

# Sequenzvergleiche ohne Alignments durch lokales Dekodieren von Sequenzen

Christian Otto

(13. Juli 2007)

## Inhaltsverzeichnis

<b>1 Grundidee</b>	<b>1</b>
<b>2 Naiver Algorithmus</b>	<b>2</b>
2.1 Definitionen . . . . .	2
2.2 Lokales N-Dekodieren . . . . .	2
2.3 Ermittlung der Distanzmatrix . . . . .	4
2.4 Aussagefähigkeit der Distanzen . . . . .	5
<b>3 Verbesserungen</b>	<b>6</b>
3.1 Verbesserung der Komplexität . . . . .	6
3.2 Erweiterung des evolutionären Modells . . . . .	8
<b>4 Leistung bei HIV/SIV Subtyping</b>	<b>9</b>
4.1 Kurze Sequenzen: HIV/SIV <i>gag</i> , <i>pol</i> , <i>env</i> und <i>nef</i> . . . . .	9
4.2 Non-coding Long Terminal Repeats (LTR) . . . . .	9
<b>5 Abschließende Bemerkungen</b>	<b>9</b>

## 1 Grundidee

Normalerweise werden bei Sequenzvergleichen multiple Alignments eingesetzt. Diese weisen jedoch neben der schlechten Laufzeit einige Nachteile auf: Sie beruhen auf Editieroperation (Deletion, Substitution, Insertion), wobei z.B. Deletionen/Insertionen nicht den regulären evolutionären Modellen folgen und somit evolutionär schlecht bewertet sind.

Beim multiplen Alignment werden mehrdeutige Regionen (in Bezug auf die Editieroperationen) gelöscht. → Dies führt zu Informationsverlust, da z.B. Permutationen und Inversionen dadurch vernachlässigt werden, obwohl sie evolutionär bedeutsam sind.

Um solche Regionen zu berücksichtigen, wird manuelles Alignment (Expert Editing) durchgeführt, welches jedoch zu einem großen Zeitaufwand vor allem bei langen Sequenzen und zu Problemen bei der Reproduzierbarkeit führt. Dabei wird zusätzliche Information verwendet, um auch Regionen zu berücksichtigen, die mehrfach aufeinanderfolgende Deletionen/Insertionen beinhalten, aber trotzdem evolutionär bedeutsam sind und im Alignment auftauchen sollten. Somit kann ein multiples Alignment verbessert werden (um z.B. ein lokales Minimum zu verlassen).

### Idee der Sequenzvergleiche ohne Alignments:

**Intuitiver Ansatz:** Bei den Sequenzvergleichen ohne Alignments ist der intuitive Ansatz, dass man die Nukleotidhäufigkeiten der Sequenzen vergleicht, da Sequenzen mit signifikant unterschiedlichen Nukleotidzusammensetzungen nicht nah verwandt sein können.

Problematisch ist, dass es auch nicht nah verwandte Sequenzen mit identischen Nukleotidhäufigkeiten existieren.

**Nächster Schritt:** Man vergleicht N-Wörter (Teilsequenzen der Länge N) zwischen Sequenzen, um daraus Ähnlichkeiten zu berechnen. Die evolutionäre Idee ist, dass man somit homologe Bereiche identifiziert und dass Sequenzen mit vielen gemeinsamen homologen Blöcken nah verwandt sind.

Problematisch hierbei ist, dass die N-Wörter entweder identisch oder unterschiedlich sind. Dies kann kompensiert werden, indem man auch nicht perfekt identische N-Wörter betrachtet, wobei die im Vortrag vorgestellten Modelle und Algorithmen dies noch nicht berücksichtigen. Obwohl dann zusätzlich die Frage entsteht, wie man nicht perfekt identische N-Wörter mit einem Score bewertet und diesen einbezieht bzw. ab welchem Ähnlichkeitsgrad es überhaupt sinnvoll ist von homologen Blöcken zu sprechen.

### Lokales N-Dekodieren:

Um evolutionäre Informationen aus mehrdeutigen Regionen (laut Alignment) zu gewinnen und nicht auf manuelles Alignment angewiesen zu sein, wurde die Idee vom lokalen Dekodieren vom Grad N von Sequenzen entwickelt:

Diese beruht auf Vergleichen zwischen überlappenden Wörtern der Länge N. Für jede Sequenzposition i gibt es somit N überlappende Wörter, wobei i dann jeweils an Position 1, ..., N dieser Wörter stehen kann.

Übereinstimmungen von solchen Teilwörtern zwischen verschiedenen Sequenzen entsprechen konzeptionell homologen Blöcken zwischen diesen Sequenzen. Man kann nun auf Basis dieser Idee Sequenzvergleiche durchführen und Distanzen zwischen Sequenzen berechnen.

## 2 Naiver Algorithmus

### 2.1 Definitionen

**Definition der Relation  $\stackrel{k}{\sim}_N$ :**

Sei s eine Sequenz und  $N \in \mathbb{N}, N > 0$ .

Dann ist die Relation  $(\stackrel{k}{\sim}_N)_{-1 \leq k < N}$  für alle Paare von Positionen  $(i, j)$  von s definiert durch:

- $i \stackrel{-1}{\sim}_N j \Leftrightarrow i = j$
- Für  $0 \leq k < N$ :  
 $i \stackrel{k}{\sim}_N j \Leftrightarrow \exists l, 0 \leq l \leq k: s_{[i-l, i-l+N-1]} = s_{[j-l, j-l+N-1]}$ .

Für alle Zahlen  $-1 \leq k < N$  sei dann  $\stackrel{k}{\sim}_N$  die transitive Hülle der Relation  $\stackrel{k}{\sim}_N$ .

**Aus der Relation  $\stackrel{k}{\sim}_N$  folgende Aussagen:**

1.  $i \stackrel{k-1}{\sim}_N j \Rightarrow i \stackrel{k}{\sim}_N j$  (wichtig für Iteration im Algorithmus)
2.  $i \stackrel{k-1}{\sim}_N j \Rightarrow i + 1 \stackrel{k}{\sim}_N j + 1$  (wichtig für Nachfolgerbeziehung im Algorithmus)
3.  $i \stackrel{k}{\sim}_N j \Rightarrow s_{[i, i-k+N-1]} = s_{[j, j-k+N-1]}$

Für alle  $-1 \leq k < N$  ist die Relation  $\stackrel{k}{\sim}_N$  eine Äquivalenzrelation (unter der Annahme, dass  $|s| \geq N$ ).

$\Delta_k$  sei die Menge der zur Relation  $\stackrel{k}{\sim}_N$  gehörigen Äquivalenzklassen.

### 2.2 Lokales N-Dekodieren

**Verwandtschaft zwischen Sequenzpositionen:**

Zwei Sequenzpositionen a, b sind genau dann verwandt, wenn  $a \stackrel{N-1}{\sim}_N b$  gilt

d.h. es gibt Positionen  $k_0, \dots, k_n$  mit  $k_0 = a, k_n = b$  und  $\forall i \in \{0, \dots, n-1\} : k_i \stackrel{N-1}{\sim}_N k_{i+1}$ . Dadurch wird sozusagen der transitive Abschluss dargestellt. Es gilt nun  $i \stackrel{N-1}{\sim}_N j$ , wenn ein überlappendes N-Wort an Position i und j übereinstimmt, wobei i und j in diesem Wort an der gleichen Position stehen.

Beispiel:

```

Sequenz      ...tagacacta...tccacactg...
                i           j
Menge von    tagac          tccac
überlappenden agaca          ccaca
N-Wörtern    gacac          cacac
              acact          acact
              cacta          cactg
    
```

$$\Rightarrow i \stackrel{-1}{\sim}_5 j, i \stackrel{0}{\sim}_5 j, i \stackrel{1}{\sim}_5 j, i \stackrel{4}{\sim}_5 j$$

**Ermittlung der Äquivalenzklassen:**

- a) Menge der N-Wörter jeder Sequenzposition ermitteln (siehe Beispiel oben) und miteinander zeilenweise vergleichen.  
 → Man findet somit Relationen  $i \stackrel{N-1}{\sim}_N j$  zwischen Sequenzpositionen  $i$  und  $j$ .
- b) Transitiven Abschluss der Relation  $\stackrel{N-1}{\sim}_N$  und Äquivalenzklassen bilden. Dazu fügt man nacheinander miteinander in Relation stehende Positionen in eine Menge und vereinigt Mengen, bei denen mindestens eine Sequenzposition gleich ist. Somit erhält man die Äquivalenzklassen  $\Delta_{N-1}$ .

Beispiel:

```

a) >seq1  CATTG TCGC TGGAC CACAC CTTGT CCCTA
    >seq2  CACTT GGCAC CATAc CATGC
    >seq3  CACTT CTTC CTGGA CcTCC
    
```

```

seq1, 11 CCGCTGGAC  seq2, 5 CACTTGGAC  seq3, 5 CACTTCTTT  seq3, 12 TTCCTGGAC
  CCGCT  CACTT  CACTT  TTCCT
  CGCTG  ACTTG  ACTTC  TCCTG
  GCTGG  CTGG  CTTC  CCTGG
  CTGGA  TTGGA  TTCTT  CTGGA
  TGGAC  TGGAC  TCTTT  TGGAC
    
```

```

b) class_1 : T0
    seq1  3 CATTGTC
    seq1  22 CACTTGTGTC
class_2 : T1
    seq1  4 CATTGTCC
    seq1  23 ACCTTGTCC
class_3 : G0
    seq1  5 CATTGTCCG
    seq1  24 CCTTGTCCC
class_4 : T2
    seq1  6 ATTGTCCGC
    seq1  25 CTTGTCCCT
class_5 : C0
    seq1  7 TTGTCCGCT
    seq1  26 TTGTCCCTA
class_6 : C1
    seq1  8 TGTCGCTG
    seq1  27 TGTCCTA
class_7 : C2
    seq1  10 TCCGCTGGA
    seq3  11 TTTCTGGA
class_8 : T3
    seq1  11 CCGCTGGAC
    seq2  5 CACTTGGAC
    seq3  5 CACTTCTTT
    seq3  12 TTCCTGGAC
class_9 : G1
    seq1  12 CGCTGGACC
    seq2  6 ACTTGGACA
    seq3  13 TCCTGGACC
class_10 : G2
    seq1  13 GCTGGACCA
    seq2  7 CTTGGACAC
    seq3  14 CCTGGACCT
class_11 : A0
    seq1  14 CTGGACCAC
    seq2  8 TTGGACACA
    seq3  15 CTGGACCTC
class_12 : C3
    seq1  15 TGGACCACA
    seq2  9 TGGACACAT
    seq3  16 TGGACCTCC
class_13 : C4
    seq1  16 GGACCACAC
    seq3  17 GGACCTCC
class_14 : C5
    seq2  1  CACTT
    seq3  1  CACTT
class_15 : A1
    seq2  2  CACTTG
    seq3  2  CACTTC
class_16 : C6
    seq2  3  CACTTGG
    seq3  3  CACTTCT
class_17 : T4
    seq2  4  CACTTGGa
    seq3  4  CACTTCTT
    
```

**Lokales N-Dekodieren:**

- c) Nun Äquivalenzklassen auf Zustände abbilden:  
 Dabei ist es in dem Paper üblich gewesen, dass einelementige Äquivalenzklassen auf das ursprüngliche Symbol (A, T, C, G) abgebildet werden. Diese Symbole werden im Folgenden nicht weiter betrachtet. Andere Äquivalenzklassen werden auch auf das ursprüngliche Symbol und einem Index abgebildet. Der Index wird pro Symbol einzeln gezählt und wird nach Vorkommen der Zeichen in einer Referenzsequenzfolge verwendet.  
 ⇒ Dies wird als lokales N-Dekodieren der Sequenz bezeichnet (analog zu Hidden Markov Modell), da man Positionen einen Zustand zuordnet (nämlich die verwandtschaftliche Verbindung zu anderen Sequenzen), der normalerweise versteckt ist.

Beispiel:

```
c) >seq1   C A T0T1G0T2C0C1G C2T3G1G2A0C3C4A C A C C T0T1G0T2C0C1C T A
>seq2   C5A1C6T4T3G1G2A0C3A C A T A C C A T G C
>seq3   C5A1C6T4T3C T T T C C2T3G1G2A0C3C4T C C
```

**2.3 Ermittlung der Distanzmatrix**

- d) Erstellen einer Liste mit Anzahl des Vorkommens jedes Zustandes (der einer Äquivalenzklasse entspricht) je Sequenz

Beispiel:

```
d)   seq1   seq2   seq3
T3    1      1      2
T4    0      1      1
A0    1      1      1
A1    0      1      1
G1    1      1      1
G2    1      1      1
C2    1      0      1
C3    1      1      1
C4    1      0      1
C5    0      1      1
C6    0      1      1
```

(T0, T1, T2, G0, C0 and C1 are  
 only repeated in seq1, two times  
 each)

**Berechnung der Distanzmatrix:**

Seien Zustände  $x_1, \dots, x_n$  mit  $n \in \mathbb{N}$  gegeben.

Nun sei  $|u|_{x_i}$  die Anzahl von Vorkommen des Zustandes  $x_i$  in der Sequenz  $u$  und  $|u|$  die Länge der Sequenz  $u$ .

- e) Berechnen der Ähnlichkeiten  $\text{sim}(u, v)$  zwischen Sequenzen  $u$  und  $v$  durch:

$$\text{sim}(u, v) = \sum_{i=1}^n \min(|u|_{x_i}, |v|_{x_i})$$

Dabei summiert man die minimale Anzahl der einzelnen Zustände in  $u$  und  $v$  auf.  
 Normalisierung und Berechnung der Distanzen durch:

$$\text{dist}(u, v) = 1 - \frac{\text{sim}(u, v)}{\min(|u|, |v|)}$$

Man normiert die Ähnlichkeiten zwischen den Sequenzen  $u$  und  $v$ , indem man durch die minimale Sequenzlänge von  $u$  und  $v$  teilt. Da die maximale Anzahl von Ähnlichkeiten der Sequenzlänge der kürzeren Sequenz entspricht, normiert man die Ähnlichkeit somit auf das Intervall  $[0, 1]$ . Zum Berechnen der Distanz muss lediglich noch 1-Ähnlichkeit gerechnet werden.

Beispiel:

d)	seq1	seq2	seq3	
T3	1	1	2	
T4	0	1	1	
A0	1	1	1	
A1	0	1	1	e) Similarities
G1	1	1	1	seq2 seq3
G2	1	1	1	seq1 5 7
C2	1	0	1	seq2 9
C3	1	1	1	Distances
C4	1	0	1	seq2 seq3
C5	0	1	1	seq1 0.75 0.65
C6	0	1	1	seq2 0.55

⇒ Erstellung eines Verwandtschaftsbaumes möglich (z.B. Neighbor-Joining-Tree)

## 2.4 Aussagefähigkeit der Distanzen

Das Distanzmaß  $\text{dist}(u, v)$  ist zunächst nicht aussagekräftig bei einem neuen Verfahren, da man nicht weiß, bei welchem Wert von einem Zufallsbefund zu sprechen ist und ab wann eine konservierte Verwandtschaft vorliegen könnte. Dies gibt dann auch Auskunft darüber, wie man den Baum erstellt, da man Sequenzen mit niedriger Distanz erst weit unten aufsplittet, weil diese nah verwandt sind.

Idee: Vergleich mit Distanzwerten von randomisierten Sequenzen ⇒ Bootstrap (statistisches Verfahren)

### Allgemeine Vorgehensweise bei Bootstrap:

1. Erstellen von Bootstrapsamples aus einer Stichprobe (Resampling)
2. Berechnung einer bestimmten Größe aus dem Sample (nicht unbedingt statistisch)
3. Iteration dieser Schritte (1000-10000 Iterationen sind sinnvoll)  
⇒ Ergebnis: Verteilung der Größe aus den Samples

⇒ p-Wert Schätzung möglich

### Anwendung des Bootstraps auf diesen Sachverhalt:

Sei Distanz  $\text{dist}(u, v)$  der Sequenzen  $u$  und  $v$  gegeben.

1. Zuerst  $u$  und  $v$  resampeln:  
Dazu für jede neue Sequenz aus der alten Sequenz mit Zurücklegen ziehen bis man eine neue Sequenz gleicher Länge erhält. Dabei kann es sein, dass manche Positionen aus der alten Sequenz gar nicht oder mehrfach gewählt werden.  
⇒ neue Sequenzen  $u^*$ ,  $v^*$
2.  $\text{dist}(u^*, v^*)$  über N-lokales Dekodieren berechnen
3. mehrfach iterieren ⇒ Verteilung von  $\text{dist}$  erstellen

Wenn  $\text{dist}(u, v)$  nun extrem in der Verteilung liegt (z.B. in den 2.5% niedrigsten oder höchsten Werten), dann deutet dieser Distanzwert auf keinen Zufallsbefund hin sondern besitzt Aussagekraft (z.B. starke Verwandtschaft). Man schätzt somit quasi die Wahrscheinlichkeit eines bestimmten Distanzwertes des neuen Verfahrens ein. Man erhält eine Verteilung der Distanzwerte, wodurch man die Möglichkeit hat, einen p-Wert (wie oben bereits beschrieben) zu berechnen. Der p-Wert gibt dabei die Wahrscheinlichkeit an, einen solchen oder noch extremeren Wert unter der gegebenen Verteilung zu erhalten.

**Abschätzen der Stabilität des Baumes:**

Wenn man die Distanzmatrix zwischen den ursprünglichen Sequenzen erstellt hat, kann man neben der Aussagekraft der einzelnen Distanzen zusätzlich die Aussagekraft des aus der Distanzmatrix erstellten Baumes abschätzen. Dies funktioniert, indem man den Bootstrap-Algorithmus etwas abändert:

1. Resampling auf alle Sequenzen anwenden
2. Distanzmatrix berechnen und Baum erstellen
3. Iteration
  - ⇒ Aus den Bäumen kann man mit dem Programm Phylib einen Konsensusbaum erstellen, der Bootstrap-Werte bei Aufspaltungen enthält [3]. Diese sagen aus, in wieviel Prozent aller Bootstrap-Bäumen, diese Teiltopologie vorhanden war.
  - ⇒ Man erhält dadurch Auskunft über die Stabilität und Aussagekraft des Baumes.

**3 Verbesserungen****3.1 Verbesserung der Komplexität**

Der naive Algorithmus hat eine hohe Komplexität in CPU-Zeit.

Das Problem dabei ist die Berechnung der Äquivalenzklassen  $\Delta_{N-1}$ .

⇒ Verbesserung auf lineare Komplexität in Zeit und Speicher möglich.

**Dazu weitere Definitionen:**

- $\nabla_k \subset \Delta_k$  ist definiert durch:
 
$$\delta \in \nabla_k \Leftrightarrow \delta \in \Delta_k, \exists i, j \in \delta \text{ mit } (i+1) \not\sim_N^k (j+1)$$
- Die Abbildung  $F_k: \Delta_k \rightarrow \Delta_{k+1}$  ist definiert durch:
 
$$F_k(\gamma) = \delta \Leftrightarrow \gamma \in \Delta_k, \delta \in \Delta_{k+1}, \forall i \in \gamma : i+1 \in \delta$$
 Man nennt  $\delta$  dann k-Nachfolger von  $\gamma$ .

**Funktionsweise des verbesserten Algorithmus:**

1. Ermittlung von  $\Delta_{-1}$ ,  $\Delta_0$ ,  $\nabla_0$  und  $F_{-1}$
2. Iterative Berechnung von  $\Delta_k$ ,  $\nabla_k$  und  $F_{k-1}$  aus  $\Delta_{k-1}$ ,  $\nabla_{k-1}$  und  $F_{k-2}$  für  $0 < k < N-1$ 
  - ⇒  $\Delta_{N-1}$

**Ermittlung von  $\Delta_0$ ,  $\nabla_0$  und  $F_{-1}$ :**

Gegeben sei die Sequenz s.

Beispiel: (N=2)

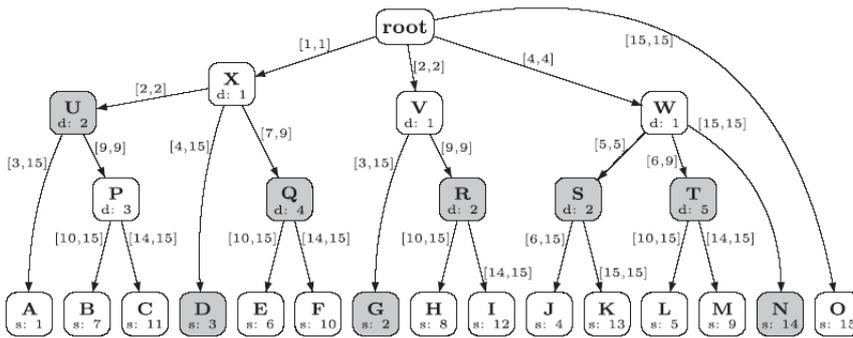
```
positions : 1 2 3 4 5 6 7 8 9 10 11 12 13 14
s         = a g a c c a a g c a a g c c
```

1. Erstellung des Suffixbaumes aus der Sequenz s:
 

Ein Suffixbaum T für einen String s mit m Symbolen ist ein gerichteter Baum mit m Blättern. Dabei sind die Blätter jeweils die Sequenzpositionen von s. Der Baum T enthält alle Suffixe von S, wobei mehrfach auftretende Teilstrings nur einmal in T enthalten sind. Somit werden gleiche Präfixe zusammengefasst und innere Knoten sind gemeinsame Substrings der darunterliegenden Suffixe.

Beschriftung: alle Knoten sind mit einem Label versehen; d bezeichnet die gemeinsame Präfixlänge bei einem inneren Knoten; die gerichteten Pfade sind mit dem Substring von s markiert, welcher identisch ist (meist erstes Auftreten in Sequenz ist Referenz); die Blattknoten sind mit der Sequenzposition bezeichnet, für welche der Pfad von der Wurzel zum Blatt den Suffix von s ab dieser Position darstellt.

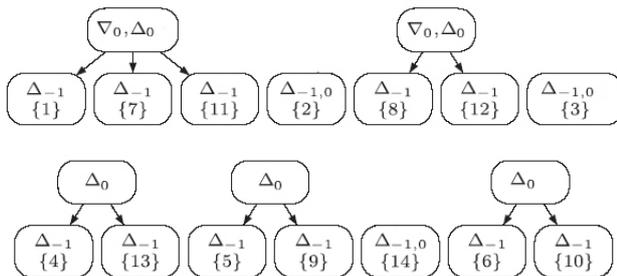
Beispiel:



2. Ablesen von  $\Delta_0$  und  $\nabla_0$ :

Dazu werden die Kenntnisse aus dem Suffixbaum (siehe 1.) genutzt. Man sucht alle Knoten, die von der Wurzel über einen Pfad der Länge größer gleich N erreicht werden (dabei die jeweils die ersten). Als Länge des Pfades wird hier die Summe der Längen der Substring auf den Kanten genommen. Alle Blätter, die unterhalb dieser Knoten liegen (evtl. nur der Knoten selbst, wenn dieser ein Blatt ist), sind dann in einer Äquivalenzklasse von  $\overset{0}{\cup}_N$ . Die Menge aus diesen Klassen ist dann  $\Delta_0$ . Alle Blätter unterhalb von Knoten, die von der Wurzel über einen Pfad der Länge genau N erreicht werden können, bilden jeweils Mengen, deren Vereinigung die Menge  $\nabla_0$  ergibt.

Beispiel:



3. Erstellen der Funktion  $F_{-1}$ :

Aus den Mengen  $\Delta_0$  und  $\nabla_0$  kann die Funktion  $F_{-1}$  erstellt werden. Dabei ist für ein Element  $\delta \in \Delta_{-1}$  ( $\delta = \{i\}$  laut Definition)  $F(\delta) = \gamma \in \Delta_0 \Leftrightarrow i + 1 \in \gamma$ .

4. Iteration (grobe Vorgehensweise, einige Details wurden vernachlässigt):

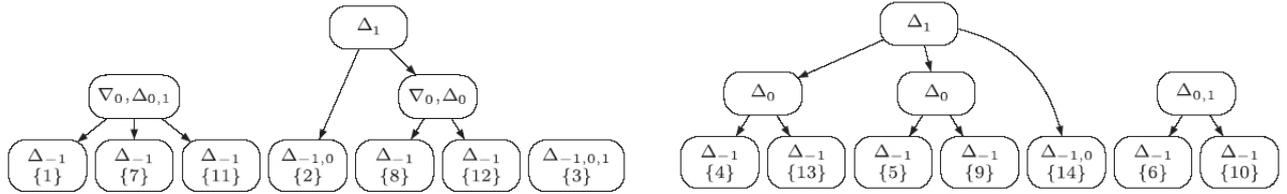
- (a) Für jedes Element  $\delta$  aus  $\nabla_{k-1}$  ein neues Element  $\omega$  in  $\nabla_k$  erstellen,  $F(\delta) = \omega$ .
- (b) Solange  $F^{-1}(\omega) \cap \nabla_{k-1} \neq \emptyset, \gamma \in F^{-1}(\omega) \cup \nabla_{k-1}$ :  
 Kinder Elemente von  $\omega$  sind Nachfolger der Kinder von  $\gamma$ .  
 Nachfolger der Eltern der Vorgänger der Kinder von  $\omega$  ist  $\omega$ .  $\Rightarrow$  ggf. Iteration
- (c) Falls bei  $k = 1$  kein Nachfolger von einem Element aus  $\Delta_0$  gesetzt ist, wird der Nachfolger auf den Nachfolger von einem Kind des Elementes gesetzt.
- (d) Iteration mit  $k = k + 1$  solange  $k < N - 1$  gilt.

**Bemerkungen zum Algorithmus:**

Der Algorithmus beruft sich auf einige theoretischen Aussagen zu  $\Delta, \nabla$  und  $F$ . Genaueres dazu (inklusive der Beweise) findet sich in dem Paper [2].

Im verbesserten Algorithmus wird immer nur eine Sequenz s betrachtet. Um mehrere Sequenzen zu vergleichen, kann dies mit einer einfachen Veränderung geschehen, indem man vorher die Sequenzen zu

Beispiel:



einer Sequenz zusammenfasst, wobei zwischen den Sequenzen jeweils ein einzigartig paarweise verschiedenes Trennsymbol eingefügt wird. Somit wird gewährleistet, dass nur die einzelnen Sequenzen verglichen werden und es zu keinen Übereinstimmungen von N-Wörtern kommt, die über diese Grenzen hinweg gehen.

**Algorithmus:**

---

```

1: init  $\Delta_{-1}$  ; {create a node for each position of  $s$ }
2:  $\Delta_0$  and  $\nabla_0 \leftarrow$  suffix tree of  $s\$$  ;
3:  $F_* \leftarrow F_{-1}$  ;
4: for  $k = 1$  to  $N - 1$  do
5:    $\nabla_k \leftarrow \emptyset$  ;  $\mathcal{S}_{done} \leftarrow \emptyset$  ;
6:   while  $\nabla_{k-1} \not\subseteq \mathcal{S}_{done}$  do
7:     Let  $\delta \in \nabla_{k-1}$  and  $\delta \notin \mathcal{S}_{done}$  ;
8:     Create a new node  $\omega$  with no child;
9:      $\nabla_k \leftarrow \nabla_k \cup \omega$  ;
10:     $F_*(\delta) \leftarrow \omega$  ; {implicitly  $F_*^{-1}(\omega) \leftarrow F_*^{-1}(\omega) \cup \{\delta\}$ }
11:    while  $F_*^{-1}(\omega) \cap \nabla_{k-1} \not\subseteq \mathcal{S}_{done}$  do
12:      Let  $\gamma \in F_*^{-1}(\omega) \cap \nabla_{k-1}$  and  $\gamma \notin \mathcal{S}_{done}$  ;
13:       $\mathcal{S}_{done} \leftarrow \mathcal{S}_{done} \cup \{\gamma\}$  ;
14:      for all  $\beta \in \text{CHILDREN}(\gamma)$  do
15:         $\text{CHILDREN}(\omega) \leftarrow \text{CHILDREN}(\omega) \cup F_*(\beta)$  ;
16:        for all  $\alpha \in F_*^{-1}(F_*(\beta))$  do
17:          if  $\alpha \in \Delta_{k-1}$  then
18:             $F_*(\alpha) \leftarrow \omega$  ;
19:          else
20:             $F_*(\text{ANCESTOR}(\alpha)) \leftarrow \omega$  ;
21:          end if
22:        end for
23:      end for
24:    end while
25:  end while
26:  if  $k = 1$  then
27:    for all  $\gamma \in \Delta_0$  with  $F_*(\gamma)$  not set do
28:      Let  $\beta \in \text{CHILDREN}(\gamma)$  ;
29:       $F_*(\gamma) \leftarrow F_*(\beta)$  ;
30:    end for
31:  end if
32: end for

```

---

### 3.2 Erweiterung des evolutionären Modells

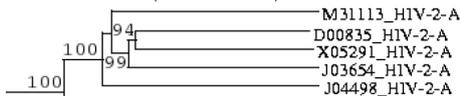
Optionen zur Erweiterung:

- Möglichkeit auch nicht perfekte Übereinstimmungen zu berücksichtigen
- Inversionen berücksichtigen:  
Invertierte Strings in Menge von überlappenden N-Wörtern übernehmen
- weitere Möglichkeiten wären denkbar (ggf. Permutationen betrachten),  
aber: Rechenzeit wird dadurch erhöht

## 4 Leistung bei HIV/SIV Subtyping

Zunächst wurde das N-lokale Dekodieren auf 70 HIV/SIV Genome angewendet, wobei 4 davon unvollständig und die anderen komplett waren. Es ergaben sich sehr gute Übereinstimmungen der Topologie des resultierenden Baumes mit existierendem Wissen. Der Baum beinhaltet auch bekannte evolutionär signifikante Ereignisse der Entwicklung der HIV/SIV Subgruppen. Jedoch ist die Topologie des Baumes abhängig von der Wahl des Parameters N, wobei für N=13-29 die Topologie gefunden wurde, welche identisch zu dem HIV/SIV Sequenz Compendium 2000 ist [4]. Man kann dadurch feststellen, dass durch geeignete Wahl von N sehr gute Ergebnisse bei den Sequenzvergleichen erzielt werden können, wenn es auf komplette Genome angewendet wird.

Ausschnitt: (mit N=15)



Der Ausschnitt enthält die Bezeichnungen der einzelnen Subgruppen und einen Bootstrap-Wert für jede Verzweigung, wobei nur Werte > 50% aufgeführt werden.

### 4.1 Kurze Sequenzen: HIV/SIV *gag*, *pol*, *env* und *nef*

Nun wird das N-lokale Dekodieren auf kurze Sequenzen innerhalb der HIV/SIV Genome angewandt. Dabei sind *gag* (1473-1569 bp), *pol* (2993-3360 bp) *env* (2499-2658 bp) Gene, die Informationen über die strukturellen Proteine für neue Viruspartikel codieren. Das Gen *nef* (292-783 bp) hingegen codiert ein Protein, welches für die effiziente Vermehrung des Virus verantwortlich ist. Es wurden *gag*, *pol*, *env* und *nef* aus 66 HIV/SIV Sequenzen verwendet und die Ergebnisse des lokalen N-Dekodierens stimmen gut mit etablierten HIV/SIV Phylogenieebäumen für diese Regionen überein. Nur bei *nef* treten einige Unterschiede auf, welche möglicherweise darauf zurückzuführen sind, dass beim N-lokalen Dekodieren mehrdeutige Regionen im Gegensatz zum multiplen Alignment berücksichtigt werden.

Als empirisch gute Werte für N haben sich herausgestellt:

*gag* - N=11-23; *pol* - N=11-30; *env* - N=12-24 und *nef* - N=11-20

### 4.2 Non-coding Long Terminal Repeats (LTR)

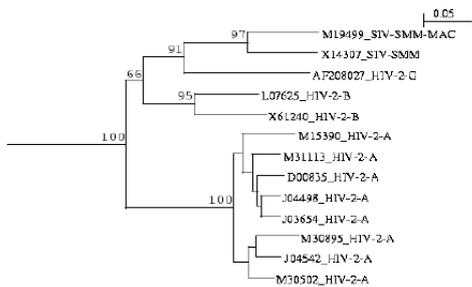
Die Enden von jedem HIV/SIV RNA-Strang enthalten eine RNA-Sequenz, die Long Terminal Repeat (LTR) genannt wird. Die Regionen in den LTR agieren als Schalter zur Kontrolle der Produktion neuer Viren und können sowohl von HIV als auch von der Wirtszelle angesteuert werden. Es wurden 43 non-coding LTR-Regionen untersucht, welche viele Duplikationen, Insertionen sowie Deletionen beinhalten und somit schwer mit herkömmlichen alignment-basierten Vergleichen analysiert werden können. Deshalb sind keine Referenzbäume bzw. Einteilungen konkret vorhanden, sodass versucht wurde, diese mit Hilfe von CLUSTAL-W sowie DIALIGN-2 zu erstellen. Jedoch brachte der Abgleich der gewonnenen Bäume mit bekanntem Wissen schlechte Ergebnisse, da CLUSTAL-W zu schlechter Trennung der HIV-1, HIV-2 und SIV-Sequenzen führte und DIALIGN-2 die Subtypen von HIV-1 schlecht auftrennte (z.B. HIV-1 N ist innerhalb von HIV-1 M). Die Performance von alignment-basierten Methoden war erwartungsgemäß schlecht.

Durch N-lokales Dekodieren mit N=10-21 wurden die HIV/SIV Subtypen gut voneinander getrennt, wobei jedoch Referenzen fehlen, um den erstellten Baum genauer und besser einzuschätzen.

## 5 Abschließende Bemerkungen

- Erweiterung auf komplexeres evolutionäres Modell wäre sinnvoll
- Wahl von Parameter N entscheidend für sinnvolle Ergebnisse:  
N=13-20 sind gute empirische Standardwerte,  
aber: weitere Untersuchungen notwendig

Ausschnitt aus dem Baum von N-lokalen Dekodieren mit  $N=11$ :



- Berücksichtigung auch von mehrdeutigen Regionen (z.b. LTR)  
aber: Problem der Validierung
- Implementation des verbesserten Algorithmus verfügbar [5]

## Literatur

- [1] Didier G, Debomy L, Pupin M, Laprevotte I,  
*Comparing sequences without using alignments: application to HIV/SIV subtyping*, BMC Bioinformatics. 2007 Jan 02
- [2] Didier G, Laprevotte I, Pupin M, Hénaut A,  
*Local decoding of sequences and alignment-free comparison*, J Comp Biol. 2006, 13: 1465-1476
- [3] Felsenstein J,  
*PHYMLIP (Phylogeny Inference Package) version 3.6*,  
Department of Genome Sciences University of Washington, Seattle. 2005
- [4] *The 2000 HIV Sequence Compendium*,  
<http://www.hiv.lanl.gov/content/hiv-db/COMPENDIUM/2000/HIV12SIVcomplete.pdf>
- [5] *SCeNE - Prototypimplementierung des N-lokalen Dekodierens*,  
<http://iml.univ-mrs.fr/~didier/laprevot/>