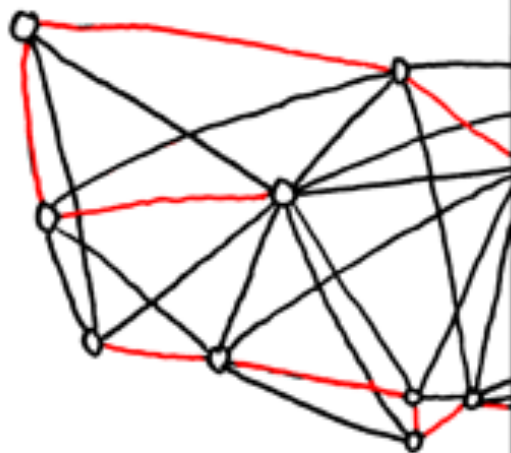
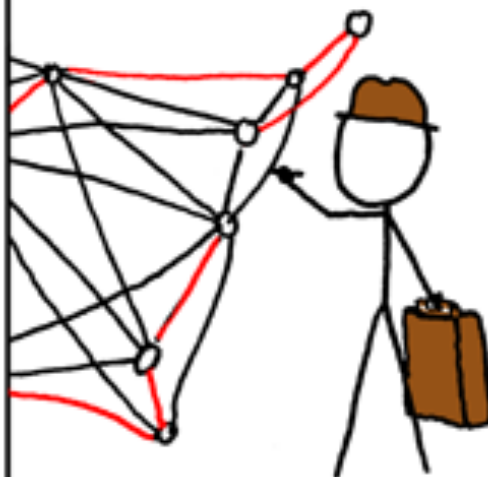


BRUTE-FORCE  
SOLUTION:  
 $O(n!)$



DYNAMIC  
PROGRAMMING  
ALGORITHMS:  
 $O(n^2 2^n)$



SELLING ON EBAY:  
 $O(1)$

STILL WORKING  
ON YOUR ROUTE?

SHUT THE  
HELL UP.



# **Algorithmen in der Bioinformatik**

## **3. Sequenz-Vergleiche - Dynamische Programmierung -**

**Prof. Dr. Gunnar Klau**

Nach Jones & Pevzner: An Introduction to Bioinformatics Algorithms,  
Kapitel 6



# DNA-Sequenz-Vergleiche

Die Suche nach Ähnlichkeit einer Sequenz mit bekannten Genen ist **die** Methode zur Funktionsbestimmung

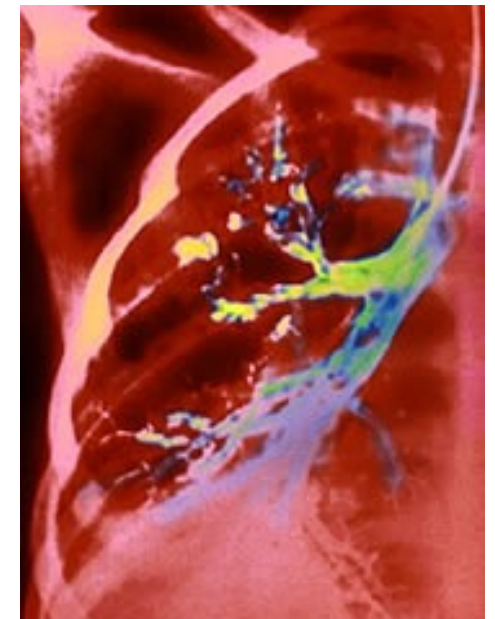
1984 fanden Russell Doolittle & Co. Ähnlichkeit zwischen einem krebsverursachenden Gen und einem normalem Wachstumsfaktor



# Zystische Fibrose/Mukoviszidose

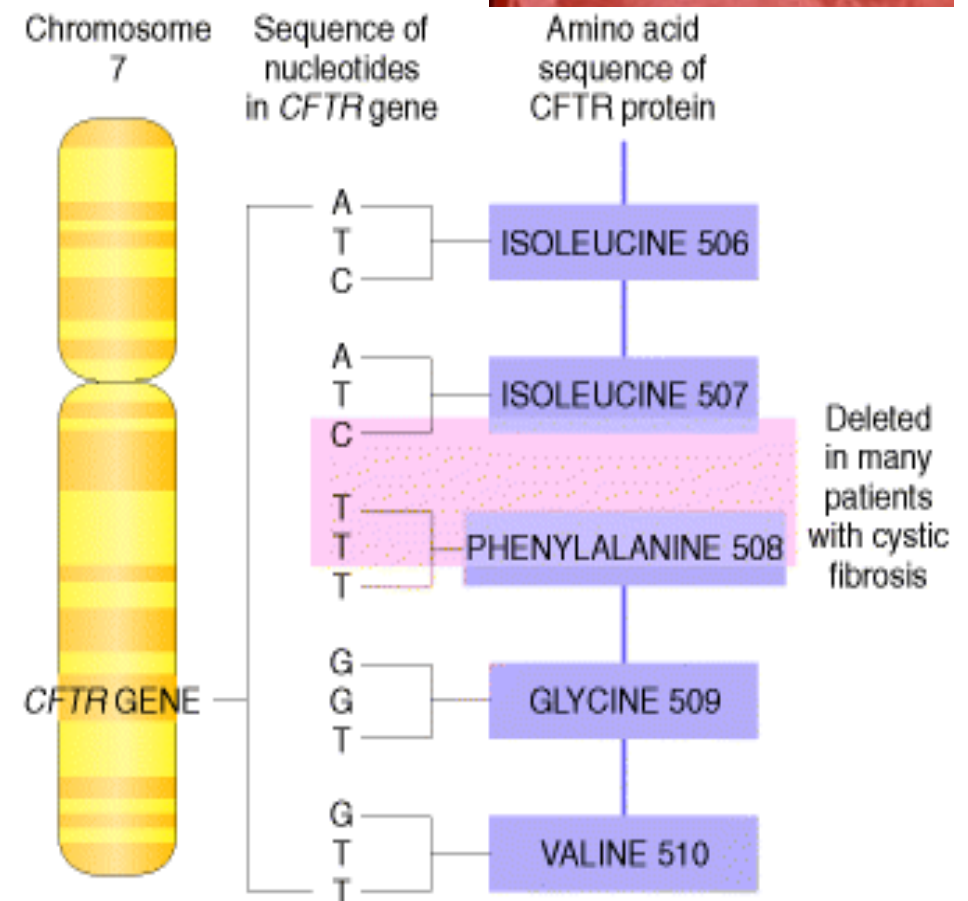
**Zystische Fibrose:** chronische, häufig tödliche Krankheit der Schleimdrüsen

- Führt zu Verstopfung der Bronchien → Atemprobleme
- autosomal rezessiv (++, +-, --)
- Verantwortliches Gen unbekannt bis 1989

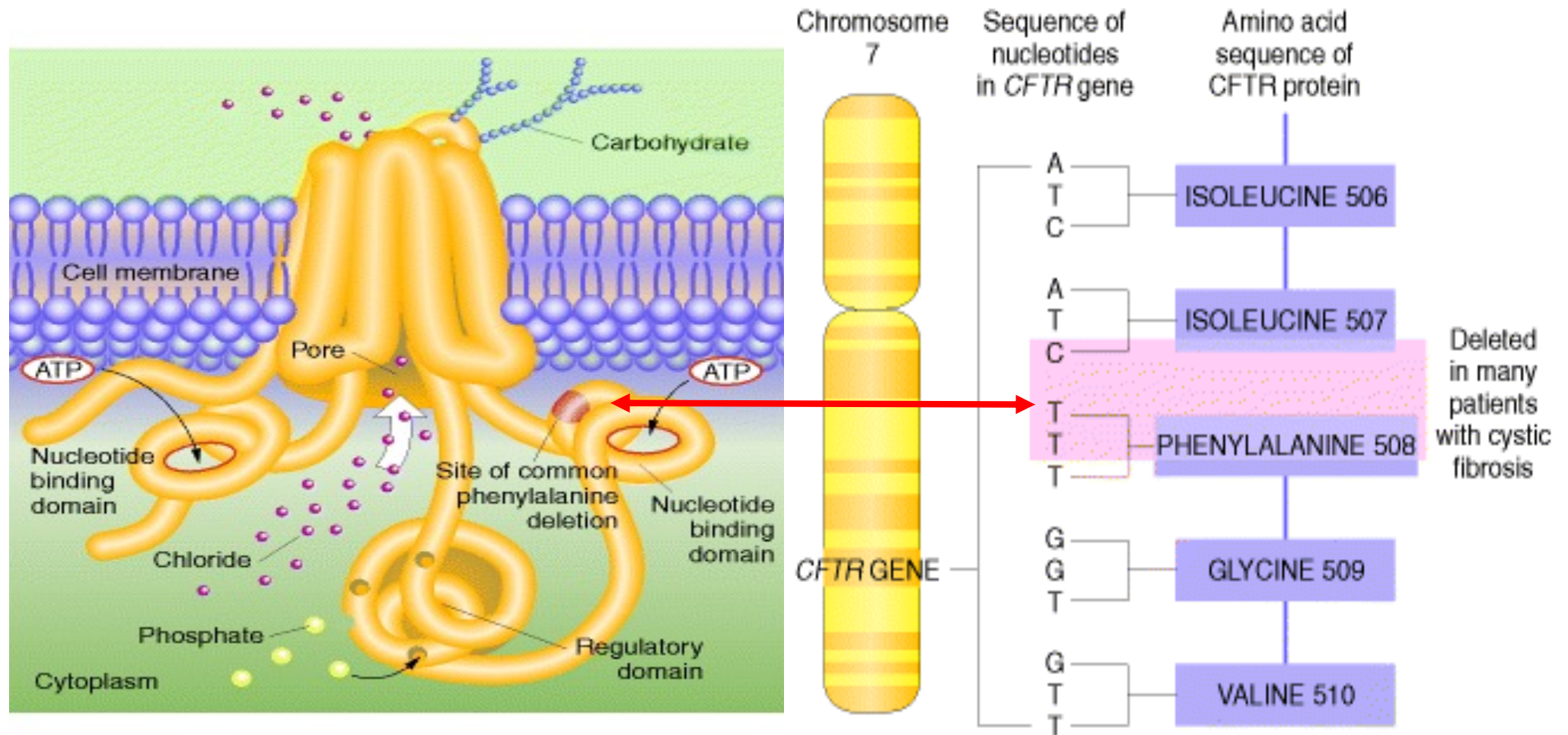


## Bioinformatik:

- 70% der Patienten haben dieselbe Mutation des CFTR-Gens
- Ähnlichkeit mit Genen, die als Transport-Kanäle in der Zellmembran fungieren
- Funktion: Cl<sup>-</sup>-Kanal (justiert Viskosität des Sekrets)



# Zystische Fibrose



# Sequenzvergleich: wozu und wie?

- Ähnlichkeit eines unbekannten Gens zu einem bekannten Gen gibt Hinweise auf Funktion
- Ähnlichkeits-Score erlaubt Aussagen über die Wahrscheinlichkeit gleicher bzw. ähnlicher Funktion
- **Dynamische Programmierung** ist eine Technik u. a. zum Aufdecken von Ähnlichkeiten zwischen Genen



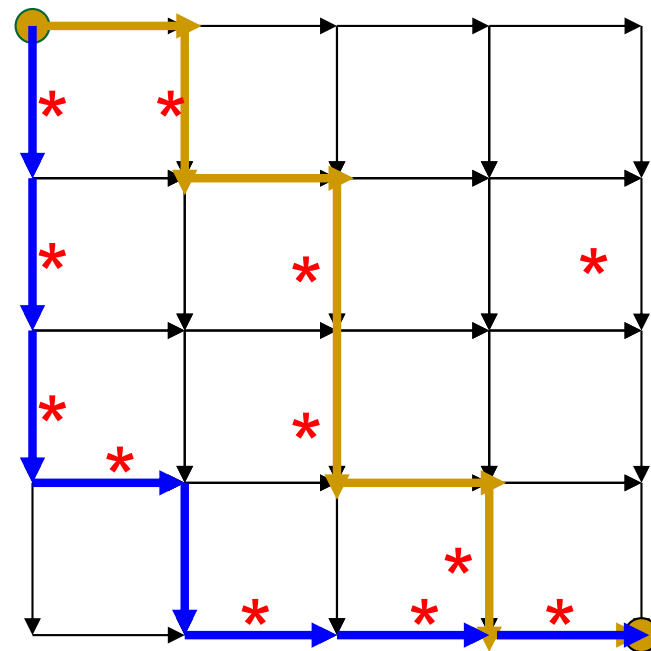


# Tourismus in Manhattan: Problem

Finde den Weg durch Manhattan (nur ost- und südwärts!), der an den meisten Attraktionen (\*) vorbeikommt.

Wenn wir Kanten Gewichte geben (\*=1, sonst 0), dann suchen wir den 'längsten' Pfad.

Quelle



Senke

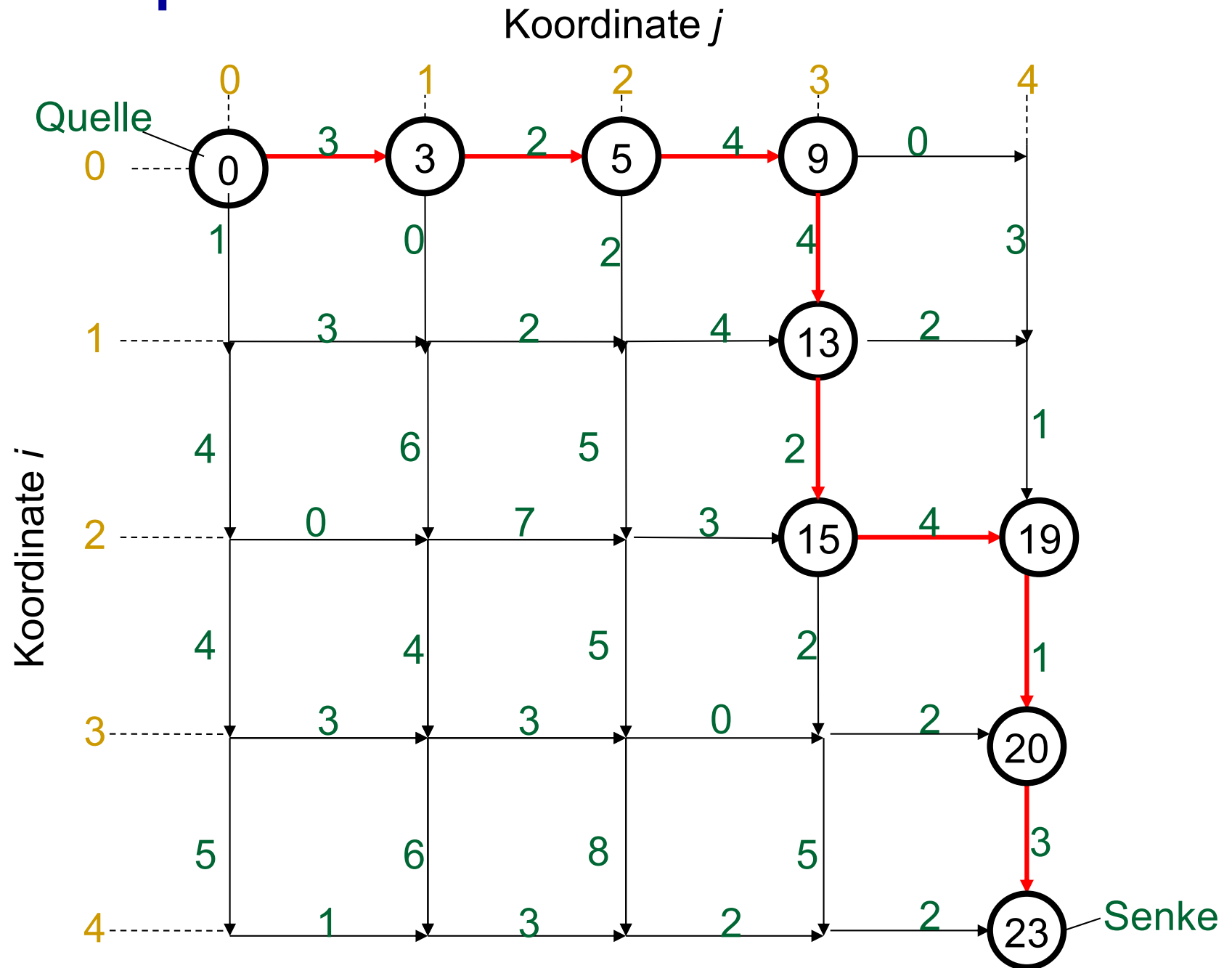
# Tourismus in Manhattan: Formulierung

Ziel: Finde den längsten Pfad in einem gewichteten “Süd-Ost”-Gitter.

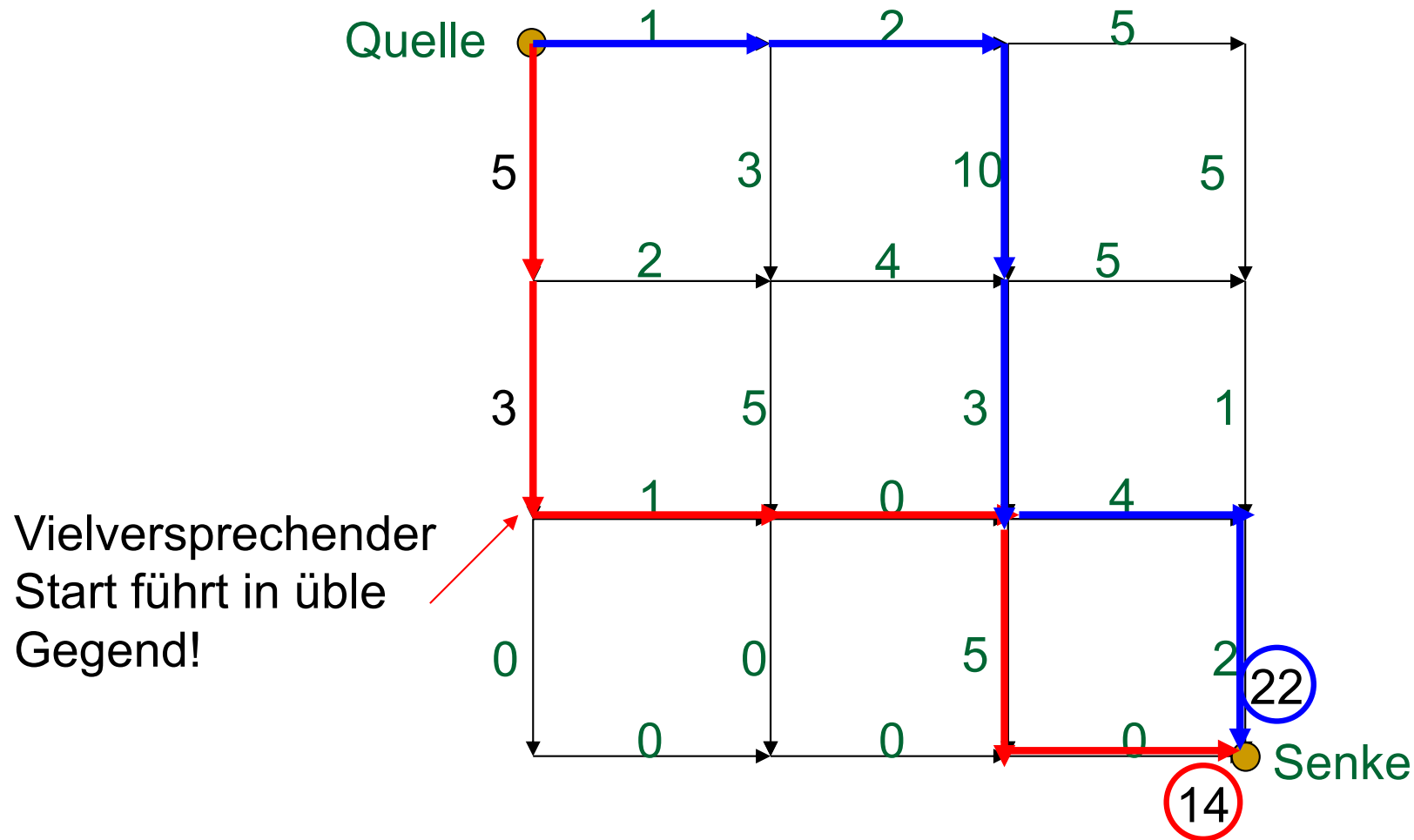
Eingabe: Ein gewichtetes Gitter  $G$  mit zwei speziellen Knoten (Quelle und Senke)

Ausgabe: Ein längster Pfad in  $G$  von Quelle zu Senke

# TIM: Beispiel



# TIM: Gieriger Algorithmus ist nicht exakt





# Warum nur ein Problem gleichzeitig lösen?

- TIM:  
Suche den längsten Pfad von  $(0,0)$  nach  $(m, n)$
- Verallgemeinertes Problem:  
Suche den längsten Pfad von  $(0,0)$  nach  $(i, j)$  für alle  $i = 1, \dots, m$ ,  
 $j = 1, \dots, n$ .

Das allgemeine Problem sieht zwar schwieriger aus, ist aber keine große zusätzliche Arbeit:

**Grundidee der dynamischen Programmierung!**

# Dynamische Programmierung (DP)

## ➤ Idee

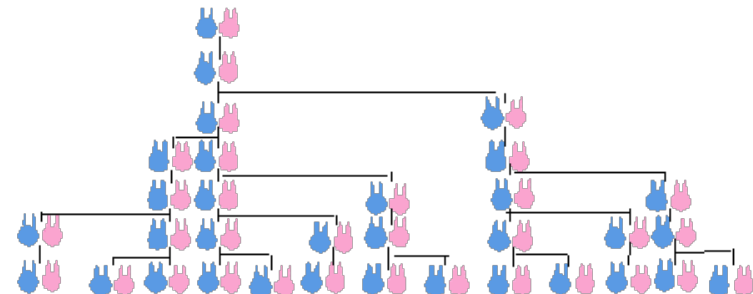
- **Zerteile** Problem in kleinere Teilprobleme
- Baue Lösungen größerer Probleme aus Lösungen kleinerer Probleme zusammen
- **Effiziente** Realisierung eines rekursiven Algorithmus durch **Speicherung** von Teilresultaten
  - Speichere diese in **DP-Tabelle**
  - Kaufe dir Rechenzeit mit Speicherplatz



Richard E. Bellman

## ➤ Wichtige Überlegungen

- Wie teile ich auf? Definiere Wert der Optimallösung rekursiv
- Bestimme diesen Wert **“bottom-up”**



DP: Fibonacci numbers

$$F_n = F_{n-1} + F_{n-2}, \quad F_0 = 0, \quad F_1 = 1$$

fib(n)

if  $n = 0$  : return 0;

if  $n = 1$  : return 1;

else return fib(n-1) + fib(n-2)

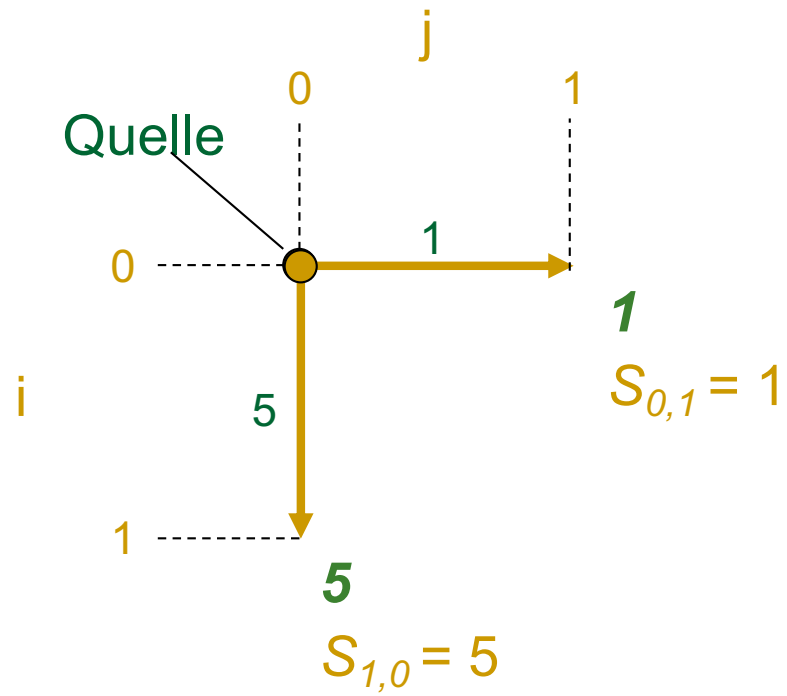
fib\_DP(n)

f[0] = 0; f[1] = 1;

for  $i = 2, \dots, n$  :  $f[i] = f[i-1] + f[i-2]$

return f[n];

# TIM: Dynamische Programmierung





# TIM: Dynamische Programmierung

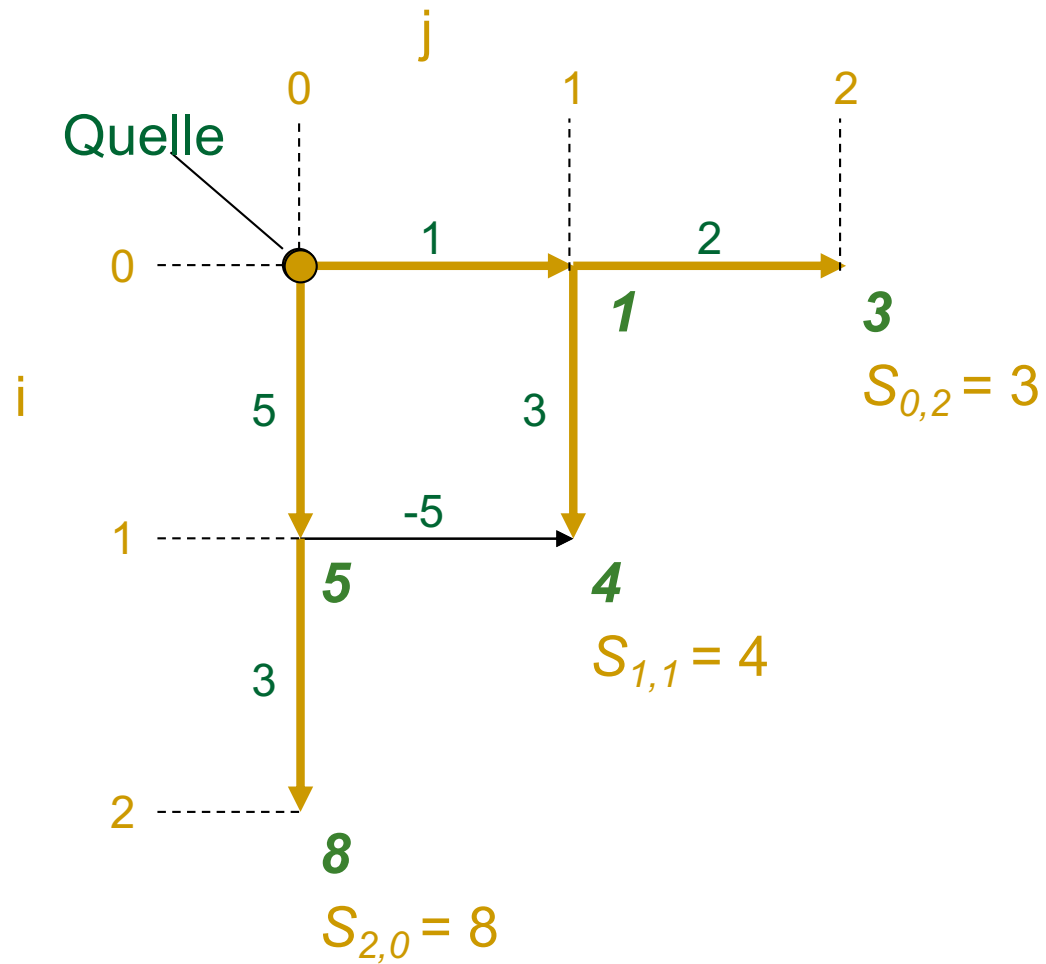


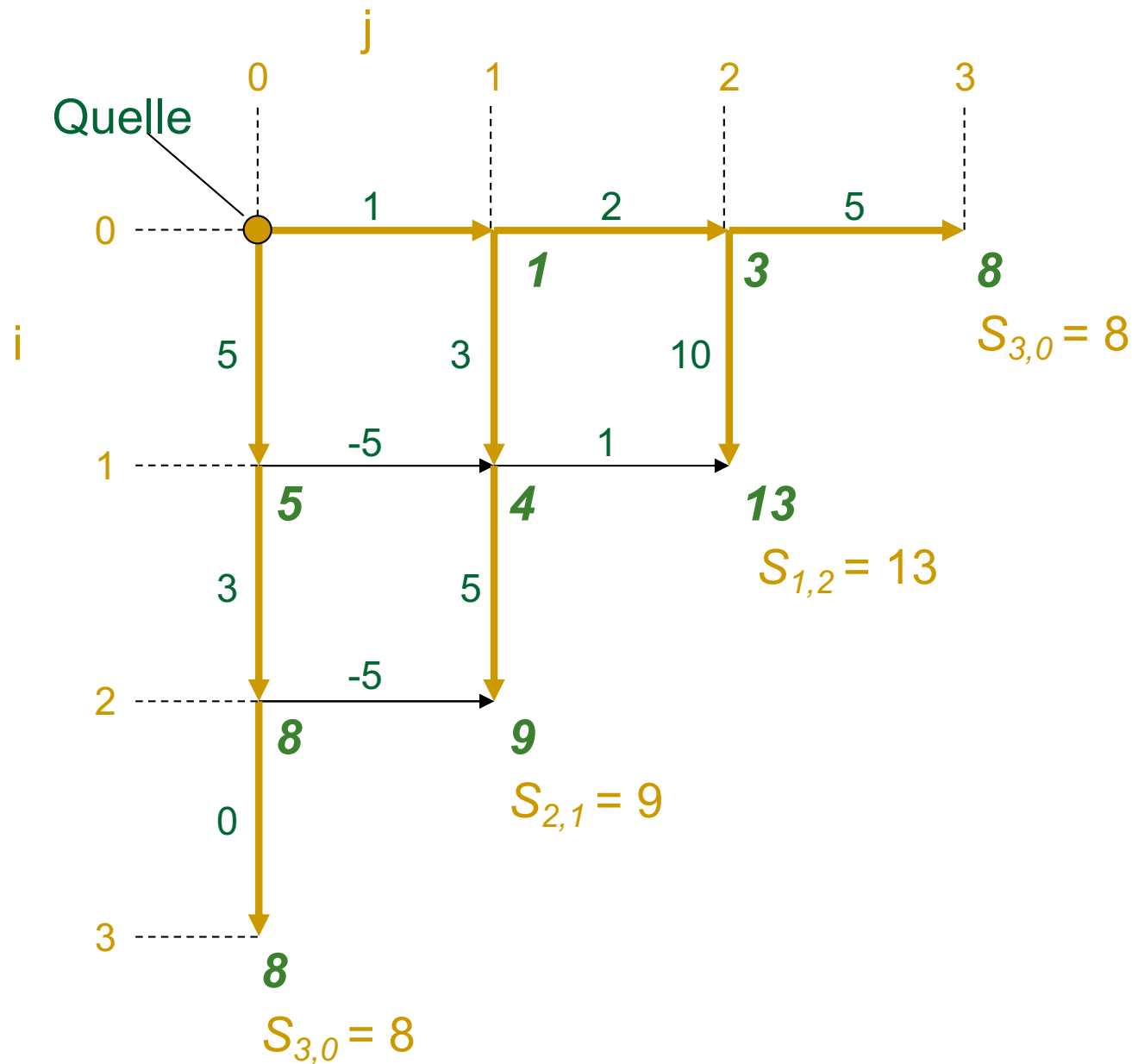
Diagram illustrating a 2D grid with nodes  $(i, j)$  and edges. The source node is labeled "Quelle" at  $(0, 0)$ .

The grid shows nodes  $(i, j)$  and the accumulated value  $S_{i,j}$  at each node:

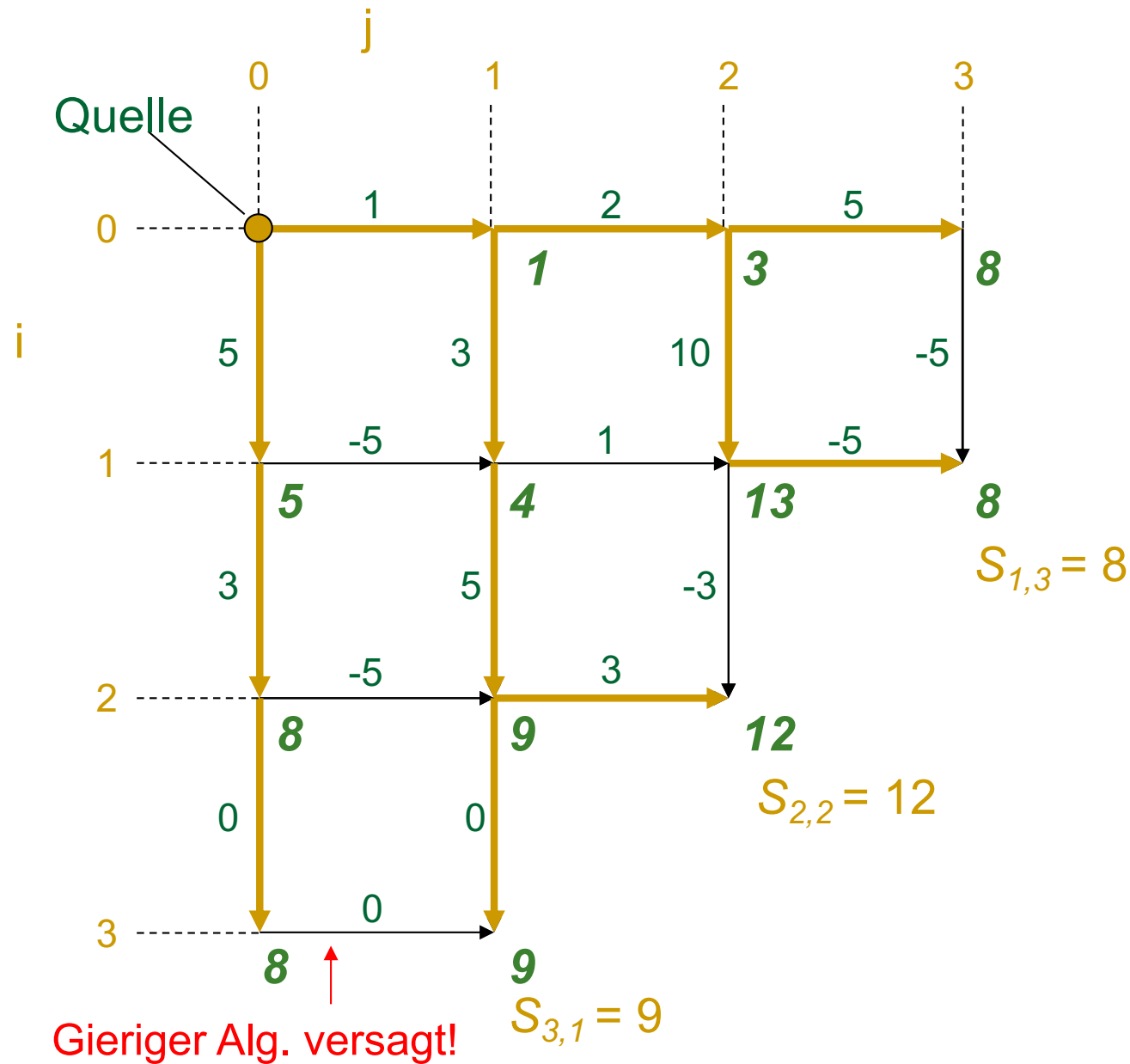
- $(0, 0)$ :  $S_{0,0} = 0$  (Quelle)
- $(0, 1)$ :  $S_{0,1} = 1$
- $(0, 2)$ :  $S_{0,2} = 3$
- $(0, 3)$ :  $S_{0,3} = 8$
- $(1, 0)$ :  $S_{1,0} = 5$
- $(1, 1)$ :  $S_{1,1} = 4$
- $(1, 2)$ :  $S_{1,2} = 9$
- $(1, 3)$ :  $S_{1,3} = 8$
- $(2, 0)$ :  $S_{2,0} = 10$
- $(2, 1)$ :  $S_{2,1} = 13$
- $(2, 2)$ :  $S_{2,2} = 9$
- $(2, 3)$ :  $S_{2,3} = 8$
- $(3, 0)$ :  $S_{3,0} = 8$
- $(3, 1)$ :  $S_{3,1} = 13$
- $(3, 2)$ :  $S_{3,2} = 9$
- $(3, 3)$ :  $S_{3,3} = 8$

Edges and their weights:

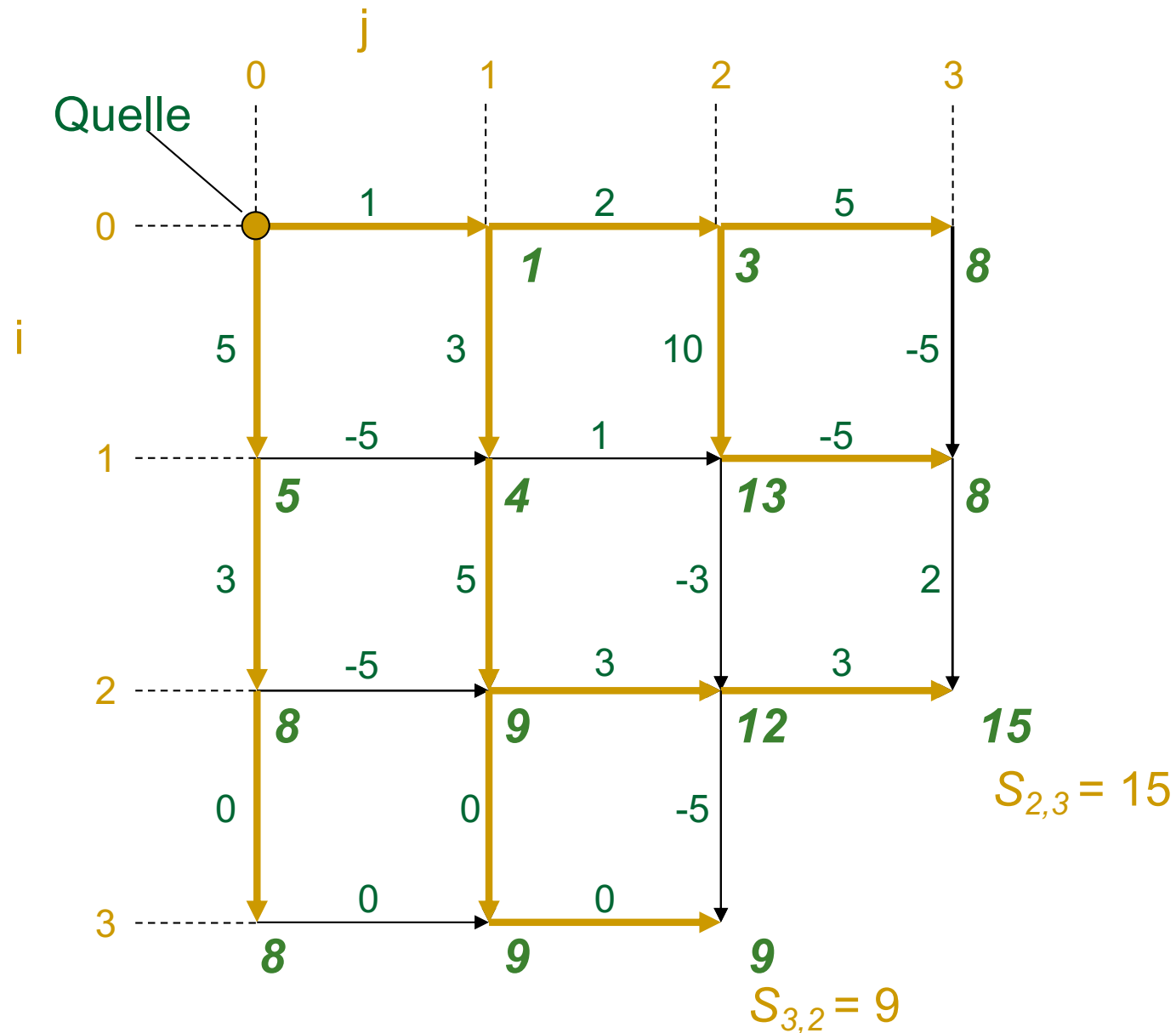
- Horizontal edges:  $(0,0) \rightarrow (1,0)$  (1),  $(1,0) \rightarrow (2,0)$  (2),  $(2,0) \rightarrow (3,0)$  (5),  $(0,1) \rightarrow (1,1)$  (-5),  $(1,1) \rightarrow (2,1)$  (1),  $(0,2) \rightarrow (1,2)$  (-5).
- Vertical edges:  $(0,0) \rightarrow (0,1)$  (5),  $(0,1) \rightarrow (0,2)$  (3),  $(0,2) \rightarrow (0,3)$  (0),  $(1,0) \rightarrow (1,1)$  (3),  $(1,1) \rightarrow (1,2)$  (5),  $(2,0) \rightarrow (2,1)$  (10),  $(2,1) \rightarrow (2,2)$  (3).



# TIM: Dynamische Programmierung

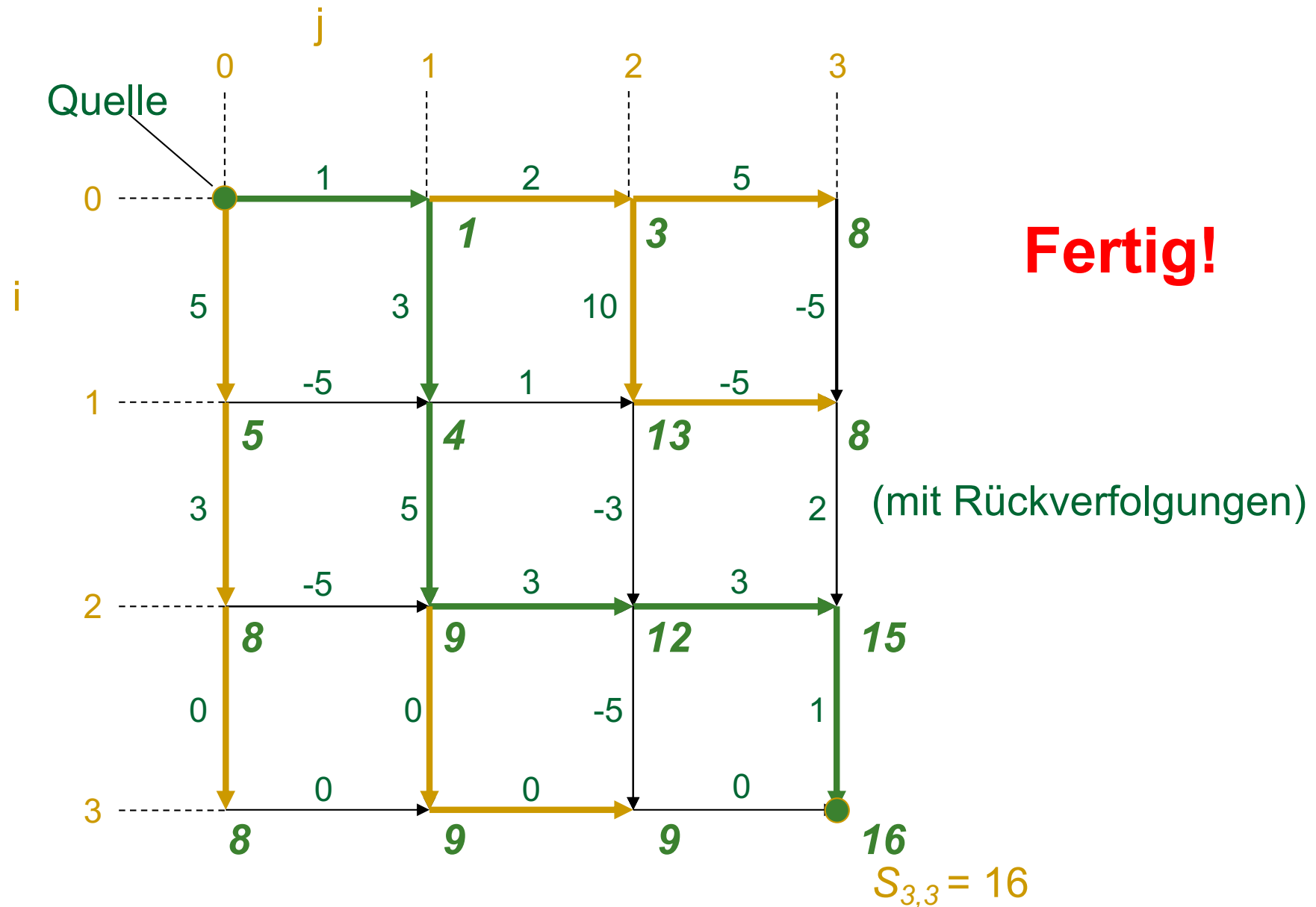


# TIM: Dynamische Programmierung





# TIM: Dynamische Programmierung



# TIM: Rekurrenz

Berechne  $\text{Score}(i,j)$  durch die Rekurrenz-Relation:

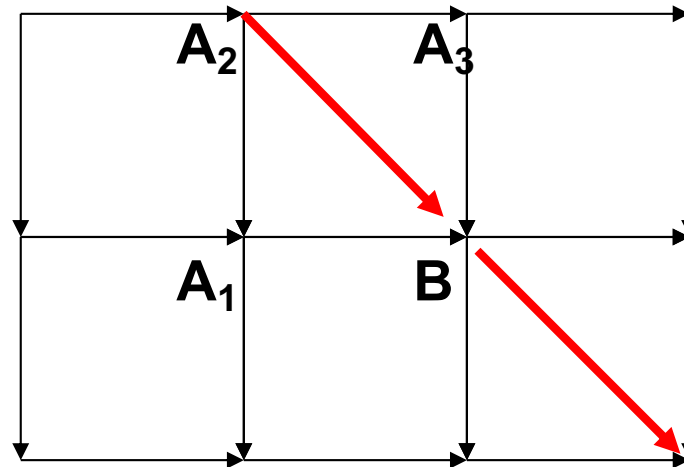
$$s_{i,j} =$$

Laufzeit:  $n \times m$  ( $n = \#$  Spalten,  $m = \#$  Zeilen)

# Gerichtete Azyklische Graphen



# Manhattan ist kein perfektes Gitter



Score( $B$ ) :

$$s_B = \max \left\{ \begin{array}{l} s_{A_1} + \text{Gewicht der Kante } (A_1, B) \\ s_{A_2} + \text{Gewicht der Kante } (A_2, B) \\ s_{A_3} + \text{Gewicht der Kante } (A_3, B) \end{array} \right.$$

Allgemeine Rekurrenz-Relation:

$$s_x = \max_{Y \in \text{Vorgänger}(X)} \{s_Y + \text{Gewicht der Kante } (Y, X)\}$$

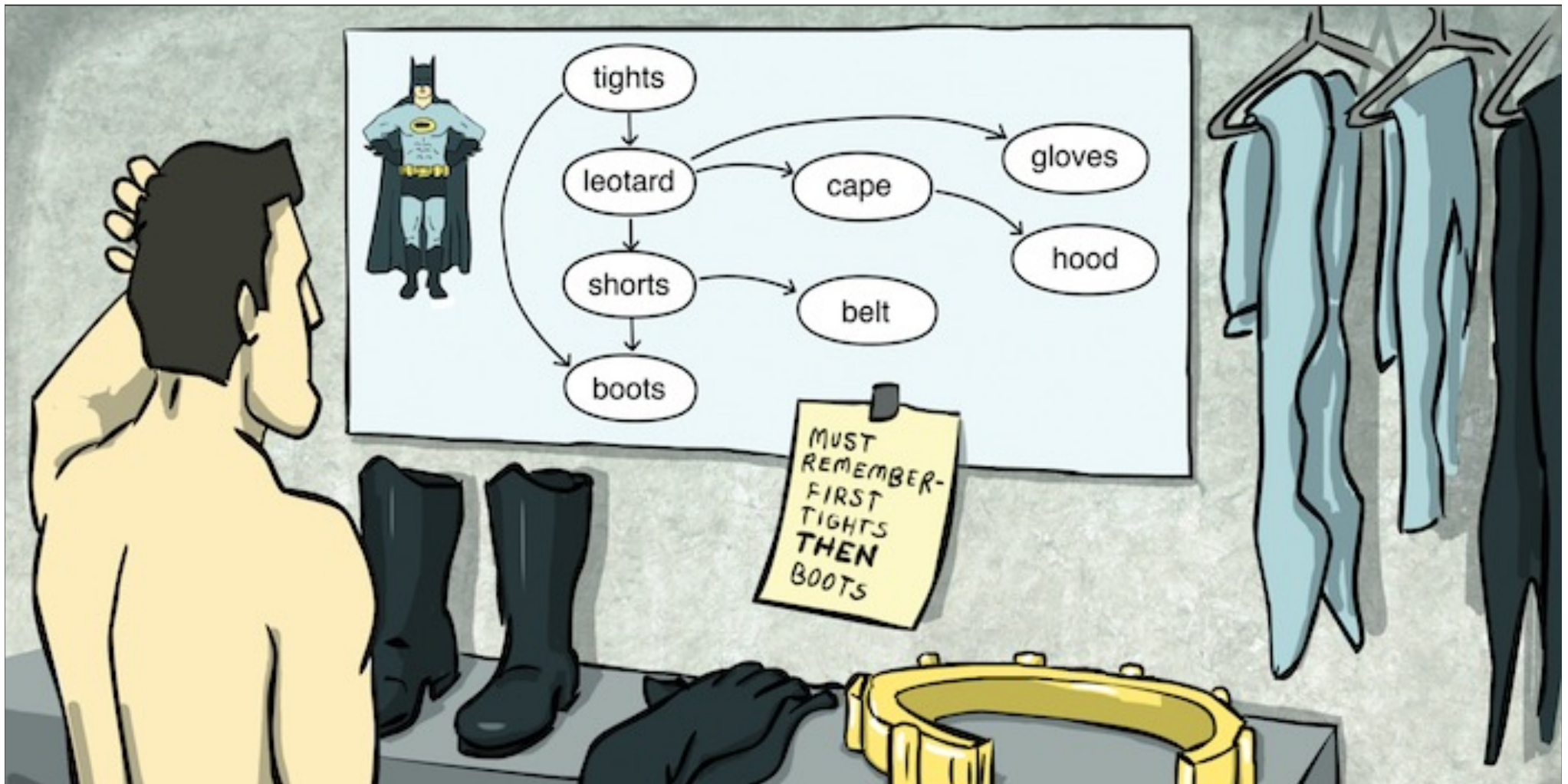
$\text{Vorgänger}(X) := \{\text{Knoten } Y \text{ mit Kante } Y \rightarrow X\}$



# DAG: Directed Acyclic Graphs

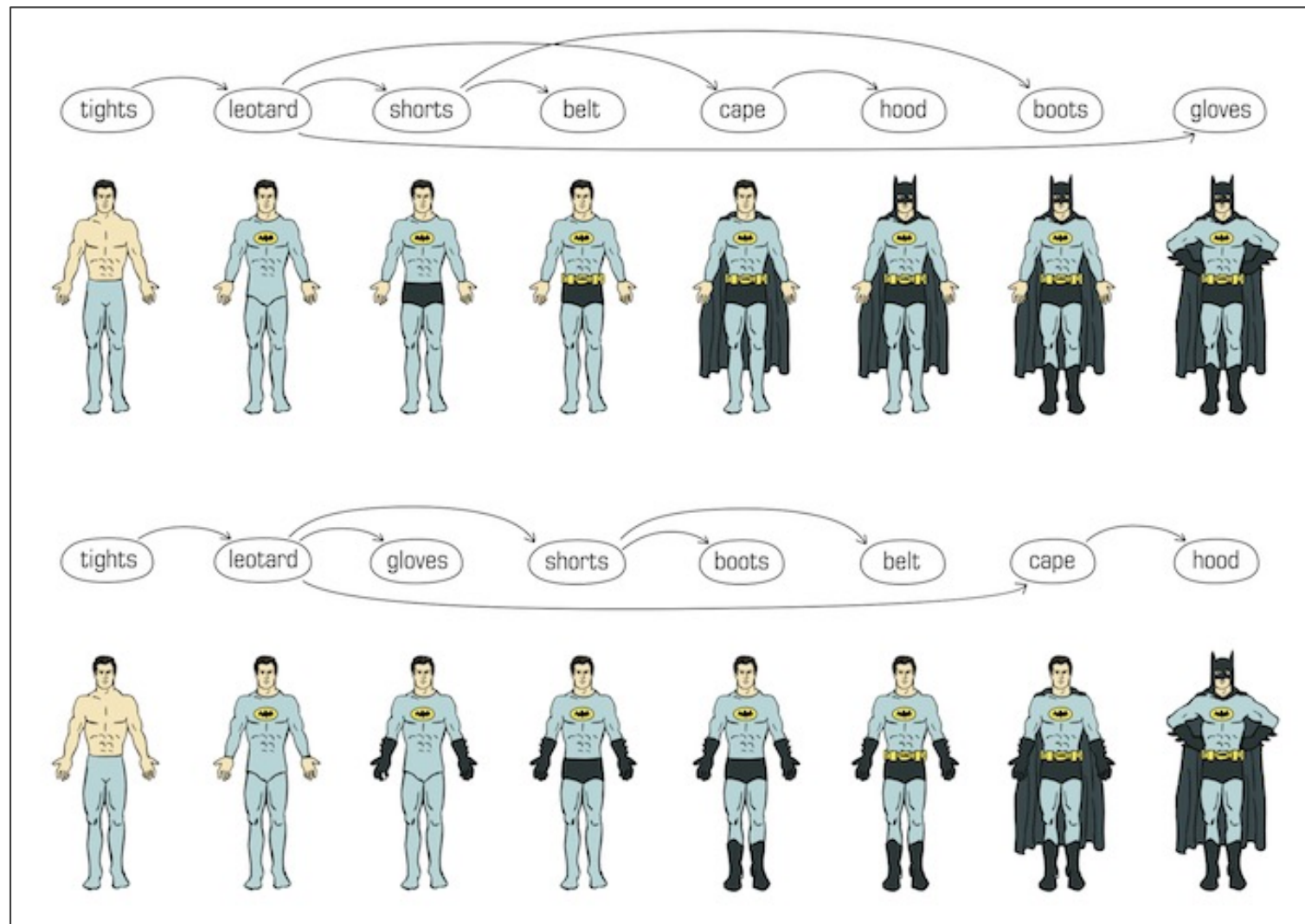
:= Graphen mit gerichteten Kanten ohne Zyklen  
(d.h.: Kein Weg von einem Knoten zu sich selbst)

Wir brauchen eine Ordnung: Wenn wir Knoten X analysieren, müssen Werte für alle Vorgänger Y bereits vorliegen.



# Topologische Ordnungen

Eine Nummerierung von Knoten heißt **topologische Ordnung** des DAG  
: $\Leftrightarrow$  jede Kante läuft von einem Knoten mit kleinerer Nummer  
zu einem Knoten mit größerer Nummer



# Der längste Pfad in einem DAG: Problem

Ziel: Finde den längsten (d.h. am stärksten gewichteten) Pfad zwischen zwei Knoten in einem gewichteten DAG

Eingabe: Ein gewichteter DAG  $G$  mit zwei ausgezeichneten Knoten (Quelle und Senke)

Ausgabe: Ein längster Pfad in  $G$  von Quelle zu Senke

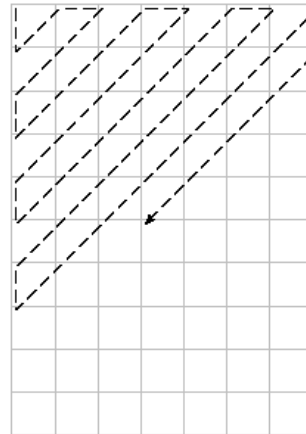
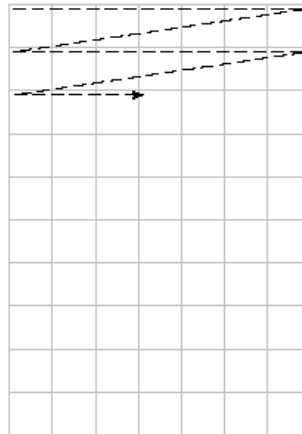
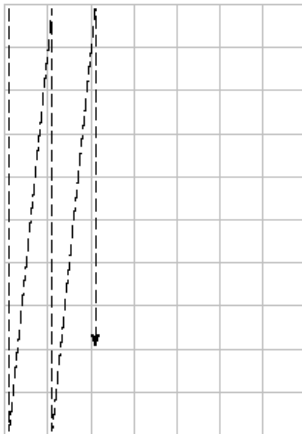
# Der längste Pfad: Dynamische Programmierung

Rekurrenz-Relation:

$$s_x = \max_{Y \in \text{Vorgänger}(X)} \{s_Y + \text{Gewicht der Kante } (Y, X)\}$$

Laufzeit ist  $O(\# \text{ Kanten})$

Verschiedene Strategien:



# Bsp: Münzenkette

Spielregeln:

- Auf dem Tisch liegt eine Kette von Münzen mit verschiedenem Wert.
- Ich kann entweder die erste oder die letzte Münze nehmen und behalten.
- Dann ist mein Gegenspieler dran...
- ...bis alle Münzen weg sind.

Ziel: möglichst viel gewinnen.

Frage: Kann ich vor Beginn des Spiels angeben, wieviel ich mindestens gewinne wenn ich bestmöglich spiele (egal, welche Züge mein Gegenspieler macht)?

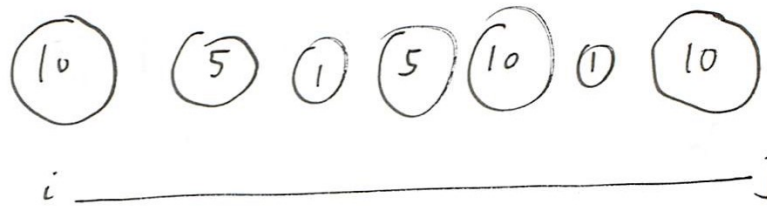
Antwort: Dynamische Programmierung!

Rekurrenz für eine beliebige Teilsequenz von  $i \dots j$ :

$$G(i,j) = \max\{ v_i + \min\{G(i+1,j-1), G(i+2,j)\}, v_j + \min\{G(i+1,j-1), G(i,j-2)\} \}$$

M

$\mathbb{I}$  Münzen mit Wert  $v_i \in \mathbb{N}$  für  $i \in M$



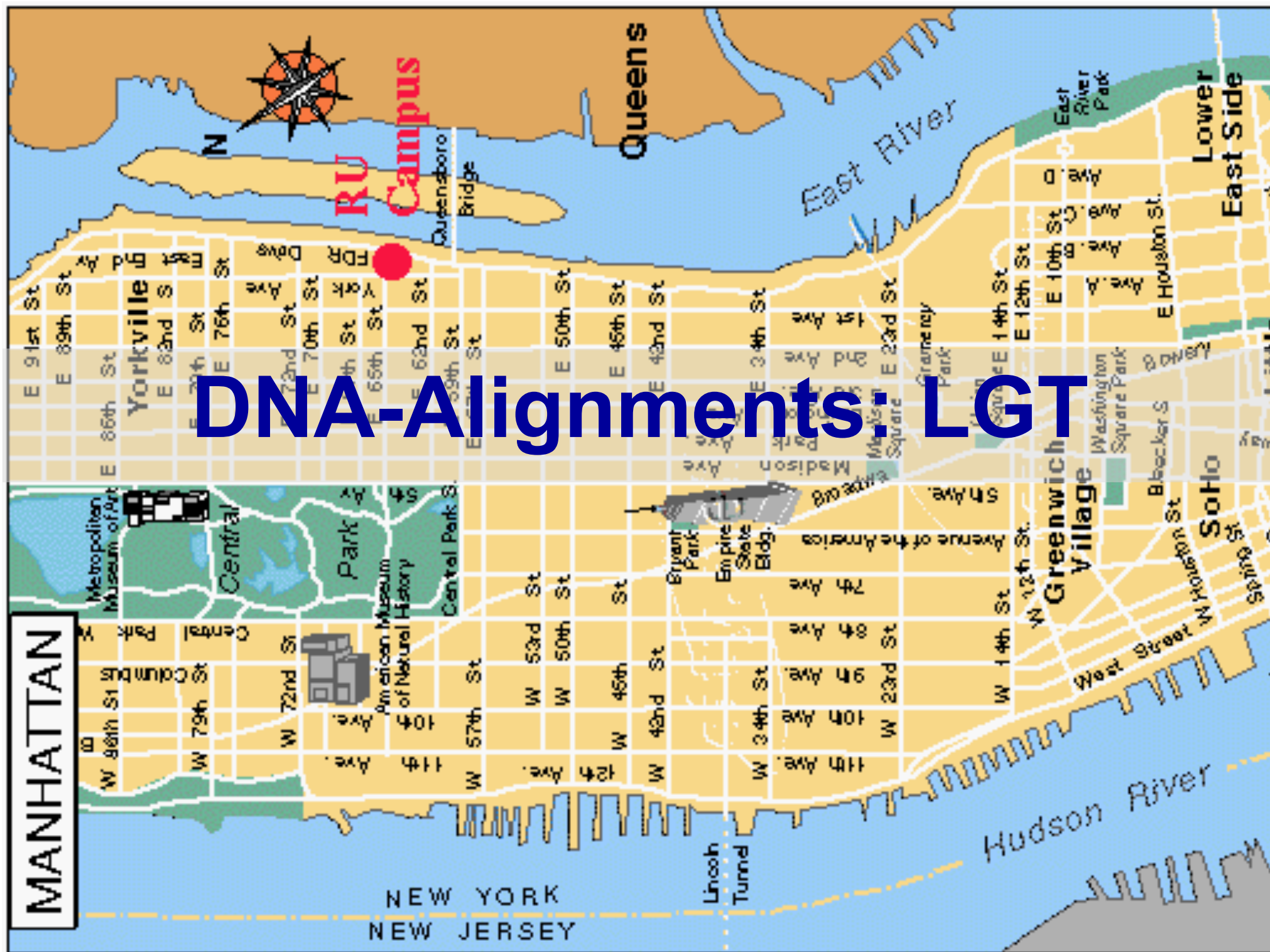
$$G(i, j) = \max \left\{ v_i + \min \left\{ G(i+2, j), G(i+1, j-1) \right\}, v_j + \min \left\{ G(i+1, j-2), G(i, j-1) \right\} \right\}$$

// Aufgabe: progr. , rek. / DP

$$G(i, i) = v_i$$

$$G(i, i+1) = \max \{ v_i, v_{i+1} \}$$





# Alignments: 2-zeilige Darstellung

Für 2 DNA-Sequenzen:

**v:**        **ATGTTAT**         $m = 7$

**w:**        **ATCGTA**         $n = 6$

**Alignment:**  $2 \times k$  -Matrix mit  $k \geq \max(m, n)$

<b>v</b>	A	T	--	G	T	T	A	T
<b>w</b>	A	T	C	G	T	--	A	--

5 Matches

1 Einfügung

2 Löschungen

Indels (= INsertions & DEletionS)

# Längste gemeinsame Teilsequenz: Problem

Ziel: Finde die längsten identischen, geordneten Teilsequenzen

Eingabe: 2 Sequenzen

$$\mathbf{v} = v_1 v_2 \dots v_m$$

$$\mathbf{w} = w_1 w_2 \dots w_n$$

Ausgabe: Eine Abfolge von Positionen  $i$  für  $\mathbf{v}$  und  $j$  in  $\mathbf{w}$  mit

$$1 \leq i_1 < i_2 < \dots < i_t \leq m$$

$$1 \leq j_1 < j_2 < \dots < j_t \leq n$$

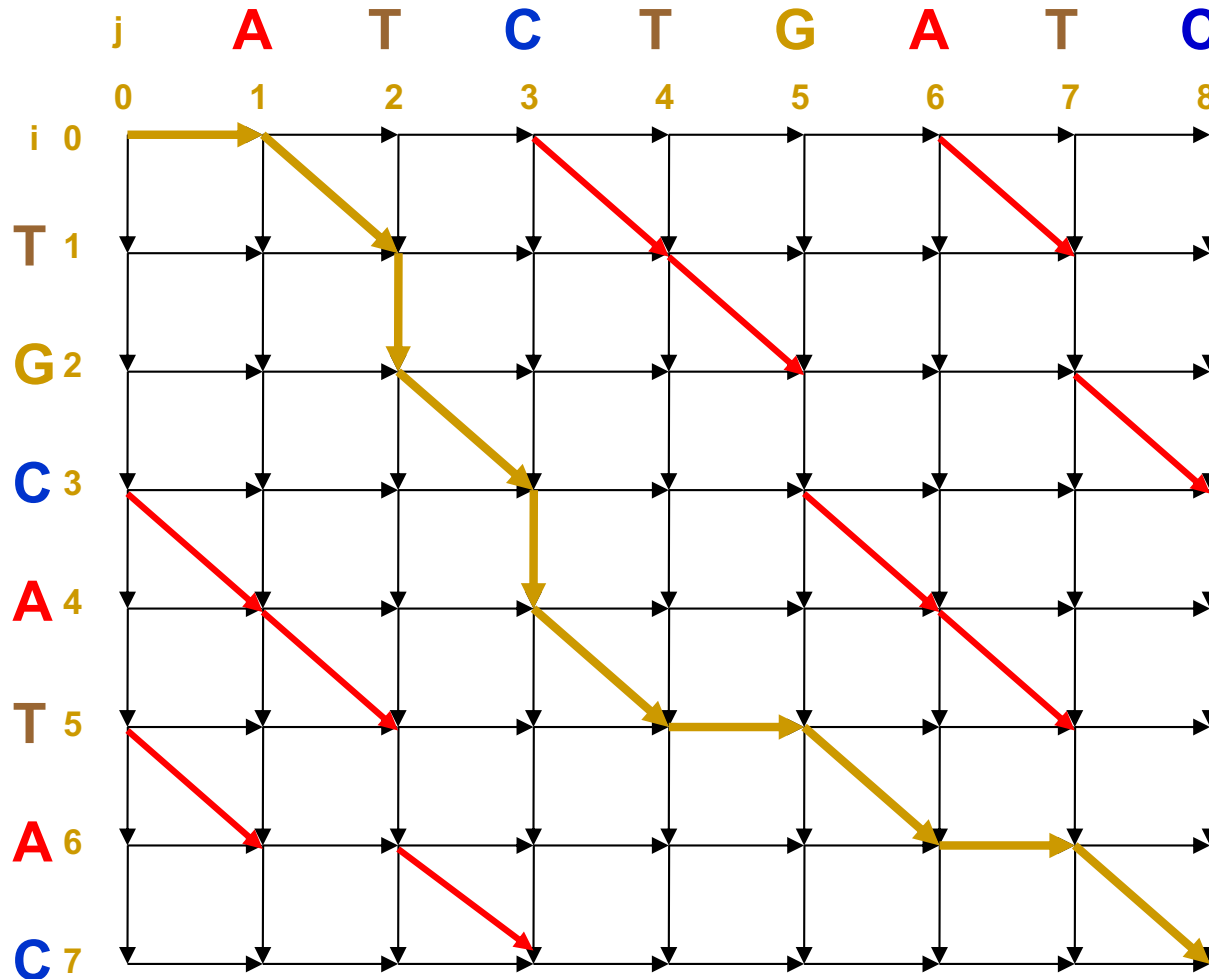
so dass  $t$  maximal ist und  $v_{i_s} = w_{j_s}$  für alle  $s$

Koordinaten $i$ :	0	1	2	2	3	3	4	5	6	7	8
Elemente von $v$ :	A	T	--	C	--	T	G	A	T	C	
Elemente von $w$ :	--	T	G	C	A	T	--	A	--	C	
Koordinaten $j$ :	0	0	1	2	3	4	5	5	6	6	7

$(0,0) \rightarrow (1,0) \rightarrow (2,1) \rightarrow (2,2) \rightarrow (3,3) \rightarrow (3,4) \rightarrow (4,5) \rightarrow (5,5) \rightarrow (6,6) \rightarrow (7,6) \rightarrow (8,7)$



# LGT & Manhattan-Tourismus



Jeder Pfad ist eine  
gemeinsame  
Teilsequenz

Jede diagonale Kante  
gibt 1 zusätzliches  
Element

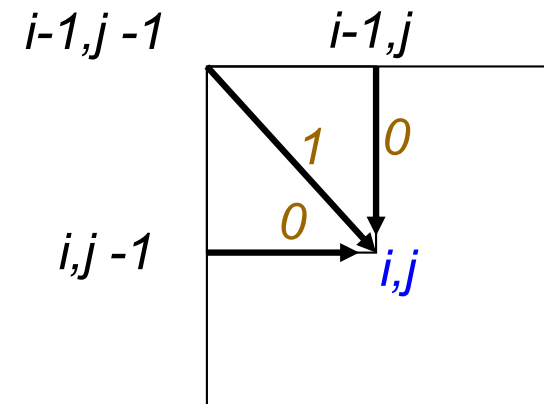
**LGT Problem:** Finde  
den Pfad mit den  
meisten diagonalen  
Kanten

# LGT: Rekurrenz

Länge der LGT bis (i, j)

:

$$s_{i,j} = \max \begin{cases} s_{i-1,j} \\ s_{i,j-1} \\ s_{i-1,j-1} + 1 \text{ wenn } v_i = w_j \end{cases}$$



# Edit-Abstand

:= Minimale Anzahl elementarer Operationen (Einfügungen, Löschungen, Substitutionen), um String **v** in String **w** zu überführen

:=  $d_E(\mathbf{v}, \mathbf{w})$

Hamming-Abstand:

Edit-Abstand:

Vergleicht immer  
**i-ten** Buchstaben von **v** mit  
**i-tem** Buchstaben von **w**

**V** = ATATATA  
**W** = TATATA

$$d_H(\mathbf{v}, \mathbf{w}) = 8$$

(trivial)

1 Einfügung  
1 Löschung →

Vergleicht (für 'beste'  $i, j$ )  
**i-ten** Buchstaben von **v** mit  
**j-tem** Buchstaben von **w**

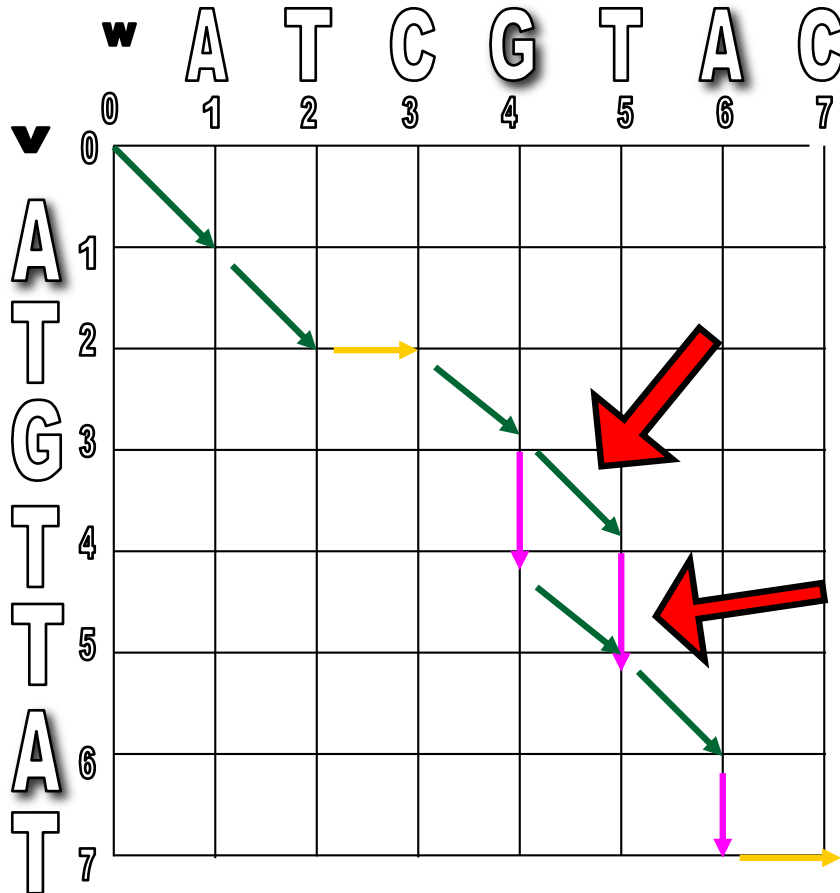
**V** = -ATATATA  
**W** = TATATA-

$$d_E(\mathbf{v}, \mathbf{w}) = 2$$

(nicht-trivial)

Beispiel: ATGTTAT → ATCGTAC

# LGT im Edit-Graphen



↓ und → sind Indels  
in **v** und **w** mit Score 0.

↘ sind Matches mit Score 1.

→ Alignment-Pfad hat Score 5

Alternative Pfade können  
gleichen Score haben!



# LGT: Dynamische Programmierung

	w	A	T	C	G	T	A	C
	0	1	2	3	4	5	6	7
v 0	0	0	0	0	0	0	0	0
A 1	0							
T 2	0							
G 3	0							
T 4	0							
T 5	0							
A 6	0							
T 7	0							

Initialisiere Reihe 0 und Spalte 0  
mit Score 0 (Indel-Präfix)

# LGT: Dynamische Programmierung

	w	A	T	C	G	T	A	C
	0	1	2	3	4	5	6	7
v 0	0	0	0	0	0	0	0	0
A 1	0	1	1	1	1	1	1	1
T 2	0	1	2	2	2	2	2	2
G 3	0	1	2	2	3	3	3	3
T 4	0	1	2	2	3	4	4	4
T 5	0	1	2	2	3	4	4	4
A 6	0	1	2	2	3	4	5	5
T 7	0	1	2	2	3	4	5	5

$$S_{i,j} = \max \begin{cases} S_{i-1,j-1} + 1 & \text{wenn } v_i = w_j \text{ (von NW)} \\ S_{i-1,j} & \text{(vom Norden)} \\ S_{i,j-1} & \text{(vom Westen)} \end{cases}$$



Manhattan-Touristen mit allen vertikalen und horizontalen Kanten = 0

# LGT Algorithmus

1. LGT(v,w)

2.     for  $i \leftarrow 1$  to  $n$

3.          $s_{i,0} \leftarrow 0$

4.     for  $j \leftarrow 1$  to  $m$

5.          $s_{0,j} \leftarrow 0$

6.     for  $i \leftarrow 1$  to  $n$

7.         for  $j \leftarrow 1$  to  $m$

8.             if ( $v_i = w_j$ )

9.                  $s_{i,j} \leftarrow \max (s_{i-1,j}, s_{i,j-1}, s_{i-1,j-1} + 1)$

10.             else

11.                  $s_{i,j} \leftarrow \max (s_{i-1,j}, s_{i,j-1})$

12.             case:

13.                  $s_{i,j} = s_{i-1,j}$  :  $b_{i,j} \leftarrow \text{“}\uparrow\text{”}$

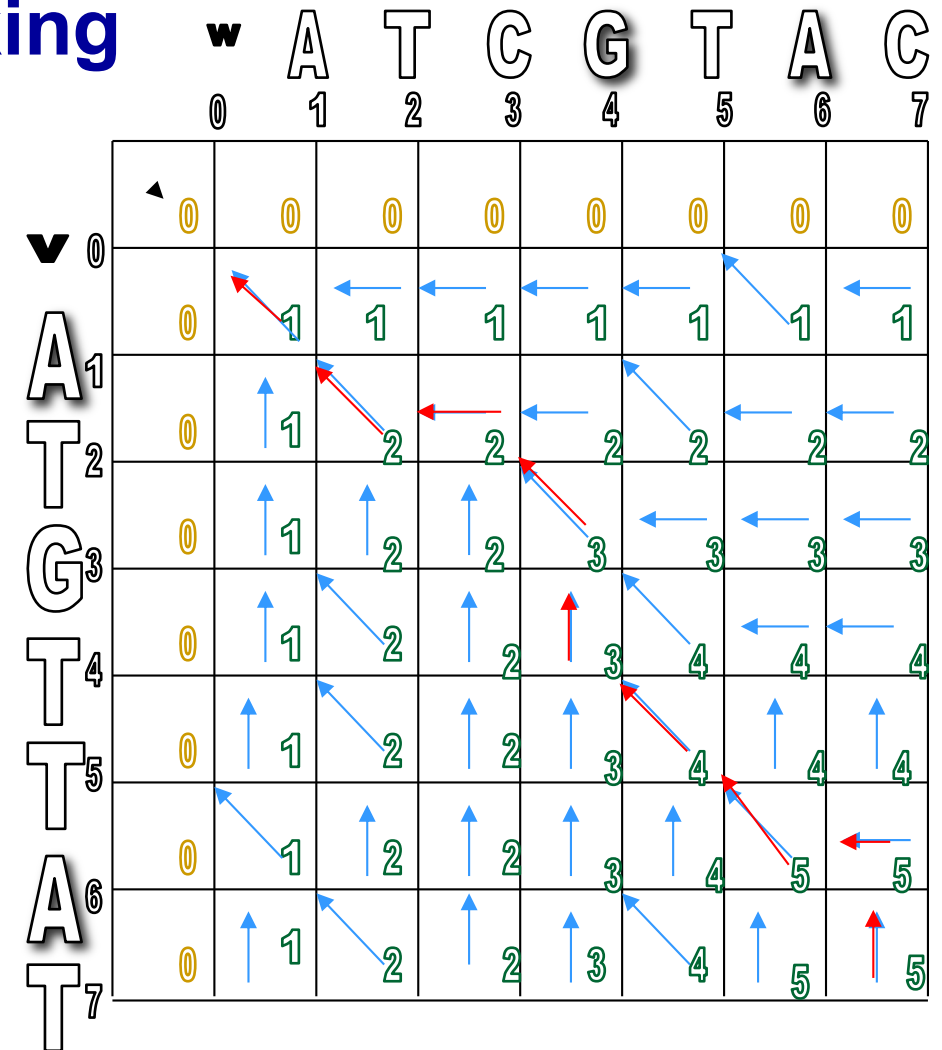
14.                  $s_{i,j} = s_{i,j-1}$  :  $b_{i,j} \leftarrow \text{“}\leftarrow\text{”}$

15.                  $s_{i,j} = s_{i-1,j-1} + 1$  and  $v_i = w_j$  :  $b_{i,j} \leftarrow \text{“}\nwarrow\text{”}$

16.     return ( $s_{n,m}, b$ )

$O(nm)$  um das Gitter zu füllen

# Ausgabe LGT: Backtracking



1. PrintLGT(b,v,i,j)
2.   if  $i = 0$  or  $j = 0$
3.     return
4.   case
5.      $b_{i,j} = \nwarrow$  : PrintLGT(b,v,  $i-1, j-1$ )
6.         print  $v_j$
7.      $b_{i,j} = \leftarrow$  : PrintLGT(b,v,  $i, j-1$ )
8.      $b_{i,j} = \uparrow$  : PrintLGT(b,v,  $i-1, j$ )