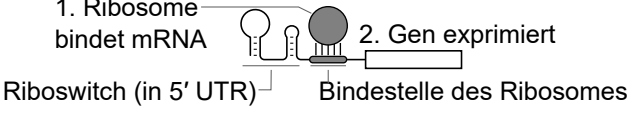
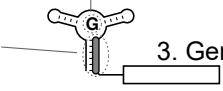


1 Über Riboswitches

Guanin-Riboswitches

Zustand	Problem	Was passiert	Ergebnis
• Zu wenig Guanin	Zelle verhungert	1. Ribosome bindet mRNA 	Mehr Guanin
• Genug Guanin	Gen-Expression vergeudet Energie	1. Guanin bindet Riboswitch, ändert Struktur 2. Ribosom bindet nicht 	Guanin bleibt, spart Energie

2 Definition Paar-HMM

Ein Paar-HMM besteht aus:

- Zuständen S_1, S_2, \dots, S_n
 - Der Startzustand ist S_1
- Übergangswahrscheinlichkeiten $T(i, j)$ (Wkt. des Übergangs von Zustand S_i auf S_j)
- 2 Ausgabealphabeten für die 2 Sequenzen. Die Alphabete können sich gleichen (z.B. beide entsprechen den 20 herkömmlichen Aminosäuren) oder abweichen (z.B. das 1. Alphabet entspricht den 20 Aminosäuren, das 2. Alphabet entspricht den 4 Nucleotiden)
 - Das Paar-HMM kann auch ein Gap in die linke bzw. rechte Sequenz ausgeben, was wir mit dem Symbol ‘-’ schreiben.
 - Das Paar-HMM kann auch gleichzeitig in beide Sequenzen ein Gap ausgeben. Aber wir gehen davon aus, dass alle doppelten Gaps automatisch gelöscht werden, als wäre nichts ausgegeben.

- Ausgabewahrscheinlichkeiten $E(i, c, d)$ (Wkt. in Zustand S_i Symbol c in die linke Seq. und Symbol d in die rechte Seq. auszugeben) c oder d könnte auch ein Gap sein. z.B. $E(i, M, -)$ entspricht der Wkt. ein ‘M’ in die linke Sequenz und ein Gap in die rechte Sequenz in Zustand S_i auszugeben.

Im Context des Viterbi-Algorithmus, definieren wir:

- Die Eingabesequenzen: x^L (linke Seq.) und x^R (rechte Seq.)
- Tabelle $A(i, k, l)$: Wkt. des wahrscheinlichsten Pfades, der
 - mit dem Startzustand S_1 anfängt,
 - auf den Zustand S_i endet,
 - die ersten k Symbole von x^L ausgibt (d.h. $x_1^L x_2^L \dots x_k^L$), und
 - die ersten l Symbole von x^R ausgibt (d.h. $x_1^R x_2^R \dots x_l^R$)

3 Viterbi-Algorithmus mit Scores/Logarithmen

Viterbi-Algorithmus für Paar-HMM mit Scores oder Logarithmen. Das heißt, wir gehen davon aus, dass die Werte in $T(i, j)$ und $E(i, c, d)$ Scores sind, die zusammen addiert werden sollten.

Der hier dargestellte Algorithmus kann mit dem Fall nicht umgehen, in dem sowohl c als auch d Gaps sind. Doch es ist möglich, einen solchen Algorithmus zu machen.

- 1: {Zunächst alle Zellen bekommen den Wert Null}
- 2: **for all** $i \in \{1, 2, \dots, n\}, k \in \{0, \dots, |x^L|\}, l \in \{0, \dots, |x^R|\}$ **do**
- 3: $A(i, k, l) = -\infty$
- 4: **end for**
- 5: {Initialisierung vom Startzustand}
- 6: $A(1, 0, 0) = 0$
- 7: {Initialisierung erster Zeile und erster Spalte in A .}
- 8: **for** $k = 1, 2, \dots, |x^L|$ **do**
- 9: **for** $i = 2, \dots, n$ {Startzustand S_1 ist hier nicht nötig} **do**
- 10: $A(i, k, 0) = \max_{j \in \{1, \dots, n\}} A(j, k-1, 0) + T(j, i) + E(i, x_k^L, -)$ {das entspricht dem normalen Viterbi-Algorithmus}
- 11: **end for**
- 12: **end for**
- 13: **for** $l = 1, 2, \dots, |x^R|$ **do**
- 14: **for** $i = 2, \dots, n$ **do**
- 15: $A(i, 0, l) = \max_{j \in \{1, \dots, n\}} A(j, 0, l-1) + T(j, i) + E(i, -, x_l^R)$
- 16: **end for**

```

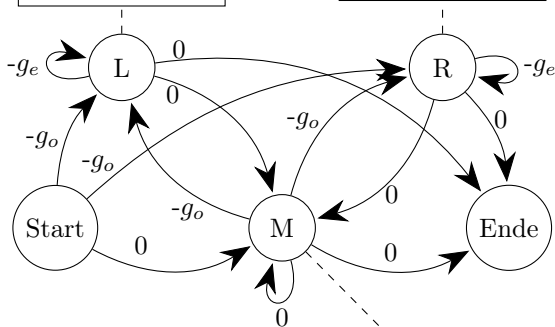
17: end for
18: {Endlich ergänzen wir den Rest der Tabelle A.}
19: for k = 1,2,...,|xL| do
20:   for l = 1,2,...,|xR| do
21:     for i = 1,...,n do
22:       Folgende Formel erinnert an Algorithmen wie Needleman-Wunsch oder Smith-Waterman. Das ist kein Zufall.
23:       A(i, k, l) = maxj ∈ {1,...,n} {
                A(j, k - 1, l - 1) + T(j, i) + E(i, xkL, ylR)
                A(j, k, l - 1) + T(j, i) + E(i, -, ylR)
                A(j, k - 1, l) + T(j, i) + E(i, xkL, -)
                }
24:     end for
25:   end for
26: end for
    
```

4 Paar-HMM: Gotoh-Algorithmus für Proteine

$\sigma(x, y)$ ist eine Score-Funktion für Aminosäuren x und y . Z.B. $\sigma(x, y) = +1$ falls $x = y$ und $\sigma(x, y) = -1$ falls $x \neq y$.

Ausgabescorcs:	
$\binom{A}{-}$	0
$\binom{C}{-}$	0
$\binom{D}{-}$	0
$\binom{E}{-}$	0
...	
Sonst	$-\infty$

Ausgabescorcs:	
$\binom{-}{A}$	0
$\binom{-}{C}$	0
$\binom{-}{D}$	0
$\binom{-}{E}$	0
...	
Sonst	$-\infty$



Ausgabescorcs:	
Score $\binom{x}{y}$	=
Falls $x = \text{gap}$ or $y = \text{gap}$,	$-\infty$.
Sonst,	$\sigma(x, y)$

5 Genetischer Code

Diese Codone benutzen wir in der Vorlesung:

Nuk.	Aminosäure	
	Buchstabe	Name
AAA	K	Lysin
AAC	N	Asparagin
ACA	T	Threonin
ACC	T	Threonin
CAA	Q	Glutamin
CCA	P	Prolin
CCC	P	Prolin

6 Beispiel bzgl. (Bundschuh, 2004)

Organismus	Beschreibung	Sequence
Schleimpilz-2	Proteinseq. <i>nad 2</i>	KTQ (In Eingabe)
Schleimpilz-2	Proteinseq. <i>nad 2</i> mit Gaps	--K--T--Q
<i>P. polycephalum</i>	Proteinsequenz von vorhergesagter mRNA-Sequenz	--K--T--P
<i>P. polycephalum</i>	mRNA-Sequenz mit vorhergesagten C-Insertionen	AAAACACCAAA
<i>P. polycephalum</i>	DNA-Sequenz vom Genom mit Gaps	AAAA-A--AAA
<i>P. polycephalum</i>	DNA-Sequenz vom Genom	AAAAAAAAA (In Eingabe)

Wir alignen eine Nukleotidsequenz mit einer Proteinsequenz, und müssen berücksichtigen: (1) den genetischen Code, (2) C-Insertionen und (3) normale Gaps wie im Gotoh-Algorithmus.