

Freie Universität



Berlin



MAX-PLANCK-GESellschaft

RNA folding

W1-High-throughput Genomics, FU Berlin
OWL RNA Bioinformatics, MPI Molgen Berlin
30.11.16



- Einführung von RNA-Molekülen
- Konzept der RNA-Sekundärstruktur
- Lernen wie eine RNA-Sekundärstruktur grafisch darzustellen
- RNA-Strukturvorhersage
 - **Nussinov algorithmus**
 - Die Motivation
 - Wie dynamische Programmierung in diesem Fall angewendet wird
 - Praktisches Beispiel

RNA Moleküle

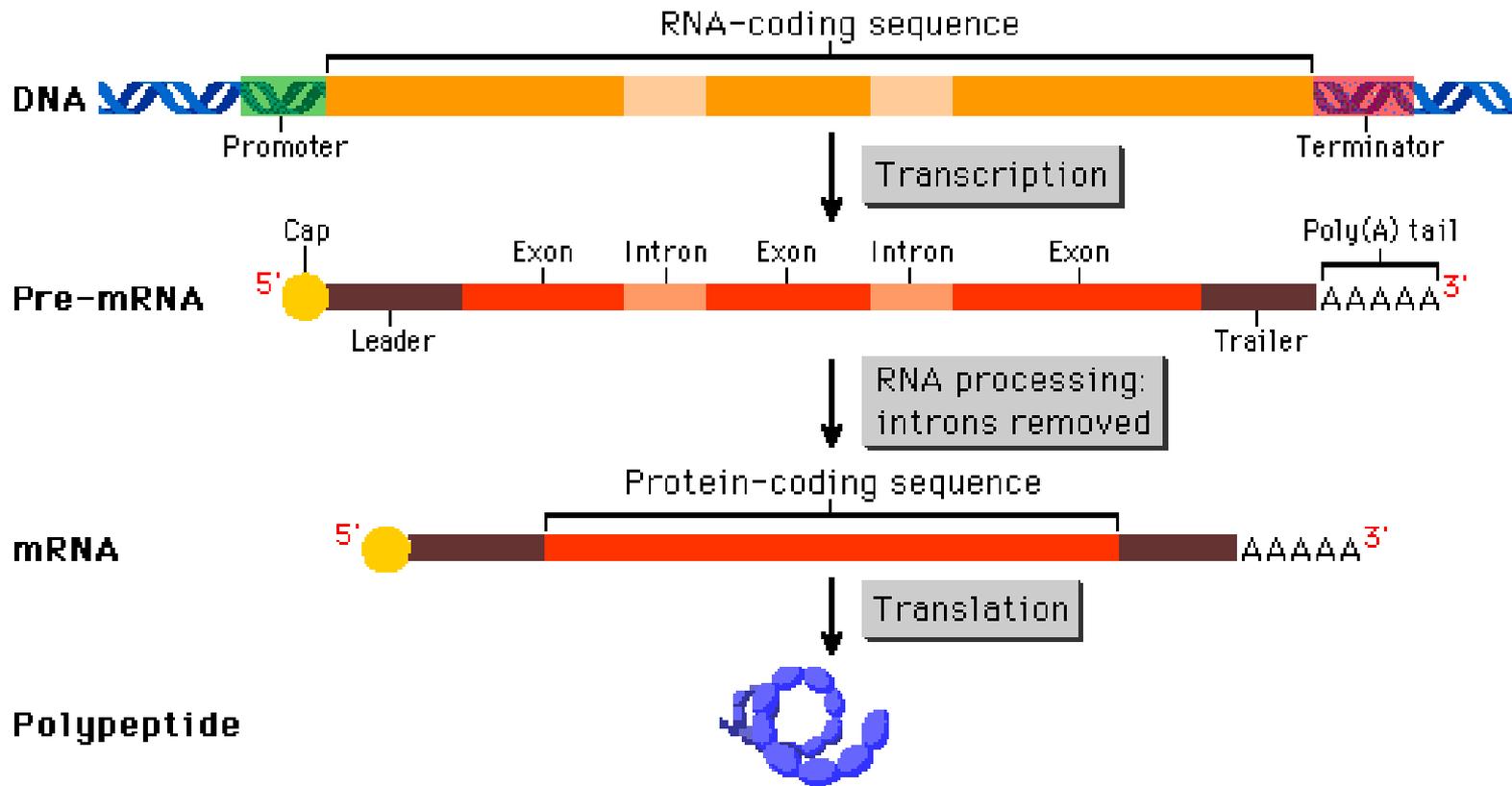


RNA, DNA und Proteine sind die Grundmoleküle des Lebens

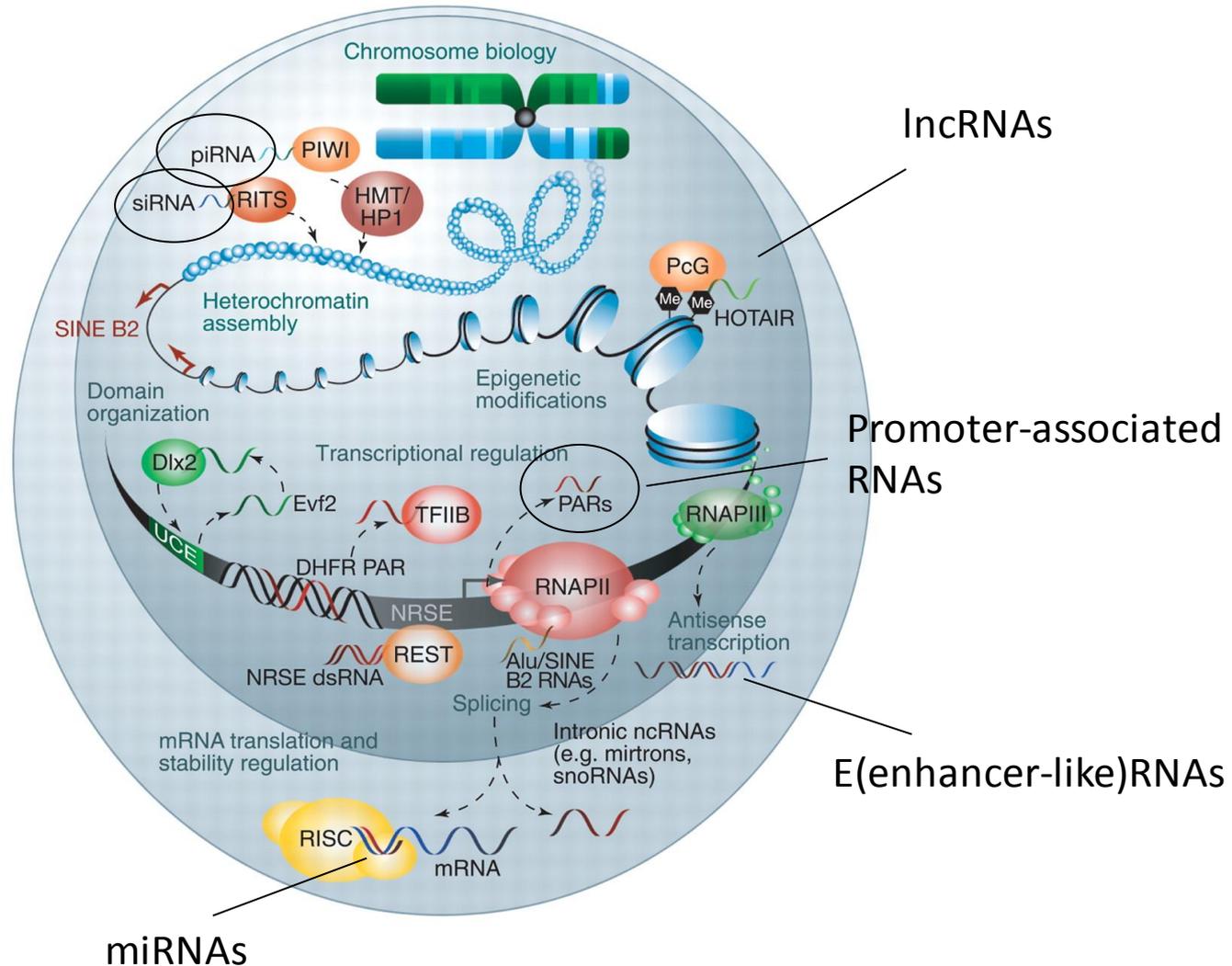
- DNA wird benutzt um die genetische Information zu speichern und replizieren
- Proteine sind die Bausteine der Zelle
- RNAs sind die Zwischenstufen zwischen DNA und Proteinen -> Transkription

Nach der **RNA-Welt-Hypothese** basierte das Leben ursprünglich auf RNA und im Laufe der Zeit delegierten die RNAs das Datenspeicherproblem zur DNA und die katalytische Funktionalität zu Proteinen.

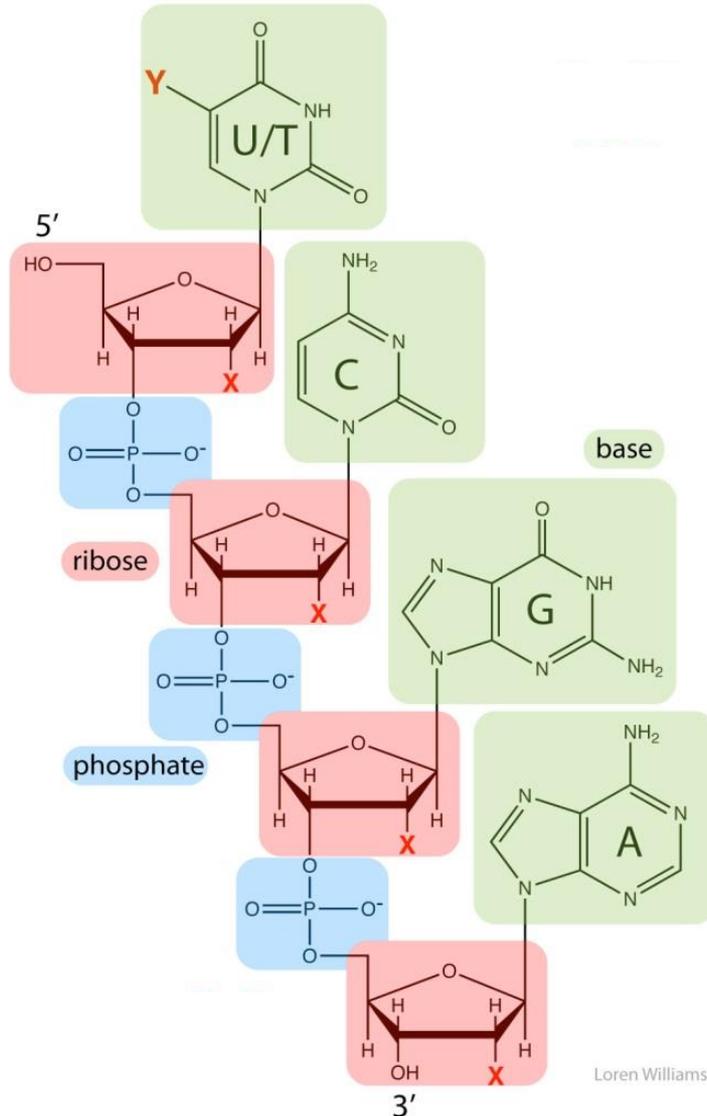
Beispiel: Der Prozess der Transkription



Die wachsende Bedeutung der RNAs: die RNA-Welt



RNA Rückgrat



Ein RNA-Molekül ist ein Polymer aus vier Arten von Ribonukleotiden, jeweils durch eine der vier Basen angegeben.

A -> adenine

C -> cytosine

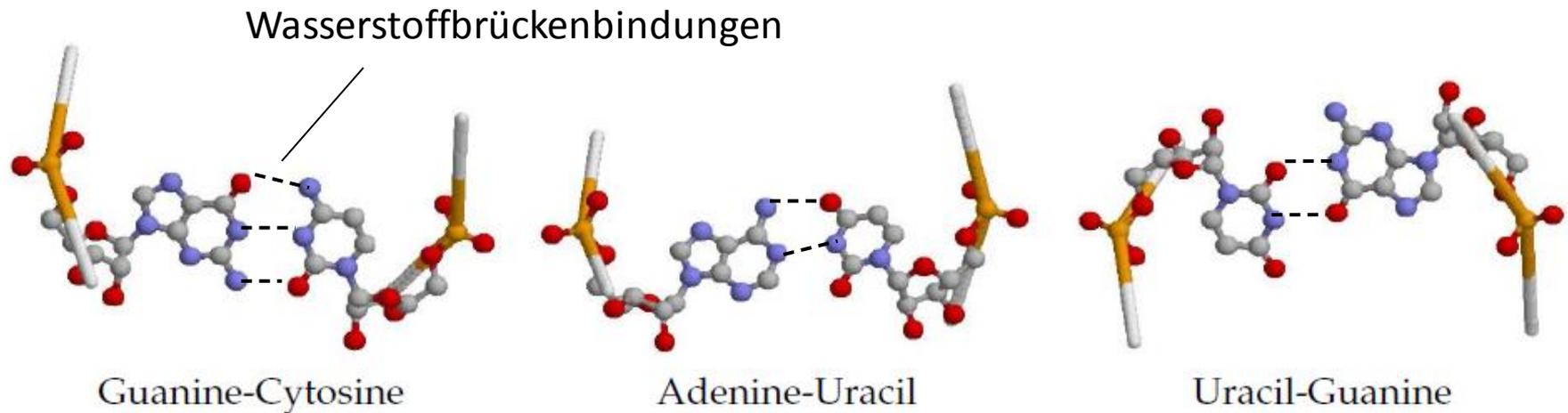
G -> guanine

U(T) -> uracil

RNA-Sekundärstruktur



Im Gegensatz zu DNA, ist RNA einzelsträngig. Allerdings formen die komplementären Basen C-G und A-U stabile Basenpaare über Wasserstoffbrückenbindungen. Diese werden Watson-Crick-Paare genannt. Wichtig sind auch die schwächeren U-G Wobble-Paare. Zusammen werden sie „kanonische Basenpaare“ genannt.

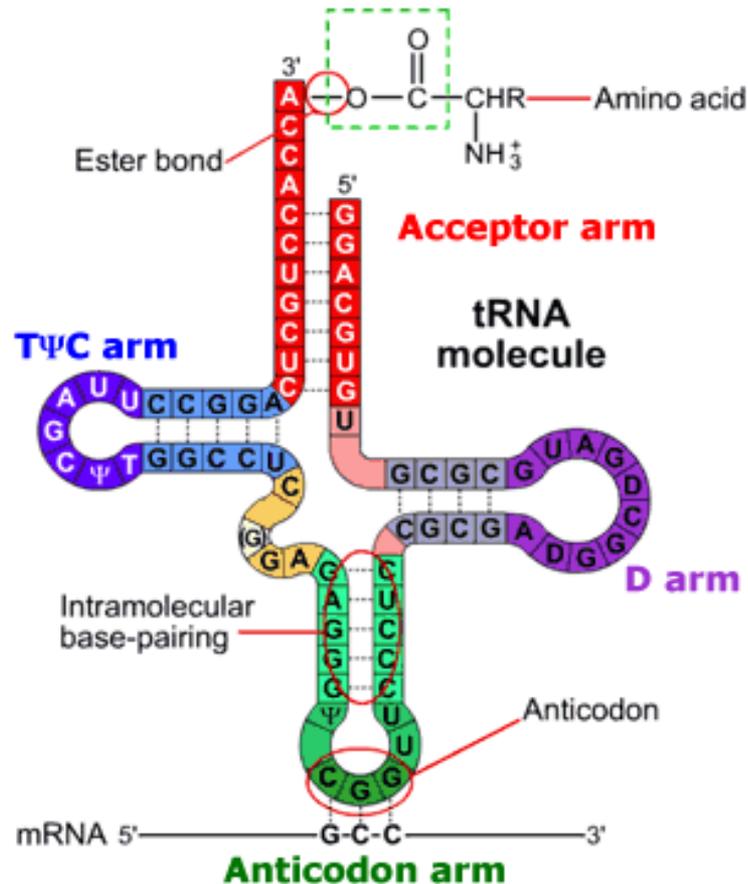


RNA-Sekundärstruktur



Die Basenpaare, die zwischen den verschiedenen Teilen eines RNA-Molekül gebildet werden, definieren die Sekundärstruktur des RNA-Moleküls.

Hier ist die Sekundärstruktur von einem tRNA:

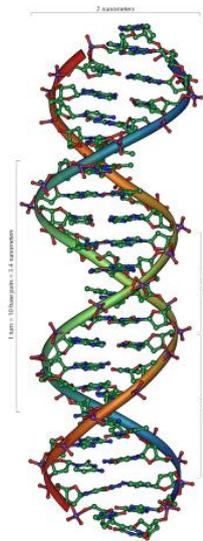


Sekundärstruktur: Satz von Basenpaaren, die auf einer Ebene abgebildet werden können

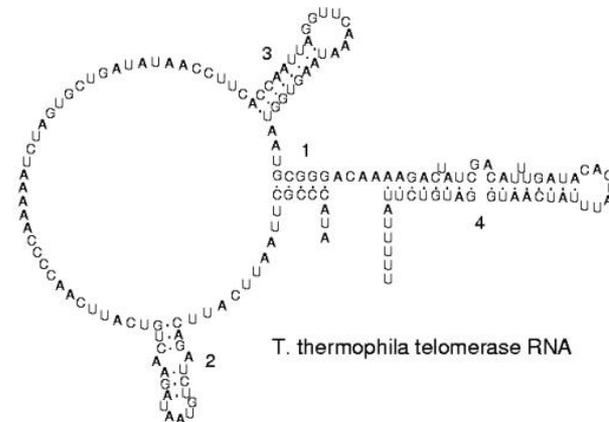
Zusammenfassen: DNA vs RNA



1. DNA doppelsträngig, **RNA einzelsträngig**, oft Sekundärstruktur durch intramolekulare Wasserstoffbrücken
2. **Länge:** DNA Millionen von Basenpaaren, RNA ~20 bis mehrere tausend Nukleotide
3. Biochemie: RNA Ribose/ Uracil (U), DNA 2-Desoxyribose / Thymin (T)
4. Biologischen Rollen: DNA Trägerin der Erbinformation, RNA unterschiedliche Funktionen. Im Gegensatz zur DNA spielt die Struktur der RNA bei deren Funktion eine wesentliche Rolle.



DNA



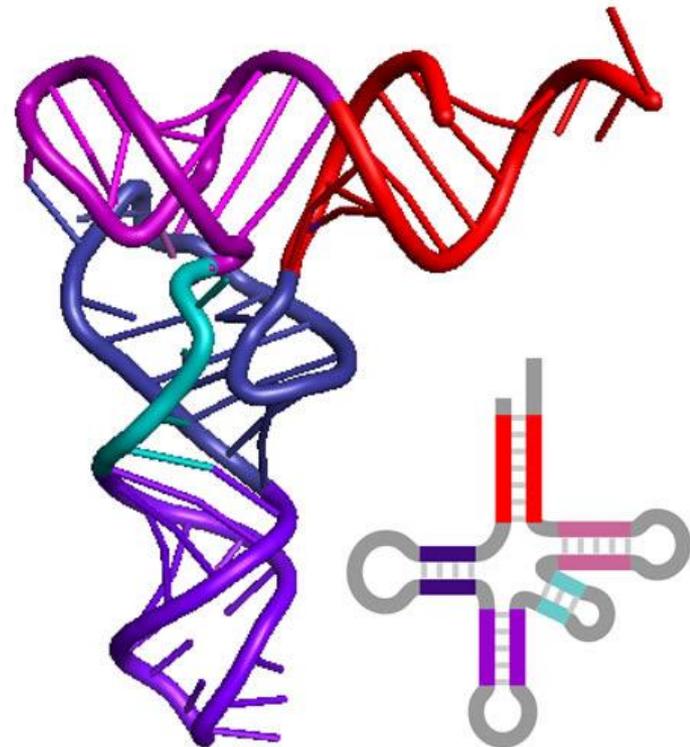
T. thermophila telomerase RNA

RNA

Struktur Konformationen von RNA



- Primärstruktur: Sequenz von Monomeren ATGCCGTCAC..
- Secondärstruktur: 2D-Faltung, durch Wasserstoffbrückenbindungen definiert
- Tertiärstruktur: 3D-Faltung
- Quartärstruktur: komplexe Anordnung von mehreren gefalteten Moleküle





- Die RNA Sequenz faltet sich selbst zurück aufgrund der Komplementarität der Basen.
- RNA-Strukturbestimmung ist experimentell schwierig, **daher ein wichtiges Thema für die Bioinformatik.**
- Die Methoden zur 2D-Vorhersage können in zwei Gruppen eingeteilt werden:
 - Methoden, die ihre Vorhersagen aus MSA ableiten
 - Methoden, die ihre Vorhersagen für **einzelne Sequenzen bestimmen** (Maximierung Anzahl Basenpaare oder Minimierung freie Energie)



Die echte Sekundärstruktur eines RNA-Moleküls ist die Menge der Basenpaare im dreidimensionalen Raum.

Definition Für unsere Zwecke ist ein RNA-Molekül einfach ein String

$$x = (x_1 \cdot x_2, \dots, x_L)$$

mit $x_i \in \{A, C, G, U\}$ für alle i

Definition Eine Sekundärstruktur für x ist eine Menge P von geordneten Basenpaaren (i, j) , mit $1 \leq i \leq j \leq L$ und:

1. Primary proximity Constraint:

$|i-j| > 3$, d.h. die Basenpaare dürfen nicht zu nahe beieinander liegen

2. Nesting Constraint:

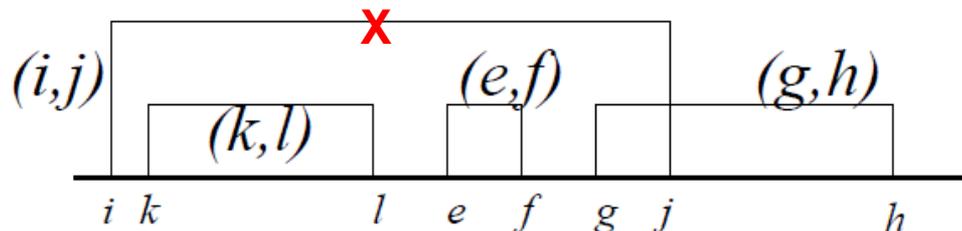
Sind (i, j) und (i', j') zwei Paare von Nukleotiden, wobei $i < i' < j$, dann gilt $j < j'$.
d.h. die Basenpaare dürfen sich nicht überschneiden.



Verschachtelte Strukturen

Im Folgenden werden wir nur verschachtelten Sekundärstrukturen behandeln, da komplizierteren nicht-verschachtelten Strukturen nicht mehr mit unsere Methoden lenkbar sind.

Diese "Schachtelungsbedingung" verbietet überkreuzte Wasserstoffbrücken, erlaubt dagegen geschachtelte Wasserstoffbrücken. Überkreuzte Wasserstoffbrücken, so genannte Pseudoknoten, kommen relativ selten vor. Algorithmen, welche Pseudoknoten zulassen, sind wesentlich weniger effizient als solche, die sie verbieten.



Strukturvorhersage durch Maximierung Von Basenparungen



- Lass uns mit diese Sequenz spielen: **UUGACAUCG**
- Ziel: die Sekundärstruktur mit der maximalen Anzahl an Basenpaaren finden, wobei zwischen zwei paarenden Basen mindestens eine Ungepaarte stehen soll



RNA-Struktur: Klammern und Punkte



Wir können die RNA-Struktur als Strings mit balanzierten Klammern und Punkten mit der entsprechenden Schachtelebene (nesting level) darstellen

(. .) (. . .)
UUGACAUCG
└───┘ └───┘

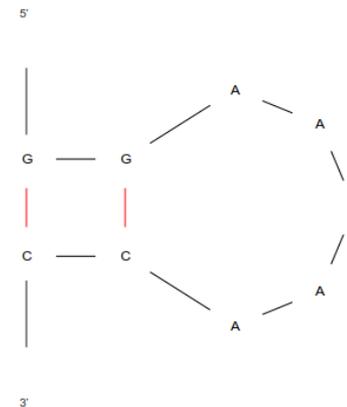
UUGACAUCG
└───┘ └───┘

(. (. . .))



- Die Sekundärstruktur kann von der Sequenz und Klammerndarstellung ermittelt werden
- Beispiel:

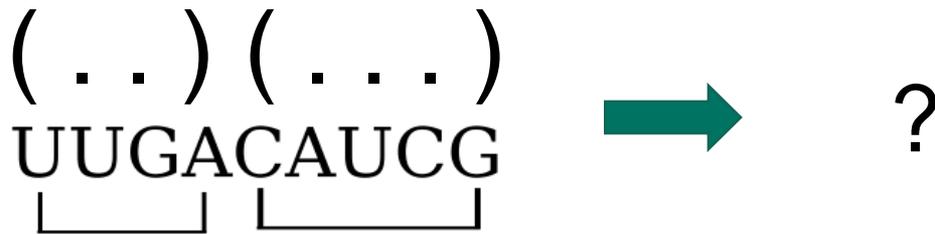
GGAAAACC
((. . . .))



Graph Darstellung



- Die Sekundärstruktur kann von der Sequenz und Klammerndarstellung ermittelt werden

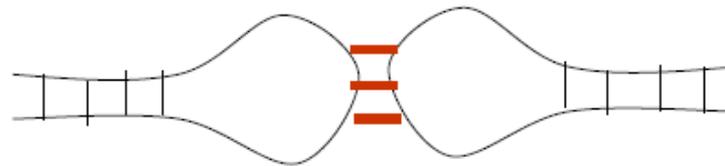
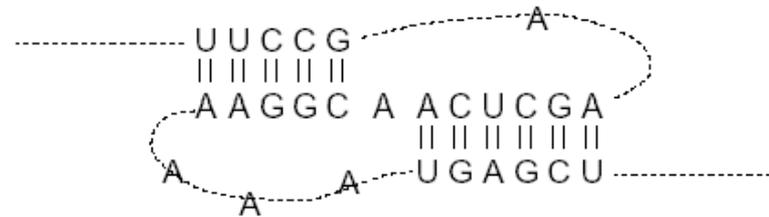


Nicht-verschachtelte Interaktionen



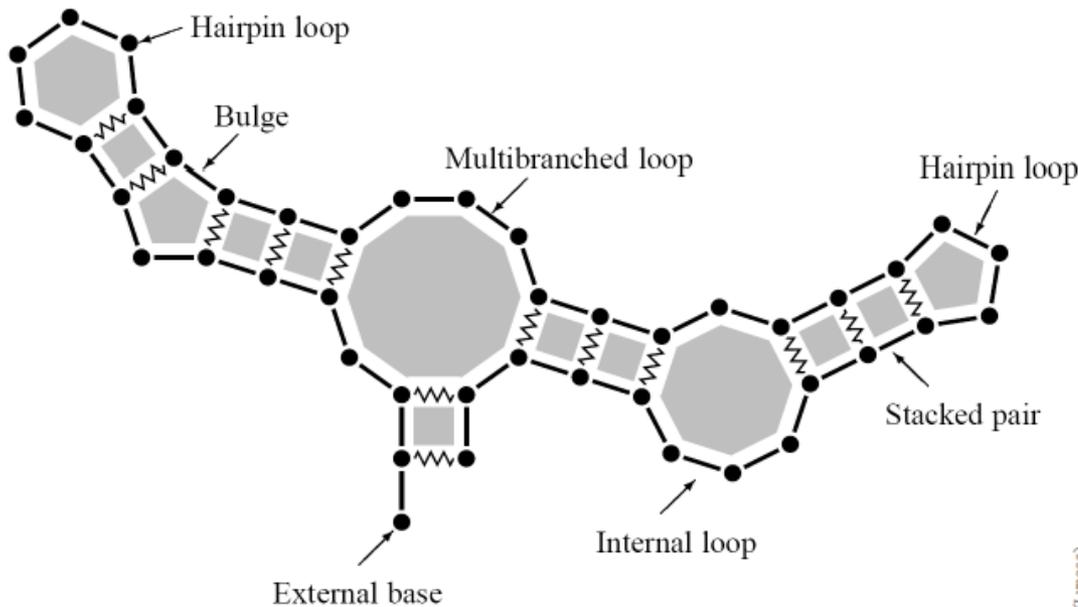
Interaktionen, die nicht verschachtelt sind, führen zu einer ‚Pseudoknoten‘-Struktur oder ‚kissing hairpins‘: Segmente der Sequenz, die in die gleiche Richtung verbunden sind oder dreidimensionale Kontakte haben:

Pseudoknot



Kissing hairpins

RNA-Sekundärstrukturelemente

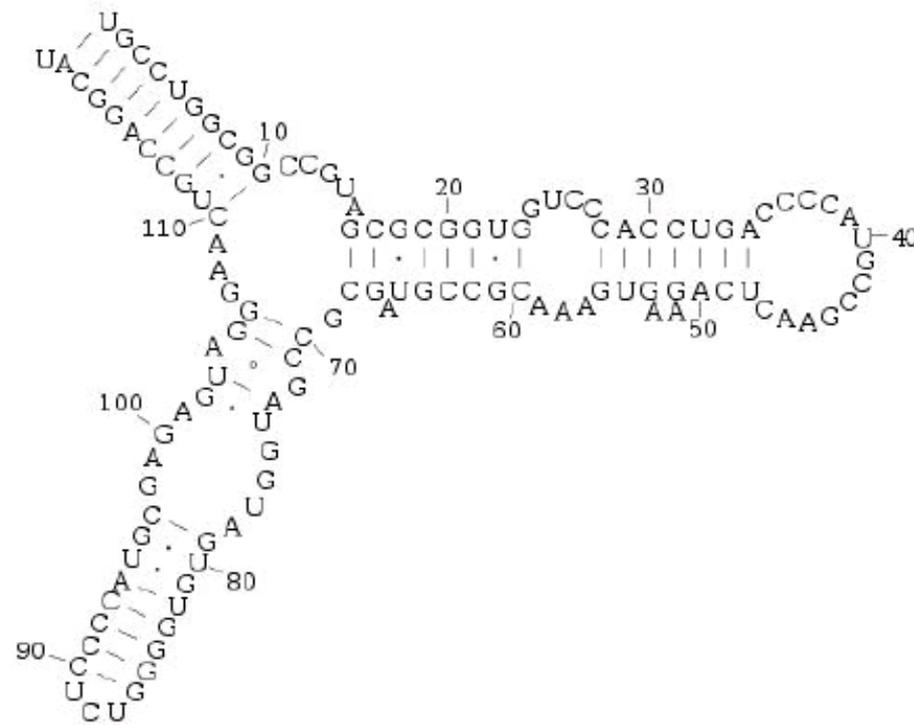


(Lyngsø)

- single stranded RNA
- stacked base pairs
- stem & loop (hairpin loop)
- bulge loop
- interior loop
- junction or multi-loop

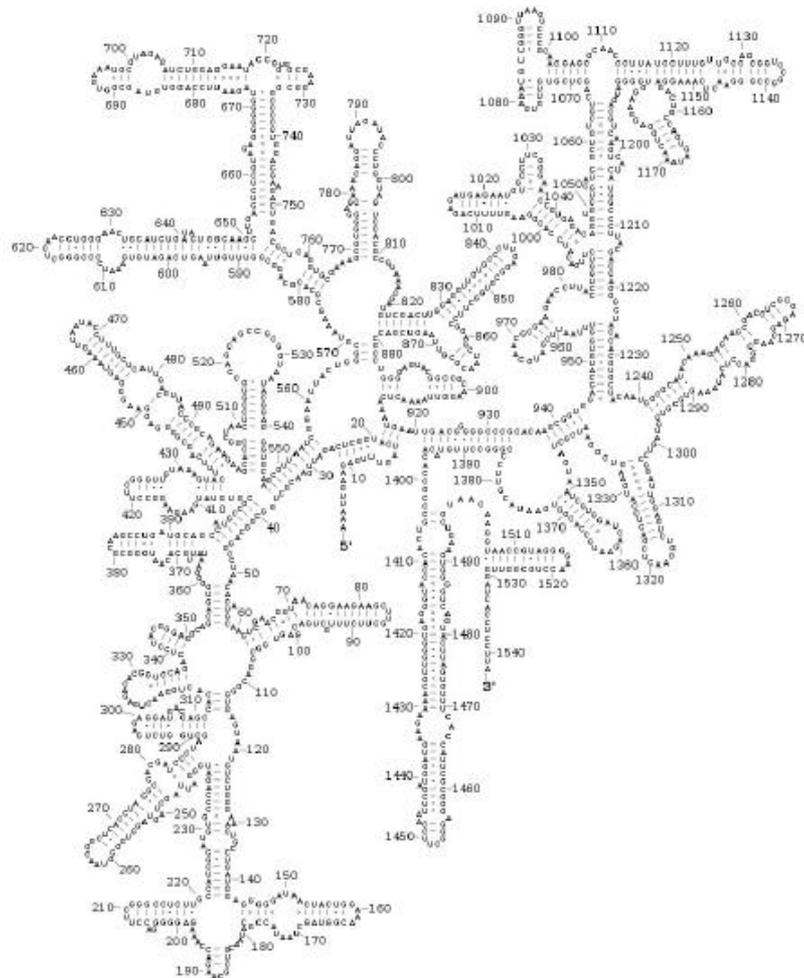
Typen von Einzel- und Doppelstrang-Regionen
Dies nennt man Basenpaar **Graphdarstellung**

RNA-Struktur Beispiel 1: 5s rRNA



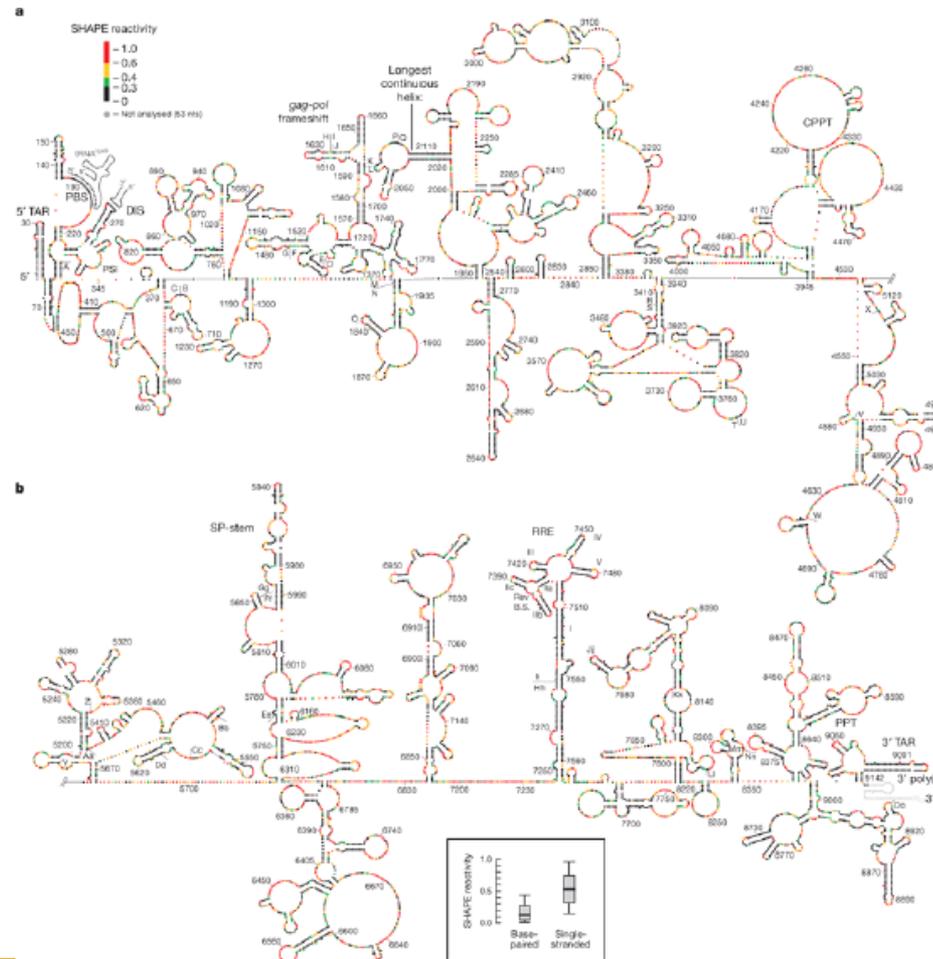
E. coli 5S
120 bases

RNA-Struktur Beispiel 2 : E.coli 16S rRNA



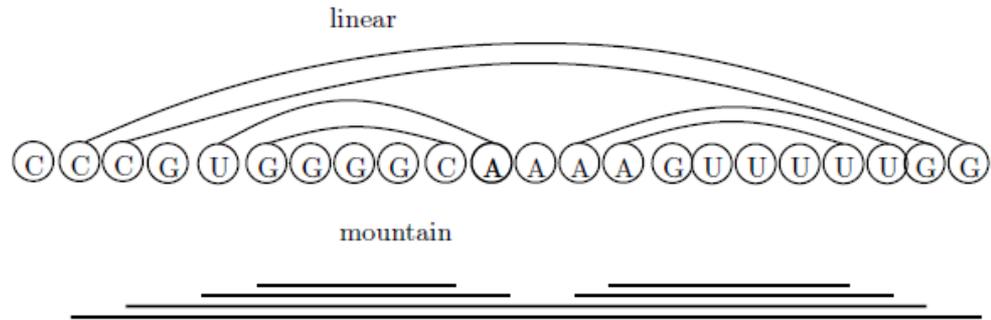
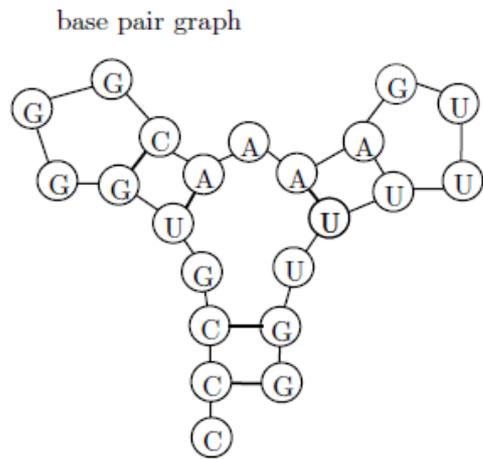
1542 bases

RNA-Struktur Beispiel 1 : HIV



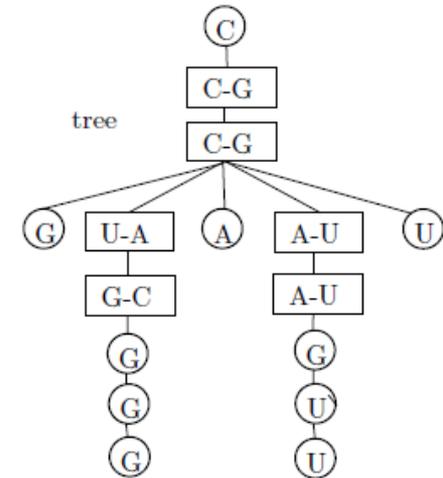
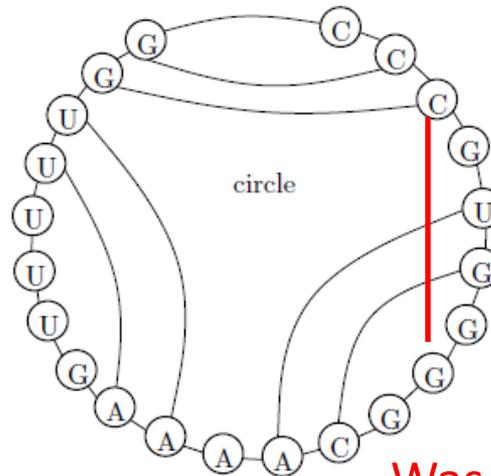
9173 basis

Darstellungstypen einer RNA-Sekundärstruktur



bracket

$(((((**)))*((**))*))$



Was ist das?



Ansatz:

Finde eine Konfiguration, die die **Anzahl von Basenpaaren maximiert**

Für eine Sequenz der Länge L , wächst die Anzahl der möglichen Konfigurationen exponentiell mit der Länge der Sequenz.

Dabei ist es unmöglich, alle möglichen Strukturen zu aufzählen!

Zum Glück, können wir ***Dynamische Programmierung*** verwenden, um die effizienteste Lösung zu finden.

Eine Methode, um das zu tun, hat Ruth Nussinov in 1978 publiziert.

Der Algorithmus ist rekursiv. Es berechnet die beste Struktur für kleinere Teilfolgen größere Sequenzen. Von der kleinsten bis zur vollständigen Sequenz.

Nussinov Faltungsalgorithmus -I (1978)



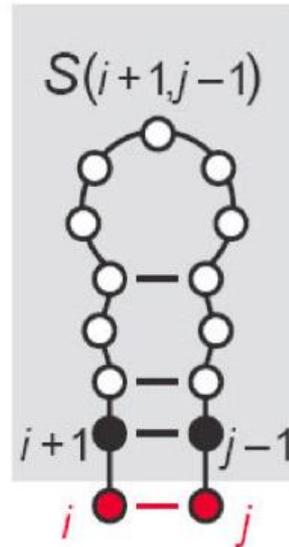
- Eine vereinfachte Version des Nussinov-Algorithmus versucht, **die Anzahl der gepaarten Basen zu maximieren**
- Unser Score: +1 für Basenpaar, 0 für alles Andere
- Sei $x_1 \dots x_L$ ist eine Sequenz von 1 zu L faltenden Nukleotiden.
- Berechne die maximale Anzahl von gebildeten Basenpaaren der Teilfolge $x[i:j]$
- **$S(i,j)$** Max score für die Subsequenz $i, i + 1, \dots, j - 1, j$
- **$S(i,j)$** kann rekursiv berechnet werden (**Dynamic Programming**)
- Und die Struktur rekursiv gefaltet werden
(gegeben das wir bereits für alle kurzen Sequenzen $x[m:l]$ $i < m < l < j$ berechnet haben)

Nussinov Faltungsalgorithmus - I (1978)



Die Struktur auf $x[i:j]$, kann auf vier Arten berechnet werden:

1) Falls i, j ein WC-Basenpaar sind



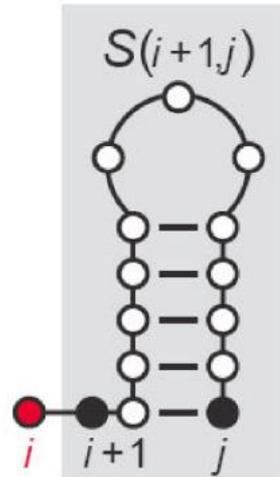
Füge die gepaarten Basen i und j der besten Struktur für Subsequenz $[i+1, j-1]$ hinzu

$$S(i, j) = 1 + S(i + 1, j + 1)$$



Die Struktur auf $x[i:j]$, kann auf vier Arten berechnet werden:

2) Falls i ungepaart bleibt



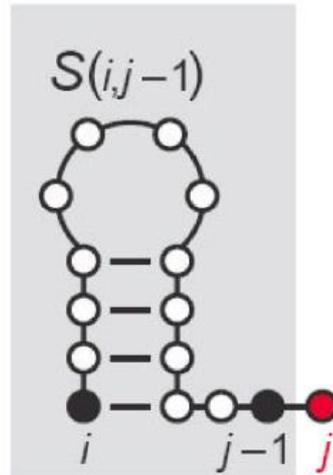
Füge die ungepaarten Base i der besten Struktur für Subsequenz $[i+1, j]$ hinzu

$$S(i, j) = S(i + 1, j)$$



Die Struktur auf $x[i:j]$, kann auf vier Arten berechnet werden:

3) Falls j ungepaart bleibt



Füge die ungepaarten Base j der besten Struktur für Subsequenz $[i, j-1]$ hinzu

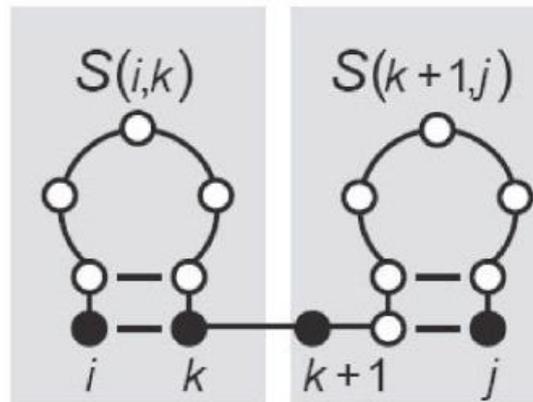
$$S(i, j) = S(i, j - 1)$$

Nussinov Faltungsalgorithmus - I (1978)



Die Struktur auf $x[i:j]$, kann auf vier Arten berechnet werden:

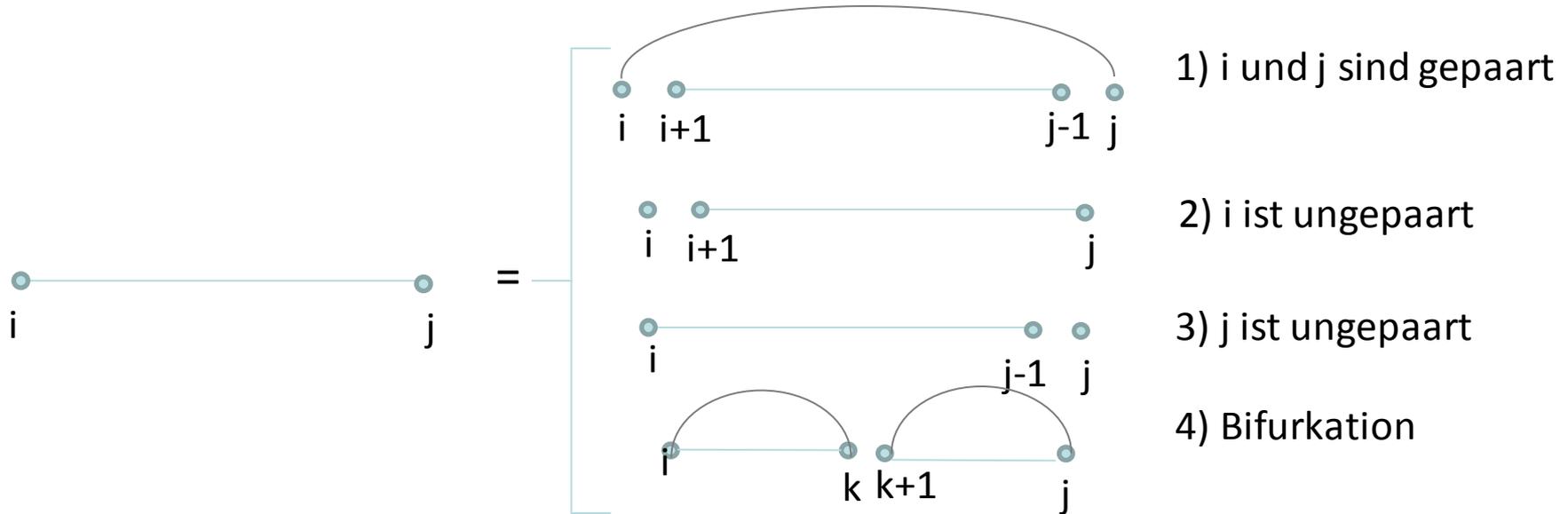
- 4) Falls i, j jeweils mit anderen Nukleotiden gepaart sind, handelt es sich um eine Bifurkation



Struktur $S(i, j)$ besteht dann aus den Strukturen für zwei Subsequenzen i, \dots, k und $k+1, \dots, j$.
(Kombiniere zwei optimale Strukturen $[i,k], [k+1,j]$)

$$S(i, j) = \max_{i < k < j} (S(i, k) + S(k + 1, j))$$

Graphische Darstellung des Faltungsalgorithmus





Der Algorithmus der dynamischen Programmierung hat zwei Phasen:

- Der **fill Schritt**, berechnet rekursiv die besten Scores $S(i, j)$, welche die maximale Anzahl an bp darstellen, die in einer Subsequenz $x_i \dots x_j$ gefunden werden
- Der **traceback Schritt**, traceback durch die berechnete Matrix, um die beste Struktur mit der maximale Anzahl an bp zu finden



Nussinov Algorithmus: Fill Phase

Algorithm (Nussinov RNA folding, fill stage) – Beispiel, Sequenz: GGGAAAUCC

Input: Sequenz $x = (x_1, x_2, \dots, x_L)$

Output: Maximale Anzahl S_{max} von Basenpaare für Sequenz x

Initialisierung der DP Matrix:

```

for i=2 to L do
    S[i, i-1]=0
for i=1 to L do
    S[i, i]=0

```

Rekursion (Fullung der DP Matrix):

```

for n=2 to L do
    for j=n to L do
        i=j-n+1
        S = max {
            S(i + 1, j)
            S(i, j - 1)
            S(i + 1, j - 1) + δ(i, j)
            max_{i < k < j} (S(i, k) + S(k + 1, j))
        }

```

	$j \rightarrow$								
	G	G	G	A	A	A	U	C	C
G	0								
G	0	0							
G		0	0						
A			0	0					
A				0	0				
A					0	0			
U						0	0		
C							0	0	
C								0	0

Die DP Rekursion prüft vier Möglichkeiten (1,2,3,4 von der vorherigen Folie)
 $\delta(i, j)=1$ for base pairs, $\delta(i, j)=0$ otherwise

Nussinov Algorithmus: Fill Phase



Rekursion (Fullung der DP Matrix):

For $n=2$ to L do

 for $j=n$ to L do

$i=j-n+1$

$$S = \max \begin{cases} S(i+1, j) \\ S(i, j-1) \\ S(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} (S(i, k) + S(k+1, j)) \end{cases}$$

$j \rightarrow$

	G	G	G	A	A	A	U	C	C	
G	0	0								
G	0	0	0							
G		0	0	0						
A				0	0					
A					0	0				
A						0	0	?		
U							0	0		
C								0	0	
C									0	0

$i \downarrow$

Nussinov Algorithmus: Fill Phase



Rekursion (Füllung der DP Matrix):

For $n=2$ to L do

 for $j=n$ to L do

$i=j-n+1$

$$S = \max \begin{cases} S(i+1, j) \\ S(i, j-1) \\ S(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} (S(i, k) + S(k+1, j)) \end{cases}$$

$j \rightarrow$

	G	G	G	A	A	A	U	C	C	
G	0	0								
G	0	0	0							
G		0	0	0						
A				0	0					
A					0	0	0	?		
A						0	0	1		
U							0	0		
C								0	0	
C									0	0

$i \downarrow$

Nussinov Algorithmus: Fill Phase



Rekursion (Fullung der DP Matrix):

For $n=2$ to L do

 for $j=n$ to L do

$i=j-n+1$

$$S = \max \begin{cases} S(i+1, j) \\ S(i, j-1) \\ S(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} (S(i, k) + S(k+1, j)) \end{cases}$$

$j \rightarrow$

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0			
G	0	0	0	0	0	0	1		
G		0	0	0	0	0	1	?	
A			0	0	0	0	1	1	
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
C							0	0	0
C								0	0

$i \downarrow$

Nussinov Algorithmus: Fill Phase



Rekursion (Fullung der DP Matrix):

For $n=2$ to L do

 for $j=n$ to L do

$i=j-n+1$

$$S = \max \begin{cases} S(i+1, j) \\ S(i, j-1) \\ S(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} (S(i, k) + S(k+1, j)) \end{cases}$$

$j \rightarrow$

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	
G	0	0	0	0	0	0	1	2	?
G		0	0	0	0	0	1	2	2
A			0	0	0	0	1	1	1
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
C							0	0	0
C								0	0

$i \downarrow$

Nussinov Algorithmus: Fill Phase



Rekursion (Fullung der DP Matrix):

For $n=2$ to L do

 for $j=n$ to L do

$i=j-n+1$

$$S = \max \begin{cases} S(i+1, j) \\ S(i, j-1) \\ S(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} (S(i, k) + S(k+1, j)) \end{cases}$$

$j \rightarrow$

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	?
G	0	0	0	0	0	0	1	2	3
G		0	0	0	0	0	1	2	2
A			0	0	0	0	1	1	1
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
C							0	0	0
C								0	0

$i \downarrow$

Nussinov Algorithmus: Fill Phase



Recursion (Fullung der DP Matrix):

For $n=2$ to L do

 for $j=n$ to L do

$i=j-n+1$

$$S = \max \begin{cases} S(i+1, j) \\ S(i, j-1) \\ S(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} (S(i, k) + S(k+1, j)) \end{cases}$$

	$j \rightarrow$								
	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G		0	0	0	0	0	1	2	2
A			0	0	0	0	1	1	1
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
C							0	0	0
C								0	0

Die letzte Zell (1,9) gibt den Score einer optimalen Struktur für die ganze Sequenz. Die Optimale Score kann hier durch Option 1 und 2 erreicht werden. Warum?

Nussinov Algorithmus: Fill Phase



$j \rightarrow$

	G	G	G	A	A	A	U	C	C
G	0								
G	0	0							
i G		0	0						
A			0	0					
A				0	0				
A					0	0			
U						0	0		
C							0	0	
C								0	0

Achtung! Wir müssen immer noch den Score für Bifurkationen berechnen:
 $k=2,3,4,5,6,7,8$

Nussinov Algorithmus: Fill Phase



$$S(1,9) = \max\{2,3,2,2\}=3$$

Traceback, um die Struktur selbst zu finden

$j \rightarrow$

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G		0	0	0	0	0	1	2	2
A			0	0	0	0	1	1	1
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
C							0	0	0
C								0	0

$i \downarrow$

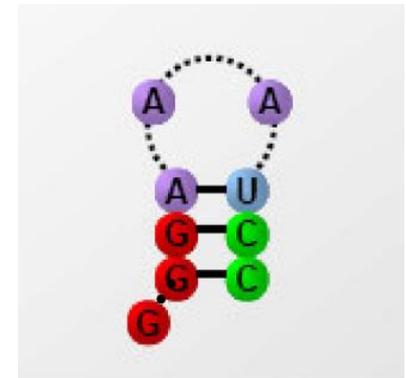
Traceback in die Nussinov Matrix



Bauen die RNA Struktur

- Der Traceback beginnt immer in der rechten oberen Ecke der Matrix und verfolgt den Pfad einer optimalen Struktur
- In diesem Fall ist der Traceback relativ einfach, weil keine verzweigte RNA-Struktur vorliegt
- man beginnt mit Zelle (1, L) die Scores der drei umgebenden zellen (links, schräg links unten und unten)
- vergleichen und so schlussfolgern, wie der Score in der aktuellen Zelle zustande gekommen sein muss und so eine optimale Struktur finden.
- Oft (wie hier) sind viele Pfade mit demselben Score möglich

	$j \rightarrow$									
	G	G	G	A	A	A	U	C	C	
$i \downarrow$	G	0	0	0	0	0	0	1	2	3
	G	0	0	0	0	0	0	1	2	3
	G		0	0	0	0	0	1	2	2
	A			0	0	0	0	1	1	1
	A				0	0	0	1	1	1
	A					0	0	1	1	1
	U						0	0	0	0
	C							0	0	0
	C								0	0



Nussinov - Der Traceback Schritt



Algorithm $\text{traceback}(i, j)$ (Nussinov RNA folding)

Input: Matrix γ and positions i, j .

Output: Secondary structure maximizing the number of base pairs.

Initial call: $\text{traceback}(1, L)$.

if $i < j$ then

 if $\gamma(i, j) = \gamma(i + 1, j)$ then // case (1)

$\text{traceback}(i + 1, j)$

 else if $\gamma(i, j) = \gamma(i, j - 1)$ then // case (2)

$\text{traceback}(i, j - 1)$

 else if $\gamma(i, j) = \gamma(i + 1, j - 1) + \delta(i, j)$ then // case (3)

 print base pair (i, j)

$\text{traceback}(i + 1, j - 1)$

 else for $k = i + 1$ to $j - 1$ do // case (4)

 if $\gamma(i, j) = \gamma(i, k) + \gamma(k + 1, j)$ then

$\text{traceback}(i, k)$

$\text{traceback}(k + 1, j)$

Komplexität des Nussinov Algorithmus



Recursion:

```
for n = 2 to L do // longer and longer subsequences
  for j = n to L do
    i ← j - n + 1
    γ(i, j) ← max {
      γ(i + 1, j),
      γ(i, j - 1),
      γ(i + 1, j - 1) + δ(i, j),
      maxi < k < j (γ(i, k) + γ(k + 1, j)).
```

Complexity analysis:

- The outer loop (for n = 2 to L) has complexity O(L).
- The middle loop (for j = n to L) has complexity O(L).
- The inner loop (i ← j - n + 1) has complexity O(L).
- The recursive call (max) has complexity O(L).

Wir haben 3 verschachtelte Schleifen, wobei jede davon O(L) Mal ausgeführt wird.

Daraus ergibt sich, dass die Gesamtzahl der ausgeführten Operationen des Algorithmus O(L³) beträgt



- Einführung von RNA-Molekülen
- Konzept der RNA-Sekundärstruktur
- Lernen wie eine RNA-Sekundärstruktur grafisch darzustellen
- RNA-Strukturvorhersage
 - **Nussinov algorithmus**
 - Die Motivation
 - Wie dynamische Programmierung in diesem Fall angewendet wird
 - Praktisches Beispiel



- R. Durbin, S. Eddy, A. Krogh und G. Mitchinson, Biological sequence Analysis, Cambridge, 1998
- Sean R. Eddy: How do RNA folding algorithms work? Nature Biotechnology, Vol 22, Num 11, pages 1457-1458, 2004
- Rune Lyngso, Lecture Notes on RNA Secondary Structure Prediction, 2010