiver Ansatz: Alle möglichen Topologien durchprobieren:  
 $\rightarrow$  exponentielle Aufwand:

Satz 11.8: Für alle  $n \geq 3$  beträgt die Anzahl der nicht isomorphen ungerichteten phylogenetischen Bäume für  $n$  Taxa

$$\prod_{i=3}^n (2i-5) = \frac{(2n-4)!}{2^{n-2} \cdot (n-2)!}$$

Bemerkung: Binärbäume ohne Wurzel mit  $n$  Blättern hat genau  $2n-3$  Knoten.

$n=3$ : genau ein Binärbaum:



Binärbaum mit  $n$  Blättern lässt sich aus Binärbaum mit  $n-1$  Blättern konstruieren:  $2(n-1) - 3 = 2n - 5$  Möglichkeiten.  $\square$

Kein effizienter Algorithmus für Min-Pair-Top-Problem bekannt

→ Heuristiken

Def. Phylogenie für Menge  $S$  von  $n$  Taxa zusammensetzen aus phylogenetischen Bäumen für Teilmengen von  $S$ .

→ Teilmengen der Größe 4

→ Quartett-Methode

Def. 11.15: Ein **Quartett** ist ein ungerichteter phylogenetischer Baum für 4 Taxa.

Für  $S = \{a, b, c, d\}$  gibt es genau 3 verschiedene Quartette.

**optimales Quartett**: Quartett mit opt. Parsimony-Bewertung.

Def. 11.16:  $S$  Menge von  $n$  Taxa,  $T$  ungerichteter phylogenetischer Baum für  $S$ , sei  $S' = \{a, b, c, d\} \in S$ ,  $Q = (a, b; c, d)$  Quartett für  $S'$ .

Sei  $P_1$  der Pfad von  $a$  nach  $b$  in  $T$ ,  $P_2$  der Pfad von  $c$  nach  $d$  in  $T$ .

$Q$  heißt **konsistent zu  $T$** , falls  $P_1$  und  $P_2$  disjunkt sind.

Satz 11.3: Sei  $S$  eine Menge von Taxa,  $T$  unger. phylog. Baum für  $S$ .

Sei  $Q_T$  die Menge aller zu  $T$  konsistenten Quartette.

Dann lässt sich  $T$  aus  $Q_T$  in polynomialischer Zeit eindeutig konstruieren.



Satz: Bestimme für jede 4-elementige Teilmenge des Taxa das Quartett mit optimaler Parsimony-Bewertung.

Dann bestimme den ungerichteten phylogenetischen Baum, der zu den meisten Quartetten konsistent ist.

Def. 11.17: Max-Quartett-Consist-Problem:

Eingabe: Menge  $S$  von  $n$  Taxa

zulässige Lösungen: jeder ungerichtete phylo. Baum  $T$  für  $S$ .

Kosten: Für einen solchen Baum  $T$  entsprechen die Kosten der Anzahl aller optimalen Quartette für Teilungen von  $S$ , die zu  $T$  konsistent sind.

Optimierungsziel: Maximierung

Satz 11.19: Das Max-Quartett-Consist-Problem ist NP-schwer.  $\square$

aber: beliebig gut approximierbar.

→ Aufwändiger Algorithmus

Zur: Heuristik: Quartett-Pursuing

Algorithmus 11.6: Quartett-Pursuing

Eingabe: Menge  $S$  von  $n$  Taxa

1. Berechne für jede 4-elementige Teilmenge  $S' \subseteq S$  das optimale Quartett  $Q(S')$ .

2. Wähle zufällige Reihenfolge  $a_1, \dots, a_n$  der Elemente in  $S$ .

3.  $T := Q(\{a_1, \dots, a_n\})$

4. für  $i := 5$  bis  $n$  do

- Initialisiere die Kosten aller Kanten in  $T$  mit 0.
- Für alle  $S = \{b_1, b_2, b_3\} \in \{a_1, \dots, a_{i-1}\}$ , so daß  $Q(\{b_1, b_2, b_3, a_i\})$  die Form  $(b_1, b_2; b_3, a_i)$  hat, erhöhe die Kantenkosten auf dem Pfad von  $b_1$  nach  $b_2$  in  $T$  jeweils um 1.
- Wähle eine Kante  $\{x, y\}$  in  $T$  mit minimalen Kosten, lösche diese und füge einen neuen Knoten ein, der mit  $x, y$  und  $a_i$  verbunden ist.

Ausgabe: der ungerichtete phyl. Baum  $T$  für  $S$ .

Beispiel 11.8:  $S = \{a, b, c, d, e\}$ ,  $(a, b; c, d)$ ,  $(a, b; c, e)$   
 $(a, d; b, e)$ ,  $(a, c; d, e)$ ,  $(b, d; c, e)$

## 10. Vergleich von Genomen

bisher: Vergleich von DNA-Sequenzen, basierend auf lokale Mutationen (löschen, einfügen oder ändern einer einzelnen Base).

jetzt: Mutationen auf höherer Ebene:

kann die DNA-Sequenz nur zwischen den einzelnen Genen

→ nur die Reihenfolge der Gene verändert

→ Genome Rearrangements

Notation: bei eng verwandten Arten sind die Gene fast identisch, aber Anordnung der Gene variiert.

Spezialfall: mitochondriale Genome  
(fast) nur Reversals als Genome Rearrangements,  
nur ein Chromosom.

⇒ Modelliere Abfolge der Gene durch Permutation

Teil: Sortiere Permutation durch Reversals

- Zwei Modelle:
- mit bekannter Leserichtung  
→ gerichtete Permutation
  - mit unbekannter Leserichtung  
→ gerichtete Permutationen

10.2. Sortieren ungerichteter Permutationen

Def. 10.1: Sei  $\pi = (\pi_1, \dots, \pi_n)$  Permutation der Ordnung  $n$

Für  $1 \leq i < j \leq n$  ist ein  $(i, j)$ -Reversal eine Permutation  $\rho(i, j)$ , so daß gilt:

$$\pi \cdot \rho(i, j) = (\pi_1, \dots, \pi_{i-1}, \pi_j, \pi_{i-1}, \dots, \pi_{i+1}, \pi_i, \pi_{i+1}, \dots, \pi_n)$$

$(i, i)$ -Reversal: identische Permutation

$(j, i)$ -Reversal =  $(i, j)$ -Reversal

Def. 10.2: Sortieren einer Permutation durch Reversals, **MinSR-Problem**

Eingabe:  $n \in \mathbb{N}$ , Permutation  $\pi = (\pi_1, \dots, \pi_n)$  der Ordnung  $n$

Zulässige Lösungen: Jede Folge  $S_1, \dots, S_k$  von Reversals, so daß  
 $\pi \circ S_1 \circ \dots \circ S_k = (1, \dots, n)$

Kosten: Anzahl  $k$  der Reversals

Ziel: Minimierung

Satz 10.1: Das MinSR-Problem ist NP-schwer.  $\square$

$\rightarrow$  Approx. alg. mit Güte 2

Def. 10.3: Sei  $\pi = (\pi_1, \dots, \pi_n)$  eine Permutation, definiere  $\pi_0 = 0$ ,  $\pi_{n+1} = n+1$

$\text{ext}(\pi) = (\pi_0, \pi_1, \dots, \pi_n, \pi_{n+1})$  erweiterte Darstellung von  $\pi$ .

Def. 10.4: Sei  $\pi = (\pi_1, \dots, \pi_n)$  Permutation. Ein **Breakpoint** von  $\pi$  ist ein  
 Paar  $(i, i+1) \in \{0, \dots, n\} \times \{1, \dots, n+1\}$  von Positionen, so daß

$$|\pi_i - \pi_{i+1}| \neq 1$$

$\text{bp}(\pi)$  Anzahl der Breakpoints von  $\pi$

Beispiel:  $0 \mid 4 \ 3 \ 2 \mid 7 \mid 1 \mid 5 \ 6 \mid 8 \ 9$

Lemma 10.1: Sei  $\pi$  eine Perm. der Ordnung  $n$ . Dann sind mindestens

$\lceil \frac{\text{bp}(\pi)}{2} \rceil$  Reversals nötig, um  $\pi$  zu sortieren.

Beweis: Die identische Permutation hat keine Breakpoints.  
 Jedes Reversal eliminiert höchstens zwei Breakpoints.

→ Ziel: Finde Folge von Reversals, die bel. Permutation sortiert und im Durchschnitt 1 Breakpoint pro Reversal eliminiert.

Def. 10.5: Sei  $\pi$  Permutation der Ordnung  $n$ . Sei  $k = \text{bcp}(\pi)$ , seien  $(i_1, i_1+1), \dots, (i_k, i_k+1)$ ,  $i_1 < \dots < i_k$  die Breakpoints von  $\pi$ .

Dann bilden die  $k+1$  Folgen

$$S_0 = (\pi_0, \dots, \pi_{i_1}), \quad S_1 = (\pi_{i_1+1}, \dots, \pi_{i_2}), \quad \dots, \quad S_k = (\pi_{i_k+1}, \dots, \pi_{n-1})$$

die Strips von  $\pi$ .

Der Strip  $S_j$  heißt aufsteigend, falls  $\pi_{i_{j-1}} < \dots < \pi_{i_j}$  gilt,

absteigend, falls  $\pi_{i_{j-1}} > \dots > \pi_{i_j}$ .

Ein-elementige Strips sind absteigend, außer  $S_0$  und  $S_k$ .

Lemma 10.2: Sei  $\pi$  eine Permutation der Ordnung  $n$ , sei  $k \in \{0, \dots, n-1\}$

(a) Falls  $k$  in einem absteigenden Strip liegt und  $k-1$  in einem aufsteigenden Strip liegt, dann ex. ein Reversal  $\rho$ , so daß  $\text{bcp}(\pi\rho) < \text{bcp}(\pi)$ .

(b) Falls  $k$  in einem absteigenden Strip liegt und  $k+1$  in einem aufsteigenden Strip, dann ex. ein Reversal  $\sigma$ , so daß  $\text{bcp}(\pi\sigma) < \text{bcp}(\pi)$ .

Lemma 10.3: Sei  $\pi$  eine Permutation mit einem absteigenden Strip.  
Dann ex. ein Reversal  $\rho$ , so daß  $\text{bnp}(\pi\rho) < \text{bnp}(\pi)$ .

Beweis: Wähle  $k$  als das kleinste Element in einem abst. Strip.  
Beh. folgt mit Lemma 10.2 (a).  $\square$

Lemma 10.4: Sei  $\pi$  Perm. ohne absteigenden Strip. Dann ist  $\pi$   
die identische Permutation, oder es ex. Reversal  $\rho$ , so daß  
 $\pi\rho$  einen absteigenden Strip enthält und  $\text{bnp}(\pi\rho) = \text{bnp}(\pi)$ .

Beweis: Sei  $\pi$  nicht die identische Perm. Dann hat  $\pi$  mind. 2 Breakpoints.  
Sowohl  $\pi_0 = 0$  als auch  $\pi_{n+1} = n+1$  liegen in aufst. Strips.  
 $\Rightarrow$  ex. mind. zwei aufst. Strips  $S = (\pi_0, \dots, \pi_i) = (0, \dots, i)$  und  
 $S_1 = (\pi_j, \dots, \pi_{n+1}) = (j, \dots, n+1)$   
mit  $j > i+1$ .

23.7.2003

Lemma 10.5: Sei  $\pi$  eine Permutation mit absteigendem Strip. Sei  $k$   
das kleinste Element in einem absteigenden Strip und  $l$  das  
größte.

$\rho$  das Reversal, daß  $k-1$  neben  $k$  platziert

$\sigma$  das Reversal, daß  $l+1$  neben  $l$  platziert

Wenn sowohl  $\pi\rho$ , als auch  $\pi\sigma$  keine abst. Strips enthält,  
dann gilt  $\rho = \sigma$  und  $\text{bnp}(\pi\rho) = \text{bnp}(\pi) - 2$

Basis: Situation (b) und (c) liegt vor.

$k$  muß in dem von  $\sigma$  umgedrehten Intervall liegen.

$l$  muß in dem von  $\rho$  umgedrehten Intervall liegen.

Annahme:  $\rho \neq \sigma$

$\Rightarrow$  es ex. ein Strip  $S$ , der nur von einem der beiden Perzeals umgedreht wird.

falls  $S$  aufsteigend in  $\pi$ , dann absteigend in  $\pi\rho$  oder  $\pi\sigma$

falls  $S$  absteigend in  $\pi$ , dann aufsteigend in  $\pi\rho$  oder  $\pi\sigma$   $\checkmark$

$\Rightarrow \rho = \sigma$  eliminiert zwei Breakpoints.  $\square$

Satz 10.2: Der Alg. 10.1 ist ein 2-Approximations Algorithmus für das MinSR-Problem.

Beweis: Perzeal am Anfang der Berechnung kann verändert werden, nicht letztem Perzeal, das zwei Breakpoints eliminieren muß, da keine Permutation, mit genau einem Breakpoint ex.

Rest folgt aus Lemma 10.1 bis 10.5.

Laufzeit:  $O(n^3)$

Beste bekannte Approximation: 1,375



bisher: Betrachtung der Moleküle als Strings

→ aber DNA/Proteine sind keine Strings, sondern räumliche Moleküle

- Funktion der Moleküle ergibt sich durch die räumliche Struktur.



## 12. Höherdimensionale Strukturen von Molekülen

- Labormethoden: Röntgen-Kristallographie  
⇒ teuer, mühsam

- (bisherigen) Molekülen gibt es einen Zusammenhang zwischen String-Darstellung und der räumlichen Struktur

⇒ Ziel: Vorhersage der räumlichen Struktur durch die Betrachtung des zugrundeliegenden Strings.

### 12.1. Sekundärstruktur von RNA

Strukturvariante:



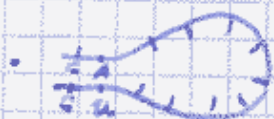
Einzelstrang mit Basen

Adenin, Guanin, Cytosin,

Uracil

→ String über  $\Sigma_{RNA} = \{A, C, G, U\}$

⇒ Primärstruktur der RNA.



Zusammenlagerung des RNA-Strangs durch

Wasserstoffbrücken zwischen den Basen.

⇒ Sekundärstruktur der RNA



Def. 12.1: Sei  $\tau = \tau_1 \dots \tau_n$  die Primärstruktur einer RNA.

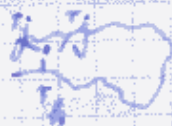
Die Sekundärstruktur von  $\tau$ ,  $\text{SecStruct}_\tau$  ist eine Menge von Index - Paaren

$$\text{SecStruct}_\tau \subseteq \{(i, j) \mid 1 \leq i < j \leq n\}$$

mit der Bedeutung Base  $\tau_i$  und  $\tau_j$  sind miteinander verbunden.

übliche Anforderungen:

(i) jeder Index  $k$  kommt in höchstens einem Paar in  $\text{SecStruct}_\tau$  vor.



nicht zugeordnet

(ii) jedes Paar  $\in \text{SecStruct}_\tau$  ist ein Watson-Crick-Paar  $\{(A, U), (C, G)\}$  oder  $\{(G, C)\} \Rightarrow$  gültige Basenpaare

(iii)



Für jedes Paar in  $\text{SecStruct}_\tau (i, j)$  gilt  $j - i \geq 4$ .

• Tertiärstruktur  $\hat{=}$  Lage der Atome im 3-dim. Raum.

### 12.1.1. Minimierung der freien Energie

- freie Energie  $\hat{=}$  Potential eines Moleküls durch Bindungen

Energie freisetzen

- stabile Moleküle besitzen geringe freie Energie

Ziel: Bestimme Sekundärstruktur mit minimaler freier Energie.

1. Ansatz: Algorithmus von Maxam

Annahme: freie Energie  $\sim$  Anzahl der Basenpaarungen in der Sekundärstruktur.

$\Rightarrow$  desto mehr Basenpaarungen, desto stabiler

Idee: dynamische Programmierung

- Betrachte optimale Strukturen für Teilstrengen
- Setze daraus sukzessive die optimale Struktur für größere Teilstrengen zusammen.

- Notation:
- $T = T_1 \dots T_n$  String (Primärdarstellung der RNA)
  - $S_{i,j} \triangleq$  opt. Sekundärstruktur für den Teilstreng  $T_i \dots T_j$
  - $BP(S_{i,j}) = |S_{i,j}|$  (Anzahl der Basenpaare in  $S_{i,j}$ )
  - $\delta(T_i, T_j) = \begin{cases} 1 & \text{falls } (T_i, T_j) \text{ gültiges Basenpaar} \\ 0 & \text{sonst} \end{cases}$

Beachte: - Der Algorithmus R.1 berechnet lediglich die Anzahl der Basenpaare in der optimalen Sekundärstruktur.

- Ermittlung der Sekundärstruktur ist durch ein Trace-Back-Verfahren möglich.

- Verfeinerung:
- gewichte die Basenpaarungen entsprechend ihrer Bindungsstärke.
  - $f_e(T_i, T_j) \triangleq$  freie Energie der Basenpaarung  $(T_i, T_j)$  (negativ)
  - $E(S_{i,j}) \triangleq$  minimale freie Energie von  $S_{i,j}$
  - $\rightarrow$  Minimierung der freien Energie.

Ersetze in Alg. R.1 die Rekurrenz durch

$$E(S_{i,j}) := \min \begin{cases} E(S_{i+1,j}) & \text{(i)} \\ E(S_{i,j-1}) & \text{(ii)} \\ E(S_{i+1,j-1}) + f_e(T_i, T_j) & \text{(iii)} \\ \min_k \{ E(S_{i,k}) + E(S_{k+1,j}) \} & \text{(iv)} \end{cases}$$

„Beschränkung auf gültige Basenpaare entfällt“

Laufzeit: - Berechnung einer  $(n \times n)$ -Matrix  $\rightarrow O(n^2)$

- jeder Eintrag: (i), (ii), (iii)  $\rightarrow O(1)$

(iv)  $\rightarrow O(n)$

$\Rightarrow O(n^3)$

Erweiterung Ansatz: Algorithmus von Zuker

- Verfeinerung der Berechnung der freien Energie

$\rightarrow$  Betrachtung der freien Energie von Substrukturen

Def. 122:

Sei  $\tau$  Primärstruktur und  $(i, j) \in \text{SecStruct}_\tau$

Wir sagen

• eine Base  $k$  ist erzielbar von  $(i, j)$ , falls  $i < k < j$  und kein Basenpaar  $(i', j') \in \text{SecStruct}_\tau$  existiert mit  $i < i' < k < j' < j$

• ein Basenpaar  $(l, m)$  ist erzielbar von  $(i, j)$ , falls  $i < l < m < j$  und kein Basenpaar  $(i', j') \in \text{SecStruct}_\tau$  existiert mit  $i < i' < l < m < j' < j$

Def. 123: -  $\tau$  Primärstruktur,  $\text{SecStruct}_\tau$  (Sekundärstruktur) 28.7.2003

$(i, j) \in \text{SecStruct}_\tau$

durch  $(i, j)$  induzierte Teilstrukturen

• Stacked Pair: falls  $(i+1, j-1) \in \text{SecStruct}_\tau$

$\Rightarrow$  unmittelbar aufeinanderfolgende Folge von Basenpaaren (Stem)

- Clairin-loop: falls von  $(i, j)$  kein Basenpaar in  $\text{SecStruct}_T$  erreichbar ist.
- Bulge: falls ein von  $(i, j)$  erreichbares Basenpaar  $(i', j') \in \text{SecStruct}_T$  existiert, so daß entweder  $i' - i > 1$  oder  $j - j' > 1$ .
- Interior Loop: falls ein von  $(i, j)$  erreichbares Basenpaar  $(i', j') \in \text{SecStruct}_T$  existiert, so daß sowohl  $i' - i > 1$  und  $j - j' > 1$ .
- Multiple loop: falls von  $(i, j)$  mehr als ein Basenpaar erreichbar ist.

Bemerkung: - Größe einer Substruktur  $\hat{=}$  Anzahl der erreichbaren Basen  
 - Loops  $\hat{=}$  Clairin-loop, Bulge, Interior loops, Multiple loops.

Stem / Stacked Pairs  $\Rightarrow$  stabilisierend (negative freie Energie)  
 Loops  $\Rightarrow$  destabilisierend ("positive" freie Energie)

Annahme: kein Pseudoknoten

Ziele: dynamische Programmierung

2. Teil (iii) des Alg. von Mulsinos

Betrachte die Energie der induzierten Substruktur und nicht nur die der Basenpaare.

Bezeichnung:  $h_{ij}$  := durch das Basenpaar  $(i, j)$  induzierte Substruktur mit minimaler freier Energie

Basiskonvention (Alg. 122)

Bestimmung von  $E(h_{ij})$

Notation:  $f_{\text{Substruktur}} \in \{\text{Stacked Pair, Clairin-loop, bulge, interior}\}$

experimentell  $\hat{=}$  freie Energie der Substruktur

bestimmt

$b(i, j) :=$  freie Energie der Konvention  $(i, j)$

Laufzeit:

- Berechnung $(n \times n)$ -Matrix	$O(n^2)$	} $O(n^2)$
- jeder Eintrag (i), (ii)	$O(1)$	
(iv)	$O(n)$	
(iii) (a), (b)	$O(n)$	
(c), (d)	$O(n)$	
(e)	$O(n^2)$	

⇒ gesamt:  $O(n^4)$   
 → Literatur  $O(n^2)$

### 12.3 Strukturvorhersage für Proteine

- wichtig, besonders auch für med. Forschung
- Strukturebenen der Proteine:

- Primärstruktur: Darstellung als String (Abfolge der Aminosäuren entlang einer Polypeptidkette)

- Sekundärstruktur: Wechselwirkungen zwischen den Atomen des Rückgrats einer Polypeptidkette  
 (groß) 2 Typen:

- Helices
- Faltblätter
- (→ Schleifen  $\hat{=}$  alle Abschnitte, in der Primärstruktur die keine Helices oder Faltblätter sind)

- Tertiärstruktur: (3-dimensionale Struktur)

→ Zusammenfassung vom Sekundärstruktur zu

- Motive
- Domänen (bitragen spezifische Funktion, aktives Zentrum)

- Quartärstruktur: Aufbau eines Proteins aus Untereinheiten und evtl. aus zusätzlichen, molekularen Bausteinen (z.B. Hämoglobin)

### B.3.1 Gitter-Modell (Minimierung der freien Energie)

- 3dim. Modell:  $\sim$  Festlegung der Lage aller Atome im Raum, Spezifikation des Winkel, Bindungslängen

- Abstraktion der Strukturen hier:

- Betrachte Aminosäuren (punktförmig) statt Atome
- Winkel nur  $90^\circ, 180^\circ, (270^\circ)$
- Bindungslängen einheitlich

$\Rightarrow$  Einbettung der Aminosäuresequenz (Strings) in ein Gitter.

- Abstraktion der freien Energie:

- grobe Einteilung der Aminosäuren in

- hydrophob (wasserabweisend, H, 1)
- hydrophil (wasserbindend, P, 0)

- häufig kommt es zur Ausbildung eines hydrophoben Molekülkerns.

$\Rightarrow$  maximale Anzahl der Wechselwirkungen zwischen hydrophoben Aminosäuren.

Def. 12.20: Sei  $S = s_1 \dots s_n$  über  $\{0, 1\}$

Eine Abbildung  $\gamma: \{1, \dots, n\} \rightarrow \mathbb{Z}^d$  heißt Einbettung von  $S$  in ein Gitter  $\mathbb{Z}^d$  ( $d \in \mathbb{R}, \mathbb{B}$ ), wenn

- alle in  $S$  benachbarten Symbole sind auch in  $\mathbb{Z}^d$  benachbart.
- keine zwei Positionen in  $S$  werden dem gleichen Gitterpunkt zugeordnet.

$\leadsto$  self-avoiding walk

- Falls  $\gamma$  einer Position im Gitter den Buchstaben 1 zuordnet, dann heißt diese Position mit 1 belegt (Analog für 0)

Def. 12.21: Sei  $S = s_1 \dots s_n$  über  $\{0,1\}$  und  $\gamma$  Einbettung von  $S$  in  $\mathbb{Z}^d$ .

Zwei Positionen  $i$  und  $j$  sind **verbundene Nachbarn**, wenn  $|i-j|=1$ .

Zwei Positionen  $i$  und  $j$  sind **topologische Nachbarn**, wenn

$$\|\gamma(i) - \gamma(j)\| = 1$$

- Wir nennen eine Position  $i$  mit  $s_i = 1$  eine **Ein-Position**

- Wir nennen ein Positionspaar  $(i,j)$  von zwei Ein-Positionen, wobei  $i,j$  topologische Nachbarn sind **1-1-Paar, Kontakt**.

Def. 12.23: Gegeben sei ein Gitter  $\mathbb{Z}^d$ .

Das  $d$ -Max-1-1-Problem ist definiert durch

Eingabe:  $S = s_1 \dots s_n \in \{0,1\}^*$

zulässige Lösungen: jede Einbettung  $\gamma$  von  $S$  in  $\mathbb{Z}^d$

Kosten: Anzahl der 1-1-Paare

Ziel: Maximierung

Satz 12.6: Das 2-Max-1-1-Problem ist NP-schwer (Analog 3-Max-1-1).  $\square$

$\Rightarrow$  Einschränkung auf  $\mathbb{Z}^2$

$\rightarrow$  Approximationsalgorithmus.

Lemma 12.1: Sei  $S = s_1 \dots s_n \in \{0,1\}^*$  und  $\gamma$  eine Einbettung von  $S$  in das Gitter  $\mathbb{Z}^2$ .

(i) Jede Position (außer den Endpunkten) in  $S$  hat max. 2 topologische Nachbarn.

(ii) Für jedes 1-1-Paar  $(i,j)$  gilt, daß  $i$  gerade und  $j$  ungerade oder  $i$  ungerade und  $j$  gerade ist.

Beweis:

(i) ✓

(ii) Färbungsargument:

- Färbe das Gitter mit zwei Farben (benachbarte Knoten besitzen verschiedene Farben).

→ "Schachbrettmuster"

- Betrachte 1-1-Paar  $(i, j)$  und zugehörige Knoten  $\gamma(i) = (x_1, y_1)$ ,  $\gamma(j) = (x_2, y_2)$

- Da  $(x_1, y_1)$  und  $(x_2, y_2)$  topologische Nachbarn sind, besitzen sie unterschiedliche Farben.

- Auf dem Pfad von  $(x_1, y_1)$  nach  $(x_2, y_2)$  wechseln die Farben der Knoten ab.

⇒ Pfad beginnt mit der einen Farbe und endet mit der anderen (insbesondere gilt das auch für den durch  $\gamma$  induzierten Pfad).

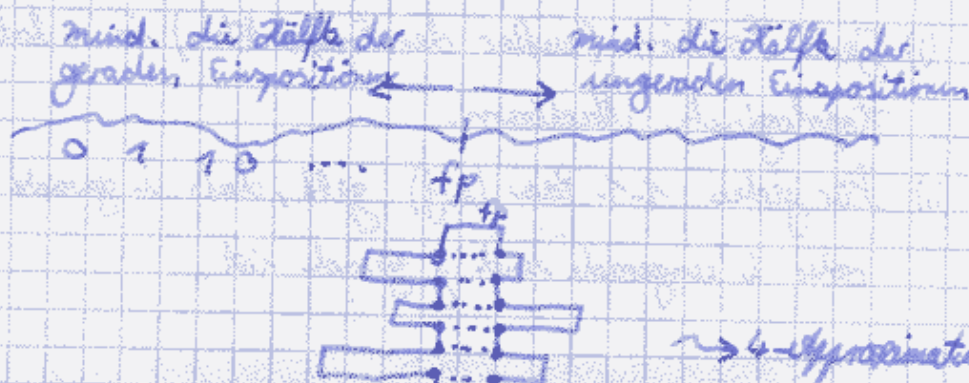
⇒ ungerade Anzahl von Kanten auf einem solchen Pfad.

⇒  $i$  genau dann gerade, wenn  $j$  ungerade ist. □

30.7.2003

Approximationsalgorithmus für 2-1aa-1-1-Problem

- Annahme: die Endpositionen des einbettenden Strings sind  $\neq 1$ .





Verteilung der Einspositionen entsprechend ihrer Parität

$X :=$  die Menge der Einspositionen mit gerader/ungerader Parität.

$Y :=$  die Menge der Einspositionen mit ungerader/gerader Parität.

wobei:  $|X| \leq |Y|$

$\Rightarrow$   $XY$ -Partitionierung

Satz 12.7: Sei  $s = s_1 \dots s_n \in \{0, 1\}^*$ .

Sei  $\text{Opt}(s)$  eine optimale Lösung des 2-Max-1-1-Problems.

$X, Y$  eine Partitionierung wie oben.

Dann gilt:  $\text{cost}(\text{Opt}(s)) \leq 2 \cdot |X|$

Beweis: folgt direkt aus Lemma 12.1.

Def. 12.23: Sei  $s = s_1 \dots s_n \in \{0, 1\}^*$ ,  $XY$ -Partitionierung.

Dann heißt eine Position  $fp \in \{1, \dots, n\}$  **Faltungspunkt** wenn gilt:

$$\bullet \underbrace{|\{i \in X \mid i \leq fp\}|}_{X'} \geq \frac{|X|}{2} \quad \text{und} \quad \underbrace{|\{j \in Y \mid fp < j\}|}_{Y'} \geq \frac{|X|}{2}$$

oder

$$\bullet \underbrace{|\{j \in Y \mid j \leq fp\}|}_{Y'} \geq \frac{|X|}{2} \quad \text{und} \quad \underbrace{|\{i \in X \mid fp < i\}|}_{X'} \geq \frac{|X|}{2}$$

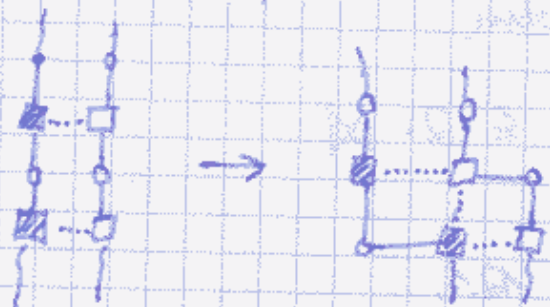
Satz 12.8: der Algorithmus 12.7 ist ein 4-Approximationsalgorithmus für das 2-Max-1-1-Problem für eine Eingabe  $\mathcal{I}$  mit "Nicht-Endpositionen"

Beweis:

- Algorithmus 12.7 bildet  $|X|/2$  1-1 Paare
- Opt hat  $2|X|$  1-1 Paare

$$\Rightarrow \frac{2|X|}{|X|/2} = 4 \text{ Approximationsgüte von } 4. \quad \square$$

Idee für bessere Approximation:



$\Rightarrow$  3-Approximation

(SODA 2002, A. Newman)

### 12.3.2 Protein - Threading

Idee: Tertiärstruktur durch Alignment der Primärstruktur mit bereits bekannten Tertiärstrukturen zu bestimmen

Annahme: Viele verschiedene Primärstrukturen, aber nur "wenige" verschiedene Tertiärstrukturen.

$\rightarrow$  Modell für Tertiärstruktur

Def. 12.24: Strukturmodell / Kernmodell

Sei  $s$  ein String. Ein Strukturmodell für  $s$  ist ein 5-Tupel

$M = (m, c, \beta, l_{\min}, l_{\max})$ , mit den folgenden Eigenschaften:

- Die Struktur von  $s$  beinhaltet  $m$  Kernregionen ( $\text{Kern} \hat{=} \text{Helix}, \text{Fallblatt}, \text{Motif}, \text{Domäne}$ )  $C_1, \dots, C_m$

- Länge der Kernregion  $C_i$  ist durch  $c_i$  gegeben

$$c = (c_1, \dots, c_m)$$

- Die Kernregionen  $C_i$  und  $C_{i+1}$  sind durch eine Schleife  $\beta_i$  miteinander verbunden.

$$\beta = (\beta_1, \dots, \beta_m)$$

- min. Länge der Schleife  $\beta_i$ :  $l_i^{\min}$ :  $l_{\min} = (l_1^{\min}, \dots, l_m^{\min})$

- max. Länge der Schleife  $\beta_i$ :  $l_i^{\max}$ :  $l_{\max} = (l_1^{\max}, \dots, l_m^{\max})$

Darüberdies gilt:

- $|S| = \beta_0 + \sum_{i=1}^m (c_i + \beta_i)$

- $l_i^{\min} \leq \beta_i \leq l_i^{\max}$

↑ Einpassen eines Strings (Primärstruktur) in ein Kernmodell  
 $\hat{=} \text{Threading}$

Def. BAS:  $S$  String,  $M = (m, c, l, l_{\min}, l_{\max})$  Kernmodell von  $S$ .

$S'$  String (mit unbekannter Struktur)

Ein Threading  $T$  von  $S'$  in  $M$  ein  $m$ -Tupel

$$T = (t_1, \dots, t_m)$$

mit

Anordnungs-  
anforderungen

$$\begin{cases} 1 + l_0^{\min} \leq t_1 \leq 1 + l_0^{\max} \\ t_i + c_i + l_i^{\min} \leq t_{i+1} \leq t_i + c_i + l_i^{\max} & \forall i, 1 \leq i < m \\ t_m + c_m + l_m^{\min} \leq |S'| + 1 \leq t_m + c_m + l_m^{\max} \end{cases}$$

$$1 + \sum_{j < i} (c_j + l_j^{\min}) \leq t_i \leq |S'| + 1 - \sum_{j > i} (c_j + l_j^{\min})$$

(Abstandsanforderung)

Bewertung eines Threadings:

- Bewertung nur der Zuordnung eines bestimmten Kerns  
 $\rightarrow$  polynomiell
- Wechselwirkungen zwischen den Kernregionen  
 $\rightarrow$  "schwieriger"
- Zuordnung zum Kern:  $g_r(i, t_i)$
- Wechselwirkungen zwischen  $r$  Kernn:  $g_r(i_1, \dots, i_r, t_{i_1}, \dots, t_{i_r})$

hier:  $r = 2$

Def. 12.26: Protein-Threading - Problem

Eingabe: Strukturmodell  $H = (m, c, \mathcal{E}, l_{min}, l_{max})$  eines Strings  $s$   
 ein String  $s'$   
 zwei Funktionen  $g_1, g_2$

zulässige Lösungen: Alle möglichen Threadings  $T = (t_1, \dots, t_m)$

Kosten:  $cost(T) = \sum_{i=1}^m g_1(i, t_i) + \sum_{i=1}^m \sum_{j=i+1}^m g_2(i, j, t_i, t_j)$

Optimierungsziel: Minimierung

→ Protein-Threading-Problem ist NP-schwer  
 (Reduktion von MaxCut)

→ Branch & Bound - Ansatz

