

Algorithmen und Datenstrukturen

Dynamische Programmierung:
Maximale Teilsumme ♦
RNA-Sekundärstruktur-Problem

Noch einmal **Maximale Teilsumme**

Maximale Teilsumme

Name: Maximale Teilsumme

Eingabe: Sequenz L von n ganzen Zahlen

Ausgabe: größte Summe von Elementen einer Teilsequenz

- zwei Algorithmen:
 - Brute-Force: $O(n^3)$
 - Teile und Herrsche: $O(n \log n)$

Maximale Teilsumme (Brute-Force)

```
def maxTeilsumme_1(L):  
    maxSumme = 0 # mindestens 0 (leere Teilsequenz)  
    for i in range(0, len(L)): # untere Grenze  
        for j in range(i, len(L)): # obere Grenze  
            summe = 0  
            for k in range(i, j+1): # Summe bilden  
                summe += L[k]  
            if summe > maxSumme:  
                maxSumme = summe  
    return maxSumme
```

Unnötige Berechnungen

- $i = 0$:
 - $j = 0$: $L[0]$
 - $j = 1$: $L[0] + L[1]$
 - $j = 2$: $L[0] + L[1] + L[2]$
 - $j = 3$: $L[0] + L[1] + L[2] + L[3]$
 - ...
- eigentlich ist die nächste Teilsumme immer die vorherige plus dem nächsten Element der Sequenz
- *Idee*: Speichern der zuvor berechneten Teilsummen in einer Tabelle (Matrix T , wobei an $T(i,j)$ die Teilsumme von $L[i]$ bis $L[j]$ gespeichert wird)

Tabellarischer Brute-Force-Algorithmus

- *Idee*: Speichern der zuvor berechneten Teilsummen in einer Tabelle (Matrix T , wobei an $T(i,j)$ die Teilsumme von $L[i]$ bis $L[j]$ gespeichert wird)

$i \backslash j$	0	1	2	...
0	$L[0]$	$L[0]+L[1]$	$L[0]+L[1]+L[2]$...
1	0	$L[1]$	$L[1]+L[2]$...
2	0	0	$L[2]$...

Tabellierter Brute-Force-Algorithmus

```
def maxTeilsumme_2(L):  
    maxSumme = 0 # mindestens 0 (leere Teilsequenz)  
    T = []  
    for i in range(0, len(L)):  
        T.append([])  
        for j in range(0, i): # j < i  
            T[i].append(0)  
        T[i].append(L[i]) # j = i  
        if T[i][i] > maxSumme: maxSumme = T[i][i]  
        for j in range(i+1, len(L)): # j > i  
            T[i].append(T[i][j-1] + L[j])  
            if T[i][j] > maxSumme: maxSumme = T[i][j]  
    return maxSumme
```

Zeit-Raum-Verschiebung

- *Bisher:* Zeitkomplexität $O(n^3)$ und Platzbedarf für sechs Variablen (*also* $O(1)$) (*ohne Platz für die Eingabe L*)
- *Jetzt:* nur eine Zuweisung für $T(i,j)$ mit $j \leq i$
- nur eine Addition und eine Zuweisung für $T(i,j) \rightarrow T(i,j+1)$
- neuen Wert sofort mit `maxSumme` vergleichen und `maxSumme` ggf. aktualisieren

→ $\Theta(n^2)$ Zeit und Platz

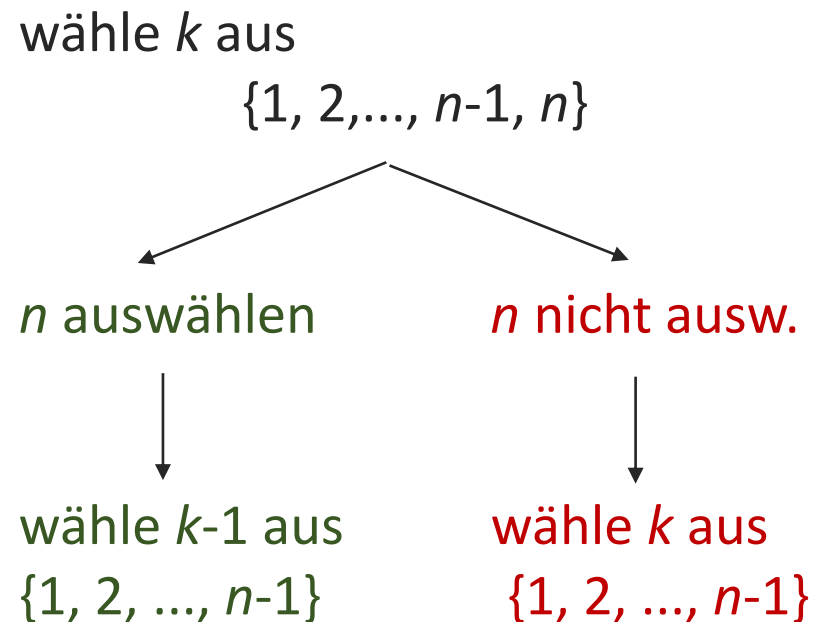
Verallgemeinerung

- Lösung des Problems durch Lösung von (kleineren) Teilproblemen
 - Lösung aller Teilprobleme *nur einmal*, Speichern der Ergebnisse in einer Tabelle
 - werden Lösungen von Teilproblemen zur Lösung eines größeren Problems benötigt:
Auslesen des Tabelleneintrags
- **Dynamische Programmierung**
- kann oft statt Rekursion verwendet werden

DP *versus* Rekursion

- *Beispiel:* Berechne Anzahl der Teilmengen mit k Elementen einer Menge mit n Elementen, also berechne $\binom{n}{k}$

- $$\binom{n}{k} = \binom{n-1}{k-1} + \binom{n-1}{k}$$



DP *versus* Rekursion

- *Beispiel:* Berechne Anzahl der Teilmengen mit k Elementen einer Menge mit n Elementen, also berechne $\binom{n}{k}$

- $$\binom{n}{k} = \binom{n-1}{k-1} + \binom{n-1}{k}$$

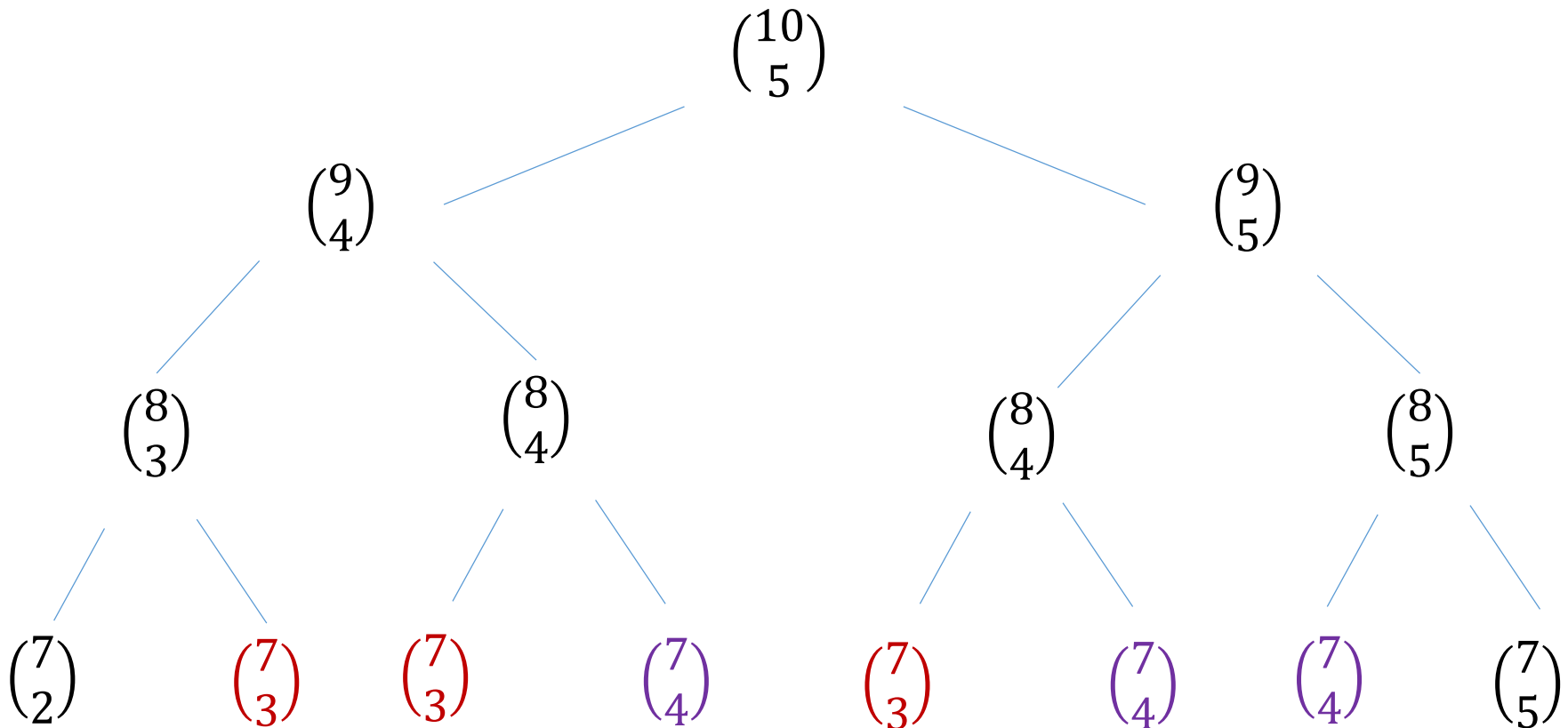
- Basisfälle:

$$\binom{n}{0} = 1 \ (n \geq 0), \quad \binom{0}{k} = 0 \ (k > 0)$$

Teilmengenproblem mit Rekursion

```
def compNChooseK(n,k):  
    if n>=0 and k==0:  
        return 1  
    elif n==0 and k>0:  
        return 0  
    else:  
        return compNChooseK(n-1,k-1)  
            + compNChooseK(n-1,k)
```

Überlappende Aufrufe bei Rekursion



...

Teilmengenproblem mit DP

Tabelle **Chs**

mit

$$\text{Chs}(n, k) = \binom{n}{k}$$

für $0 \leq k \leq n$

und

$$\text{Chs}(n, k) = 0$$

für $k > n$

n

3	1	3	3	1	0	0
2	1	2	1	0	0	0
1	1	1	0	0	0	0
0	1	0	0	0	0	0

0

1

2

3

4

5

k

DP *versus* Rekursion

- **Rekursion:** *Top-Down-Ansatz:*

Berechnung von $f(n)$ durch

- Bestimmung von m_1, m_2, \dots, m_k kleiner als n so, dass
- $f(n)$ einfach aus $f(m_1), f(m_2), \dots, f(m_k)$ berechnet werden kann

→ z.T. müssen nicht alle Teilprobleme gelöst werden,
manche dafür aber mehrfach

- **DP:** *Bottom-Up-Ansatz:*

- Berechnung der einfachsten Teilprobleme zuerst
- Errechnung der nächst größeren aus den bereits errechneten Lösungen, bis $f(n)$ berechnet ist

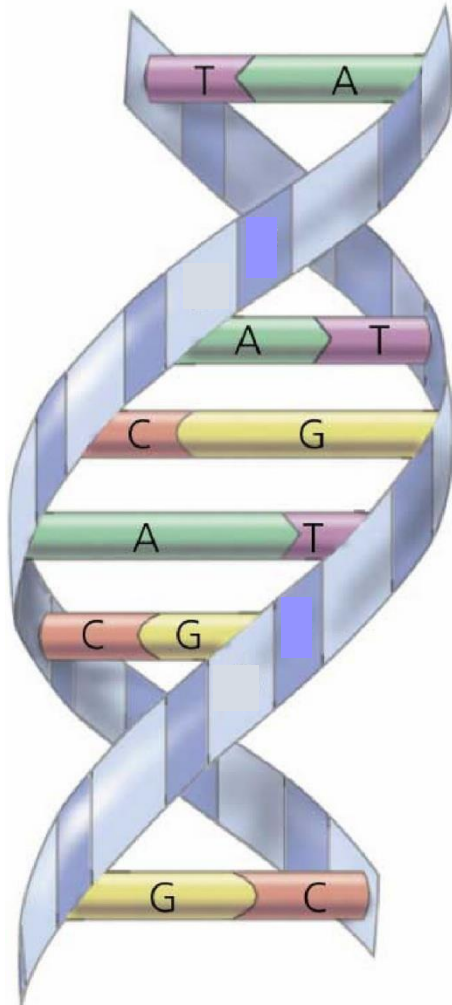
→ Berechnung aller Teilprobleme, aber jedes nur einmal

DP *versus* Rekursion (2)

- **Vorteil** von DP gegenüber Rekursion, falls viele **überlappende Teilprobleme** (*also falls viele Teilprobleme die rekursive Funktion mit denselben Parametern aufrufen*):
Keine wiederholten Aufrufe der rekursiven Funktion mit denselben Parametern bei DP.
- **Nachteil** von DP: Speicherplatz für Tabellengröße muss die gesamte Laufzeit im Speicher gehalten werden.

Das RNA-Sekundärstruktur-Problem

DNA



gewundene Doppelhelix,
bestehend aus:

- Adenin- (A),
- Cytosin- (C),
- Guanin- (G),
- Thymin- (T)

Nukleotiden

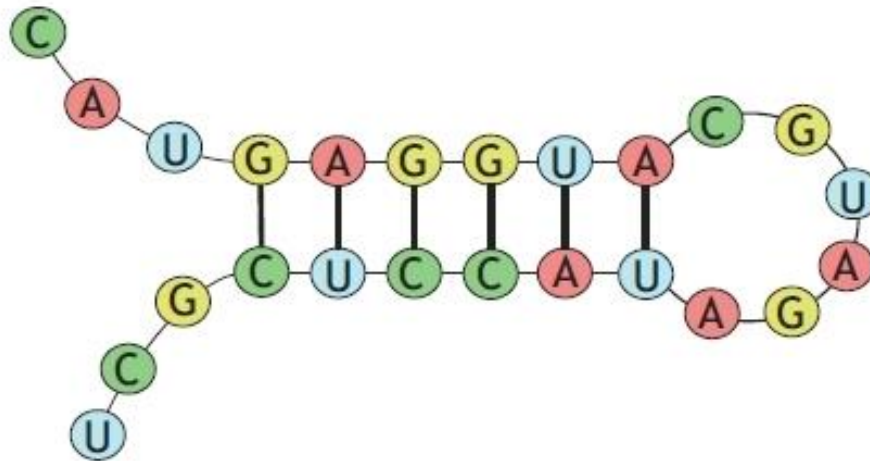
**Watson-Crick-
Komplementarität:**

nur A und T sowie C und G
können Bindungen eingehen

RNA

- Einfachstrang statt Doppelstrang
- Uracil (**U**) statt Thymin (**T**)
- verantwortlich für Proteinbildung
 - „stückweises Auftrennen“ des Doppelstrangs und
 - „Ablesen“ der DNA-Information mit Hilfe der Watson-Crick-Komplementarität
- besitzt **Sekundärstruktur** (Faltung) durch Bindung von komplementären Nukleotiden

Sekundärstruktur – Beispiel



Die stabilste Sekundärstruktur (mit minimaler freier Energie) wird angenommen, wenn eine maximale Anzahl von Nukleotid-Paaren Bindungen eingehen.

RNA-Sekundärstruktur-Problem

- Gegeben ein RNA-Strang (ohne Bindungen).
- Finde maximale Anzahl von Paaren von Nukleotiden, die eine Bindung eingehen können.
- **Eingabe:** RNA-Sequenz
- **Ausgabe:** maximale Anzahl der bindungsfähigen Nukleotiden-Paare in RNA-Sequenz

Algorithmisches Denken: Vom Problem zur Lösung

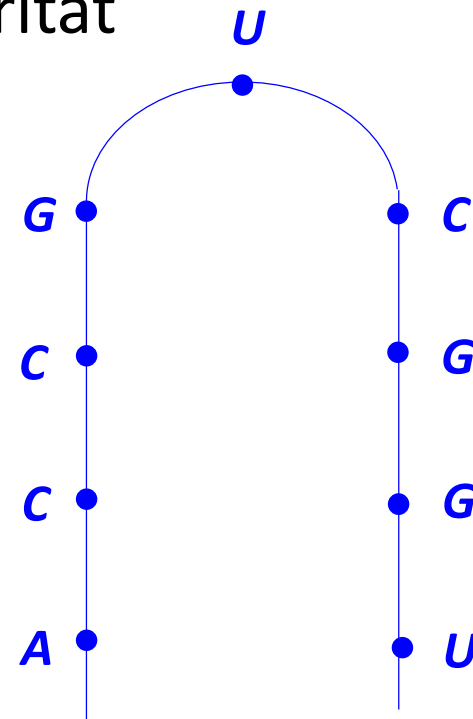
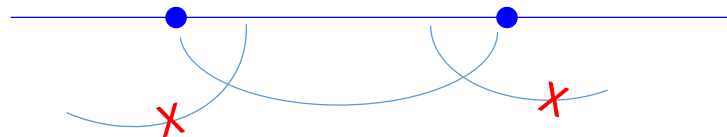
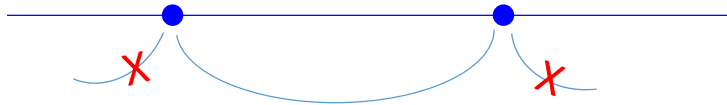
1. Identifizieren des Problems
2. Formulieren des Problems
3. Entwurf des Algorithmus
4. Implementierung des Algorithmus
5. Anwendung des Algorithmus
→ Problemlösung



*Vom Problem
zum Algorithmus*

Identifizieren des Problems: Kriterien Bindungsfähigkeit

1. Watson-Crick-Komplementarität
2. weiche Wendungen
3. Eindeutige Bindungen
4. Überlappungsfreiheit



Formulieren des Problems

- Gegeben RNA-Sequenz w der Länge n .
- **Modell** für einen RNA-Strang:
String $w \in \{A, C, G, U\}^*$
- $\text{nuc}(i) :=$ Nukleotid an Position i ($1 \leq i \leq n$)
- $\text{wcc}(i, j)$ ($1 \leq i < j \leq n$) ist **True** gdw.
Nukleotide an Positionen i und j sind
Watson-Crick-komplementär, d.h.

$$\text{nuc}(i) = A \longrightarrow \text{nuc}(j) = U$$

$$\text{nuc}(i) = C \longrightarrow \text{nuc}(j) = G$$

$$\text{nuc}(i) = G \longrightarrow \text{nuc}(j) = C$$

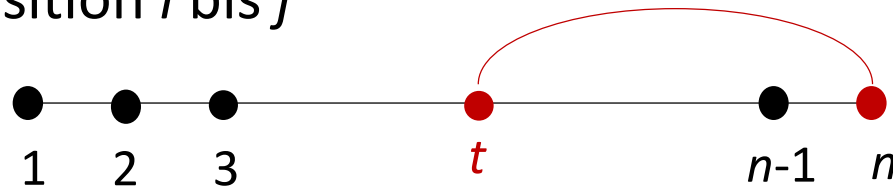
$$\text{nuc}(i) = U \longrightarrow \text{nuc}(j) = A$$

Kriterien Bindungsfähigkeit im Modell

- Gegeben RNA-Sequenz w der Länge n .
 $\text{nuc}(i) :=$ Nukleotid an Position i ($1 \leq i \leq n$)
- Ein Paar (i, j) , $1 \leq i < j \leq n$, ist **Matching Pair** wenn es folgende Kriterien erfüllt:
 1. Watson-Crick-Komplementarität: $\text{wcc}(i, j)$
 2. weiche Wendungen: $j > i + 4$
 3. Eindeutige Bindungen:
wenn (i, l) oder (k, j) Matching Pair, dann $l = j$ bzw. $k = i$
 4. Überlappungsfreiheit:
wenn (k, l) Matching Pair, dann
 $j < k$ oder $l < i$ oder $i < k < l < j$ oder $k < i < j < l$

Entwurf des Algorithmus

- $\text{OPT}(i,j)$: maximale Anzahl von Matching Pairs von Position i bis j



- Suchen $\text{OPT}(1,n)$
- **1. Fall:** n ist keine Komponente in einem Matching Pair in maximaler Menge
 $\rightarrow \text{OPT}(1,n) = \text{OPT}(1,n-1)$
- **2. Fall:** (t,n) ist Matching Pair in maximaler Menge
 $\rightarrow \text{OPT}(1,n) = 1 + \text{OPT}(1,t-1) + \text{OPT}(t+1,n-1)$
- $\text{OPT}(1,n)$ ist das Maximum der Werte aus den beiden Fällen.

Entscheidung für das Paradigma

Dynamische Programmierung

- Die Option *rekursive Aufrufe* von $\text{OPT}(i,j)$ wird zu häufigen Aufrufen der Prozedur für dieselben Parameter i, j führen
- besser: systematisch für alle Paare (i,j) aufrufen und tabellarisch die Ergebnisse speichern
- Reihenfolge der Aufrufe so, dass für die Berechnung eines neuen Wertes nur bereits berechnete Werte aus der Tabelle ausgelesen werden müssen

Pseudocode

Eingabe: $w \in \{A,C,G,U\}^*$ mit $|w| = n$

Ausgabe: $\text{OPT}(1,n)$

Für alle $1 \leq i \leq n$ und $0 \leq j \leq i+4$

$\text{OPT}(i,j) \leftarrow 0$

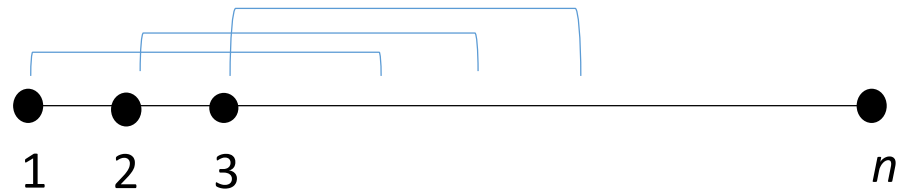
Für $5 \leq k \leq n-1$ # Distanz zwischen i und j

Für $1 \leq i \leq n-k$

$j \leftarrow i+k$

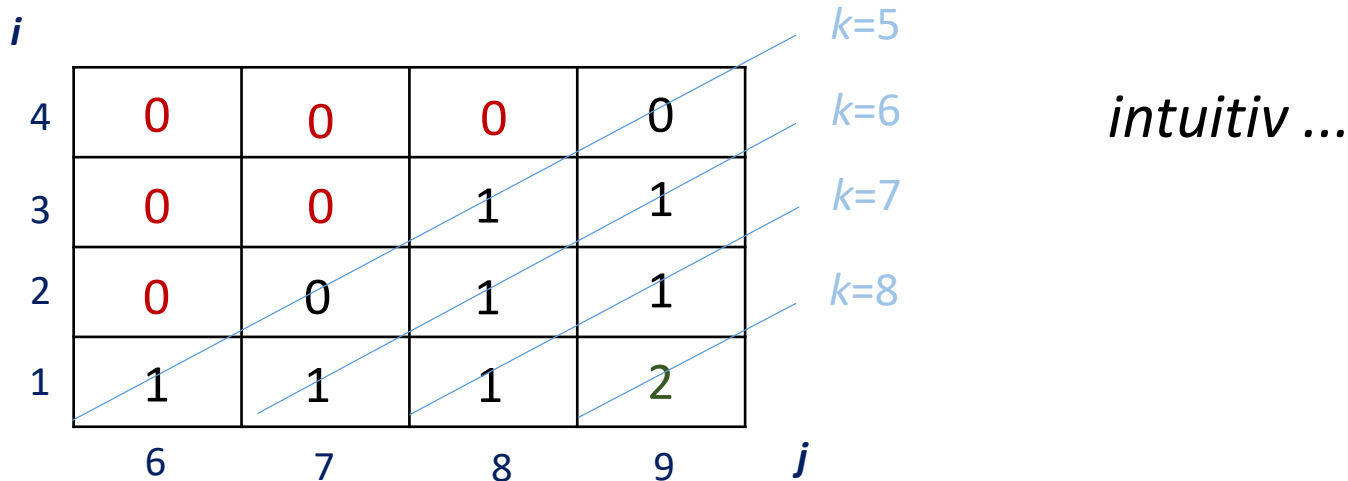
$\text{OPT}(i,j) \leftarrow \max \{ \text{OPT}(i,j-1),$
 $\max_t \{ 1 + \text{OPT}(i,t-1) + \text{OPT}(t+1,j-1) \mid (t,j) \text{ Matching Pair} \} \}$

Gib $\text{OPT}(1,n)$ aus



Beispiel: A C C G G U A G U

1 2 3 4 5 6 7 8 9 (n=9)



Für $5 \leq k \leq n-1$ # Distanz zwischen i und j

Für $1 \leq i \leq n-k$

$j \leftarrow i + k$

$$\text{OPT}(i,j) \leftarrow \max \{ \text{OPT}(i,j-1),$$

$$\max_t \{ 1 + \text{OPT}(i,t-1) + \text{OPT}(t+1,j-1) \mid (t,j) \text{ Matching Pair} \}$$

Beispiel: A C C G G U A G U

1 2 3 4 5 6 7 8 9 (n=9)

i

4	0	0	0	0
3	0	0	1	1
2	0	0	1	1
1	1	1	1	2

6 7 8 9 *j*

k=5
k=6
k=7
k=8

$$\begin{aligned} \text{OPT}(1,6) &= \max \{ \text{OPT}(1,5), \\ &\quad \max\{1+\text{OPT}(1,0)+\text{OPT}(2,5)\} \} \\ &\quad (t=1) \\ &= \max \{ 0, 1+0+0 \} \\ &= 1 \end{aligned}$$

Für $5 \leq k \leq n-1$ # Distanz zwischen i und j

Für $1 \leq i \leq n-k$

$j \leftarrow i + k$

$$\begin{aligned} \text{OPT}(i,j) &\leftarrow \max \{ \text{OPT}(i,j-1), \\ &\quad \max_t \{ 1 + \text{OPT}(i,t-1) + \text{OPT}(t+1,j-1) \mid (t,j) \text{ Matching Pair} \} \} \end{aligned}$$

Beispiel: A C C G G U A G U

1 2 3 4 5 6 7 8 9 (n=9)

i

				<i>k=5</i>	
4	0	0	0	0	<i>k=6</i>
3	0	0	1	1	<i>k=7</i>
2	0	0	1	1	<i>k=8</i>
1	1	1	1	2	
	6	7	8	9	<i>j</i>

$$\begin{aligned}
 \text{OPT}(2,8) &= \max \{ \text{OPT}(2,7), \\
 &\quad \max\{1+\text{OPT}(2,1)+\text{OPT}(3,7), \\
 &\quad (t=2,3) \\
 &\quad 1+\text{OPT}(2,2)+\text{OPT}(4,7)\} \} \\
 &= \max \{ 0, 1+0+0, \\
 &\quad 1+0+0 \} \\
 &= 1
 \end{aligned}$$

Für $5 \leq k \leq n-1$ # Distanz zwischen i und j

Für $1 \leq i \leq n-k$

$j \leftarrow i + k$

$\text{OPT}(i,j) \leftarrow \max \{ \text{OPT}(i,j-1),$
 $\max_t \{ 1 + \text{OPT}(i,t-1) + \text{OPT}(t+1,j-1) \mid (t,j) \text{ Matching Pair} \} \}$

Beispiel: A C C G G U A G U

1 2 3 4 5 6 7 8 9 (n=9)

i

				<i>k=5</i>
4	0	0	0	<i>k=6</i>
3	0	0	1	<i>k=7</i>
2	0	0	1	<i>k=8</i>
1	1	1	1	
	6	7	8	<i>j</i>

The table shows the dynamic programming table for the sequence A C C G G U A G U. The rows are indexed by *i* (1 to 4) and the columns by *j* (6 to 9). The values in the cells are: (1,6)=1, (1,7)=1, (1,8)=1, (1,9)=2 (circled in green); (2,6)=0, (2,7)=0, (2,8)=1, (2,9)=1; (3,6)=0, (3,7)=0, (3,8)=1, (3,9)=1; (4,6)=0, (4,7)=0, (4,8)=0, (4,9)=0. Blue diagonal lines are drawn from (1,6) to (4,9), (2,7) to (4,9), (3,8) to (4,9), and (4,9) to (4,9). Labels *k=5*, *k=6*, *k=7*, and *k=8* are placed to the right of the table, corresponding to the diagonal lines.

$$\begin{aligned} \text{OPT}(1,9) &= \max \{ \text{OPT}(1,8), \\ &\quad \max\{1 + \text{OPT}(1,0) + \text{OPT}(2,8)\} \} \\ &\quad (t=1) \\ &= \max \{1, 1+0+1\} \\ &= 2 \end{aligned}$$

Für $5 \leq k \leq n-1$ # Distanz zwischen *i* und *j*

Für $1 \leq i \leq n-k$

$j \leftarrow i + k$

$$\begin{aligned} \text{OPT}(i,j) &\leftarrow \max \{ \text{OPT}(i,j-1), \\ &\quad \max_t \{ 1 + \text{OPT}(i,t-1) + \text{OPT}(t+1,j-1) \mid (t,j) \text{ Matching Pair} \} \} \end{aligned}$$

Zeitkomplexität

Eingabe: $w \in \{A,C,G,U\}^*$ mit $|w| = n$

Ausgabe: $\text{OPT}(1,n)$

$O(n^3)$

Für alle $1 \leq i \leq n$ und $0 \leq j \leq i + 4$ $O(n^2)$

$\text{OPT}(i,j) \leftarrow 0$

Für $5 \leq k \leq n-1$ $O(n)$ $O(n^3)$

Für $1 \leq i \leq n-k$ $O(n)$

$j \leftarrow i + k$

$\text{OPT}(i,j) \leftarrow \max \{ \text{OPT}(i,j-1),$
 $\max \{ 1 + \text{OPT}(i,t-1) + \text{OPT}(t+1,j-1) \mid (t,j) \text{ Matching Pair} \} \}$
 $O(n)$

Gib $\text{OPT}(1,n)$ aus

Abschließende Bemerkungen

- Finden der maximalen Menge der Matching Pairs durch „Mitschreiben“ der (t,j) , die den maximalen Wert bei der Berechnung der $OPT(i,j)$ liefern
- Für Algorithmen nach dynamischer Programmierung: Tabellengröße ist untere Schranke der Laufzeit
 - beim RNA-Sekundärstrukturproblem: Tabellengröße ist quadratisch in der Länge der RNA-Sequenz
 - wäre diese z.B. exponentiell in der Größe der Eingabe: mindestens exponentielle Laufzeit