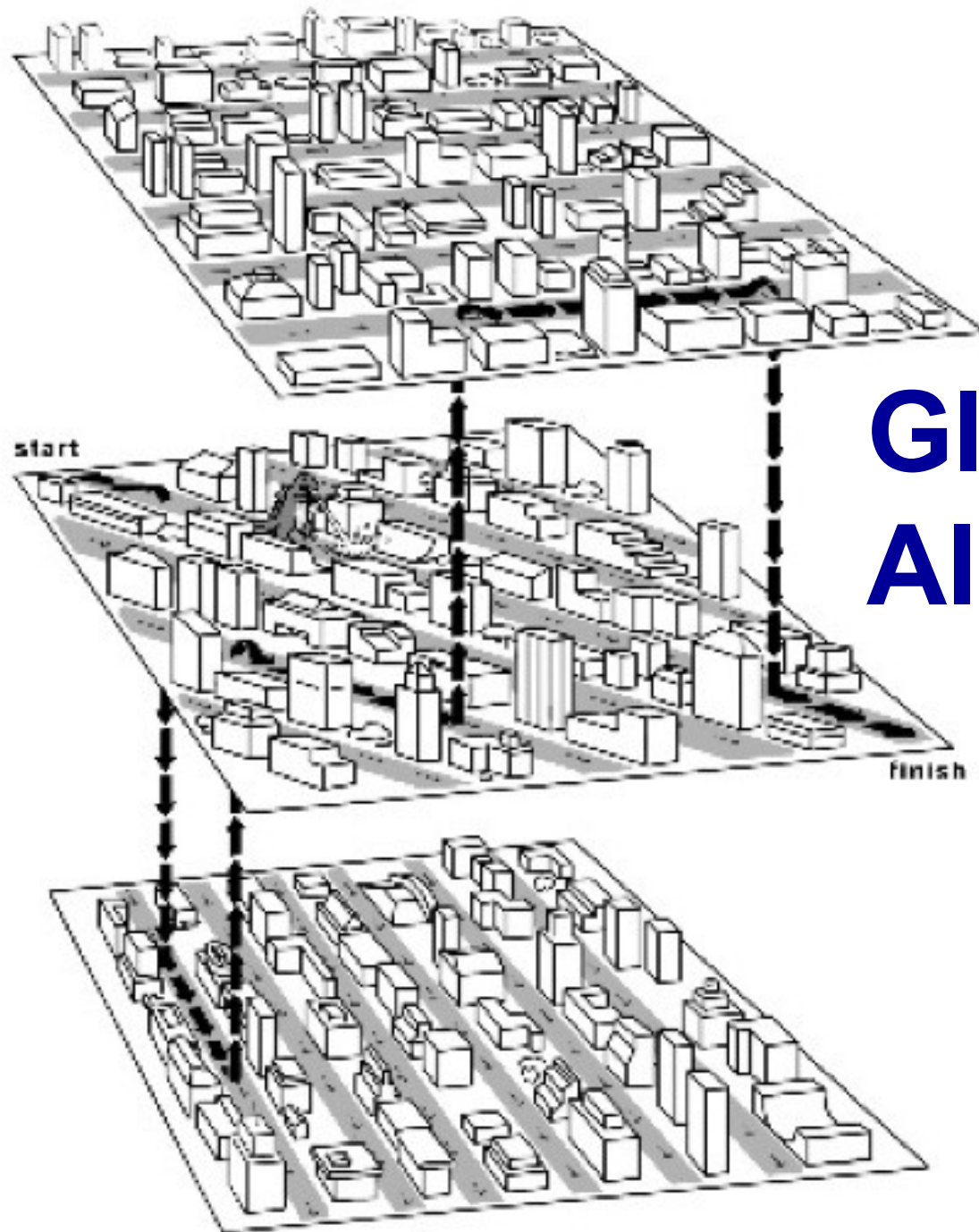


Algorithmen in der Bioinformatik

5. Sequenz-Vergleiche - Dynamische Programmierung II -

Prof. Dr. Gunnar Klau

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



Globales Alignment

LGT → Globales Alignment

- LGT erlaubt nur Matches & Indels (keine Mismatches)
- Score = 1 für Matches, 0 für Indels
- Wir können Indels mit negativem Score bestrafen
- Einfachstes *Scoringschema*:
 - $+1$: Belohnung für **Match**
 - $-\mu$: Strafe für **Mismatch**
 - $-\sigma$: Strafe für **Indel**

$$\text{Score} = \# \text{ Matches} - \mu (\# \text{ Mismatches}) - \sigma (\# \text{ Indels})$$

Rekurrenz:

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

Scoring-Matrizen

DNA [ATGC-]:

(4+1) x (4+1) **Scoring-Matrix** $\delta(v_i, w_j)$

Aminosäuren [ARNDCSEQGHILKMFPSTWYV-]:

(20+1) x (20+1) **Scoring-Matrix** $\delta(v_i, w_j)$

⇒ Vereinfachte Rekurrenz-Relation:

$$s_{i,j} = \max \begin{cases} s_{i-1,j-1} + \delta(v_i, w_j) \\ s_{i-1,j} + \delta(v_i, -) \\ s_{i,j-1} + \delta(-, w_j) \end{cases}$$

Herstellung von Scoring-Matrizen

- Basierend auf biologischen Daten
- Alinierte Sequenzen sind durch Mutationen aus einem gemeinsamen Vorfahren entstanden
- Mutationen können sehr verschiedene Auswirkungen auf Fitness haben → Scores $\delta(v_i, w_j)$ können sehr unterschiedlich sein
- Aminosäure-Austausche sind umso wahrscheinlicher, je ähnlicher die Codons sind.
- Aminosäure-Veränderungen sind weniger kritisch, wenn sie physiko-chemische Eigenschaften erhalten, z.B.:
 - polar (Aspartat → Glutamat)
 - nicht-polar (Alanin → Valin)
 - pos. Ladung (Arginin → Lysin)

AKRANR
KAAANK
-1 -1 -2 +5 +7 +3

	Alanin	Arginin	Asparagin	Lysin
	A	R	N	K
A	5	-2	-1	-1
R	-	7	-1	3
N	-	-	7	0
K	-	-	-	6

Alignment von kodierender DNA

Der genetische Code ist redundant

		Second letter				
		U	C	A	G	
First letter	U	<div>UUU</div> <div>UUC</div> <div>UUA</div> <div>UUG</div> <div>Phenyl-alanine</div> <div>Leucine</div>	<div>UCU</div> <div>UCC</div> <div>UCA</div> <div>UCG</div> <div>Serine</div>	<div>UAU</div> <div>UAC</div> <div>UAA</div> <div>UAG</div> <div>Tyrosine</div> <div>Stop codon</div> <div>Stop codon</div>	<div>UGU</div> <div>UGC</div> <div>UGA</div> <div>UGG</div> <div>Cysteine</div> <div>Stop codon</div> <div>Tryptophan</div>	U C A G
	C	<div>CUU</div> <div>CUC</div> <div>CUA</div> <div>CUG</div> <div>Leucine</div>	<div>CCU</div> <div>CCC</div> <div>CCA</div> <div>CCG</div> <div>Proline</div>	<div>CAU</div> <div>CAC</div> <div>CAA</div> <div>CAG</div> <div>Histidine</div> <div>Glutamine</div>	<div>CGU</div> <div>CGC</div> <div>CGA</div> <div>CGG</div> <div>Arginine</div>	U C A G
	A	<div>AUU</div> <div>AUC</div> <div>AUA</div> <div>AUG</div> <div>Isoleucine</div> <div>Methionine; start codon</div>	<div>ACU</div> <div>ACC</div> <div>ACA</div> <div>ACG</div> <div>Threonine</div>	<div>AAU</div> <div>AAC</div> <div>AAA</div> <div>AAG</div> <div>Asparagine</div> <div>Lysine</div>	<div>AGU</div> <div>AGC</div> <div>AGA</div> <div>AGG</div> <div>Serine</div> <div>Arginine</div>	U C A G
	G	<div>GUU</div> <div>GUC</div> <div>GUA</div> <div>GUG</div> <div>Valine</div>	<div>GCU</div> <div>GCC</div> <div>GCA</div> <div>GCG</div> <div>Alanine</div>	<div>GAU</div> <div>GAC</div> <div>GAA</div> <div>GAG</div> <div>Aspartic acid</div> <div>Glutamic acid</div>	<div>GGU</div> <div>GGC</div> <div>GGA</div> <div>GGG</div> <div>Glycine</div>	U C A G

- DNA ist weniger konserviert als Aminosäuren
- Alignments von DNA werden besser, wenn man **zuerst Aminosäuren aliniert**, und dann das Alignment zurück in DNA übersetzt!

PAM - Scoring-Matrix für Aminosäuren

PAM = Point Accepted Mutation (Dayhoff et al.)

- 1 PAM = PAM₁ = 1% aller Positionen unterschiedlich
- PAM₁₀₀ = (PAM₁)¹⁰⁰ (Mutationen wirken auf das schon mutierte!)
→ nicht jede Position ist verändert
- manche Positionen mutieren mehrmals, andere selten

		A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W
PAM ₂₅₀	A	2																	
	R	-2	6																
	N	0	0	2															
	D	0	-1	2	4														
	C	-2	-4	-4	-5	4													
	Q	0	1	1	2	-5	4												
	E	0	-1	1	3	-5	2	4											
	G	1	-3	0	1	-3	-1	0	5										
	H	-1	2	2	1	-3	3	1	-2	6									
	I	1	2	2	2	2	2	2	2	2	5								

Aminosäure	3-Buchstaben Code	1-Buchstaben Code	Molekular-gewi
Alanin	Ala	A	89
Arginin	Arg	R	174
Aspargin	Asn	N	132
Asparginsäure	Asp	D	133
Cystein	Cys	C	121
Glutamine	Gln	Q	146
Glutaminsäure	Glu	E	147
Glycin	Gly	G	175
Histidin	His	H	155
Isoleucin	Ile	I	131
Leucin	Leu	L	131
Lysin	Lys	K	146
Methionin	Met	M	149
Phenylalanin	Phe	F	165
Prolin	Pro	P	115
Serin	Ser	S	105
Threonin	Thr	T	119
Tryptophan	Trp	W	204
Tyrosin	Tyr	Y	181
Valin	Val	V	117

$M(i, j) \approx$ wie häufig ~~erscheint~~ ~~ist~~ mutiert
i-te AA zu *j*-ter AA in verwandten
 Algnm. Seq.

$M(i, j) =$ Anteil aller AA *i*,
 die mit AA *j* aligniert
 sind

C	C	C
A	C	R

$f(j) =$ Anteil AA *j* in den Seq

$$PAM^1(i, j) = \log \left(\frac{M(i, j)}{f(j)} \right)$$

eigentlich 10
 $\log_{10}(\text{fraction})$,
 hier $10 \log_{10}(1.5) =$
 ca. $10 * 0.176 \dots$

$$PAM^1(C, A) = \log \left(\frac{1/4}{1/6} \right) = + \dots$$

BLOSUM - Scoring-Matrix für Aminosäuren

BLOSUM = BLOck SUbstitution Matrix

- Scores aus beobachteten Häufigkeiten von Substitutionen in Blöcken von **lokalen** Alignments verwandter Proteine
- Subskript gibt an, wie ähnlich die Proteine einander waren (BLOSUM62 aus Sequenzen mit ca. 62% Identität)
- nicht extrapoliert
- PAM für verwandte, BLOSUM für divergente Sequenzen

BLOSUM₅₀

[illegible]

Globales Alignment: Problem

Ziel: Finde das beste Alignment für 2 Sequenzen bei vorgegebener Scoring-Matrix.

Eingabe: Strings **v**, **w** und eine Scoring-Matrix δ

Ausgabe: Ein Alignment zwischen **v** und **w** mit maximalem Score

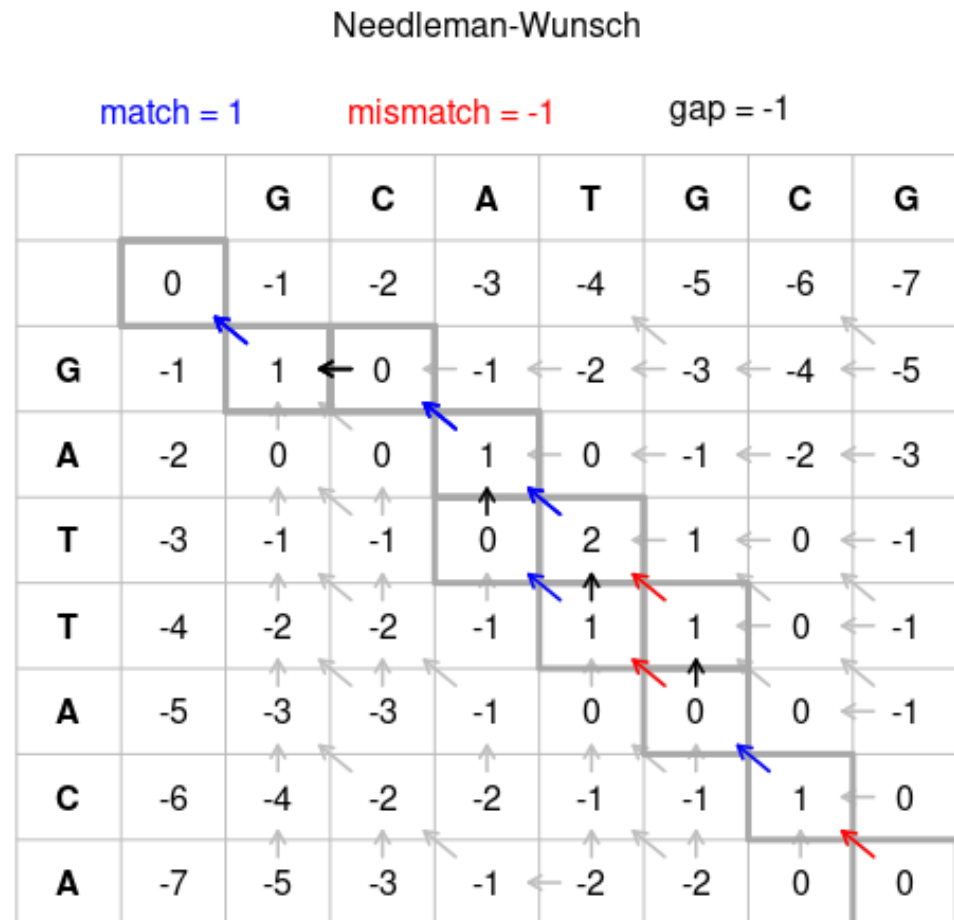
Lösung: Dynamische Programmierung (à la LGT) mit Rekurrenz

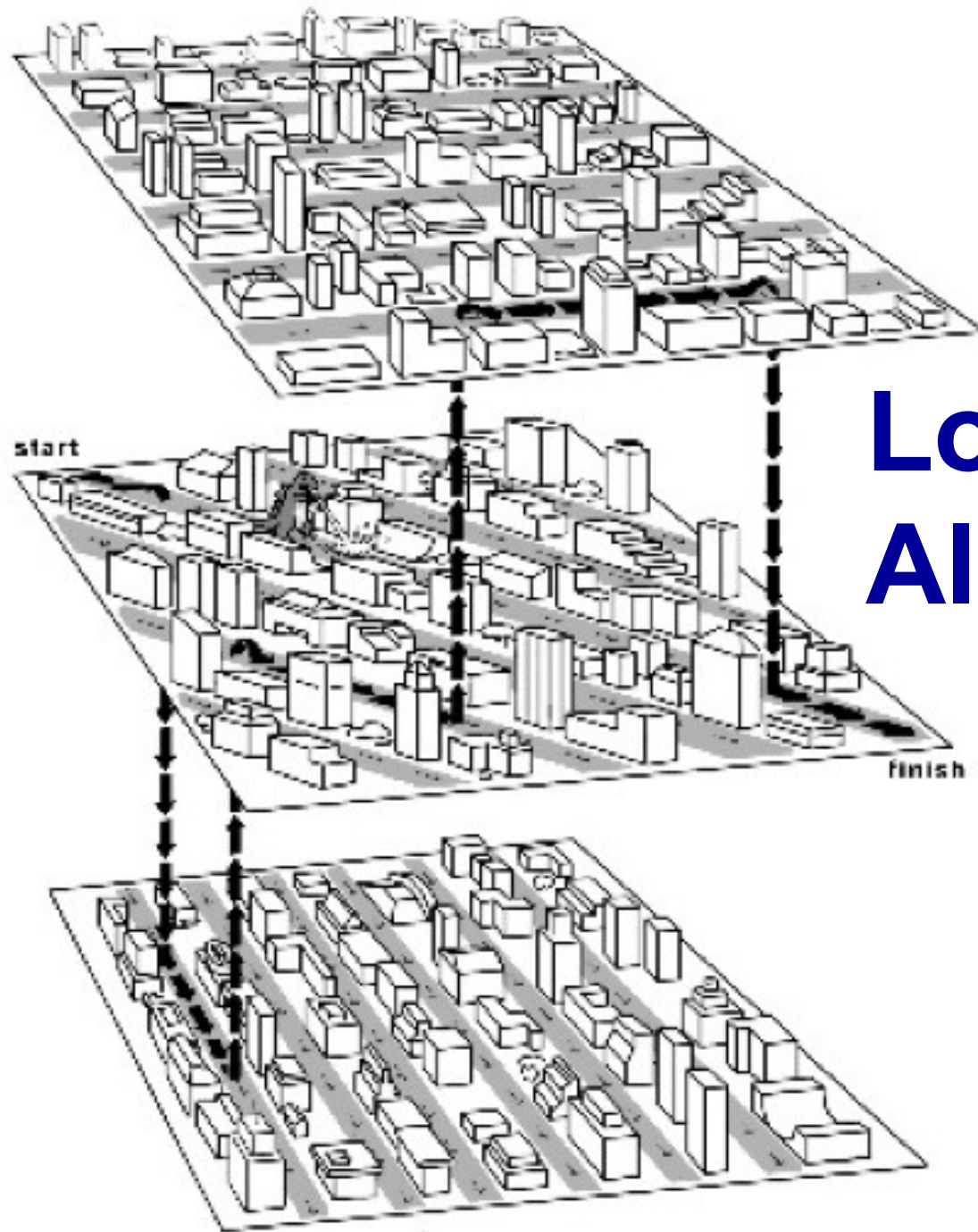
$$s_{i,j} = \max \begin{cases} s_{i-1,j-1} + \delta(v_i, w_j) \\ s_{i-1,j} + \delta(v_i, -) \\ s_{i,j-1} + \delta(-, w_j) \end{cases}$$

Needleman-Wunsch-Algorithmus

Globales Alignment: Beispiel

- $v = \text{GCATGCG}$, $w = \text{GATTACA}$. Match = 1, Mismatch = Gap = -1
- Siehe https://en.wikipedia.org/wiki/Needleman-Wunsch_algorithm





Lokales Alignment

Lokales vs. Globales Alignment

- Globales Alignment sucht den längsten Pfad zwischen $(0,0)$ und (n,m) im Edit-Graphen

```
--T--CC-C-AGT--TCTGT-CAGGGGACACG-A-GCATGCAGA-GAC
  |  ||| |  ||  |  ||  |  |||  ||  |  |  |||  |  |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG-T-CAGAT--C
```

- Lokales Alignment sucht den längsten Pfad zwischen beliebigen Knoten (i,j) und (i',j') im Edit-Graphen

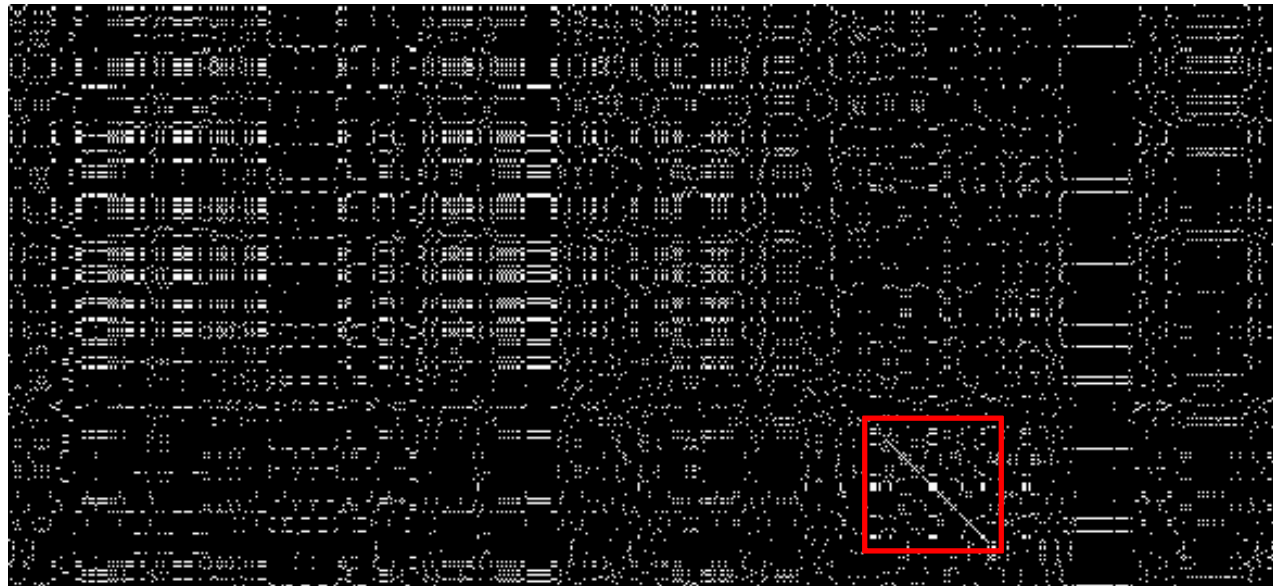
```
          tccCAGTTCTGTCAgggacacgagcatgcagagac
            |||||  |||||
aattgccgccgtcgttttcagCAGTTATGTCAGatc
```

Durch negative Scores $\delta(v_i, w_j)$ kann ein lokales Alignment den höheren Score haben!

Warum lokale Alignments?

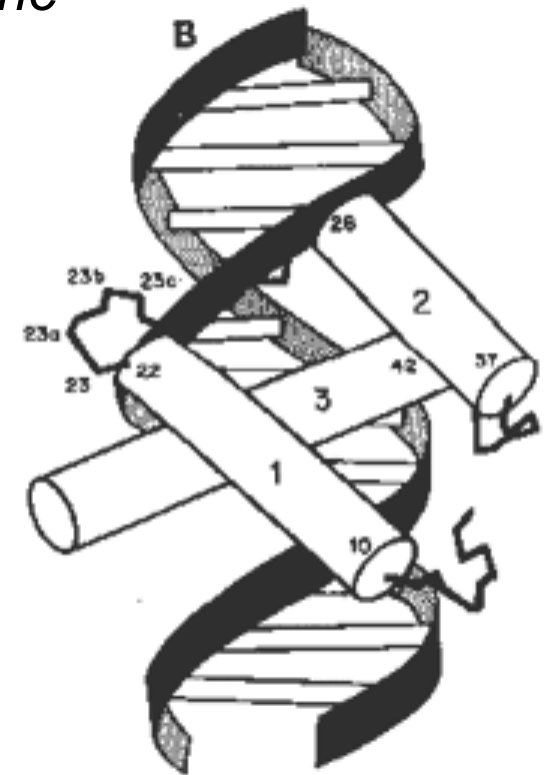
Zwei entfernt verwandte Gene mit konservierter Funktion können anhand von funktionell wichtigen Regionen erkannt werden

Beispiel: Homöobox-Gene haben einen *Homöo-Domäne*
- globales Alignment würde die übersehen!

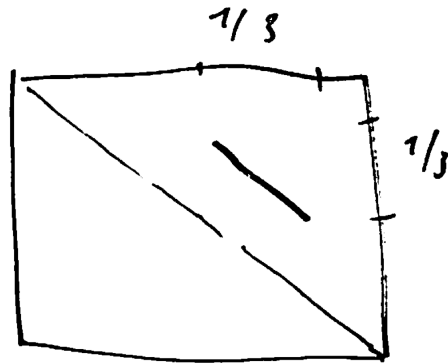


Homeobox

```
YTRFQTLELEKEFHFNHYLTRRRRIEIAHALCLTERQIKIWFQNRRMKLKK  
YTR-QTLELEKEFHFN-YLTRRRRIEIAHALCL-ERQIKIWFQNRRMK-KK
```



Fiktives Bsp.:



Perfektes lok. Ali, $\frac{1}{3}$ de Länge

Score (lok. Ali) \geq

$$\underbrace{\frac{1}{3} (1)}_{\text{Match}} + \underbrace{\frac{2}{3} (-\sigma)}_{\text{Gap}} + \frac{2}{3} (-0) =$$

Score (glob) \geq

$$\frac{1}{4} + \underbrace{\frac{3}{4} (-\mu)}_{\text{Mismatch.}}$$

$$= \frac{1}{4} - \frac{3}{4} \mu$$

$$\frac{1}{3} - \frac{4}{3} \sigma$$

$$\left| \frac{1}{3} - \frac{4}{3} \sigma < \frac{1}{4} - \frac{3}{4} \mu \right. \leftarrow \begin{array}{l} \frac{4}{3} \sigma - \frac{1}{3} > \frac{3}{4} \mu - \frac{1}{4} \\ \frac{4}{3} \sigma > \frac{3}{4} \mu + \frac{1}{12} \\ \sigma > \frac{9}{16} \mu + \frac{3}{16} \\ \text{Meistens kl.!} \end{array}$$

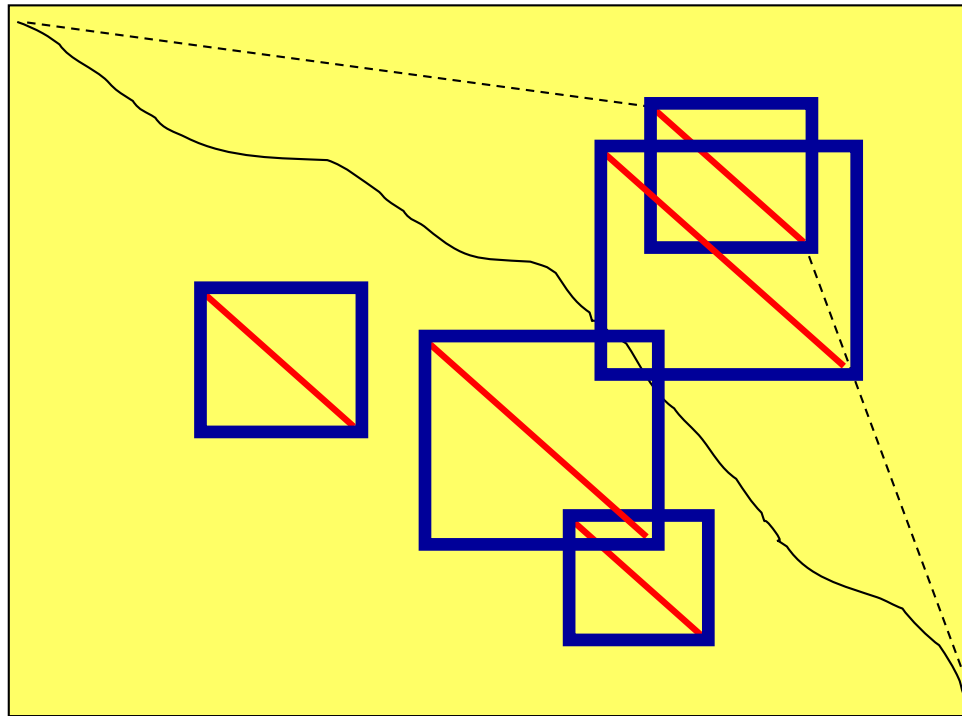
Lokales Alignment: Problem

Ziel: Finde das beste lokale Alignment zwischen 2 Sequenzen

Eingabe: Sequenzen **v**, **w** und Scoring-Matrix δ

Ausgabe: Alignment von Teilsequenzen von **v** und **w** mit maximalem Score (verglichen mit allen Alignments aller Teilsequenzen)

Lokales Alignment: Laufzeit



Laufzeit $O(n^4)$

- im $n \times n$ Gitter gibt es n^2 Knoten (i,j) , die als Quellen dienen können
- Berechnung aller Alignments von diesem Startpunkt braucht $O(n^2)$

Schnellere Lösung?
Freifahrtscheine!

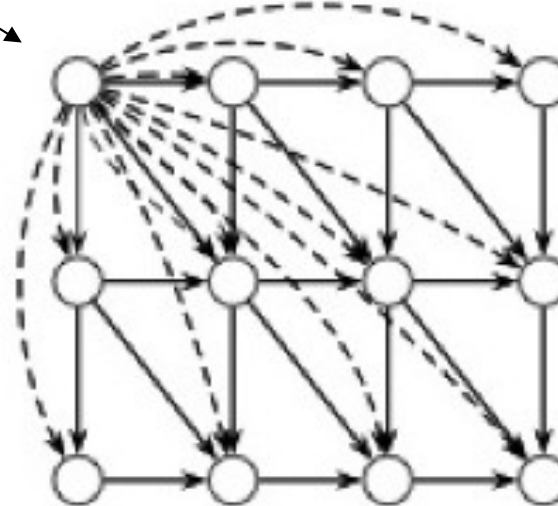
Lokales Alignment: Freifahrtscheine

Vertex (0,0)



Yiihaa, a free ride!

Gestrichelte Kanten sind Freifahrten von **(0,0)** zu allen Knoten



max. $s_{i,j}$ im Edit-Graphen entspricht dem besten lokalen Alignment

Rekurrenz:

$$s_{i,j} = \max \begin{cases} 0 \\ s_{i-1,j-1} + \delta(v_i, w_j) \\ s_{i-1,j} + \delta(v_i, -) \\ s_{i,j-1} + \delta(-, w_j) \end{cases}$$

Unterschied zur Rek. globalem Alignment

Needleman-Wunsch \rightarrow Smith-Waterman

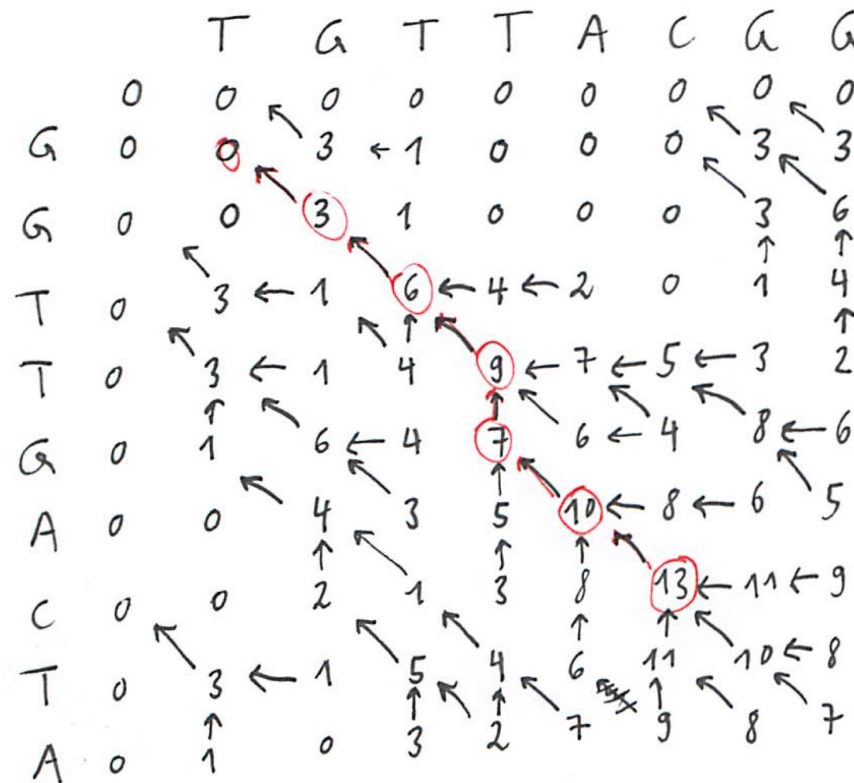
Smith-Waterman-Algorithmus

Sequenzen T G T T A C G G und
G G T T G A C T A

$$d(a, b) = \begin{cases} +3 & a=b \\ -3 & a \neq b \end{cases}$$

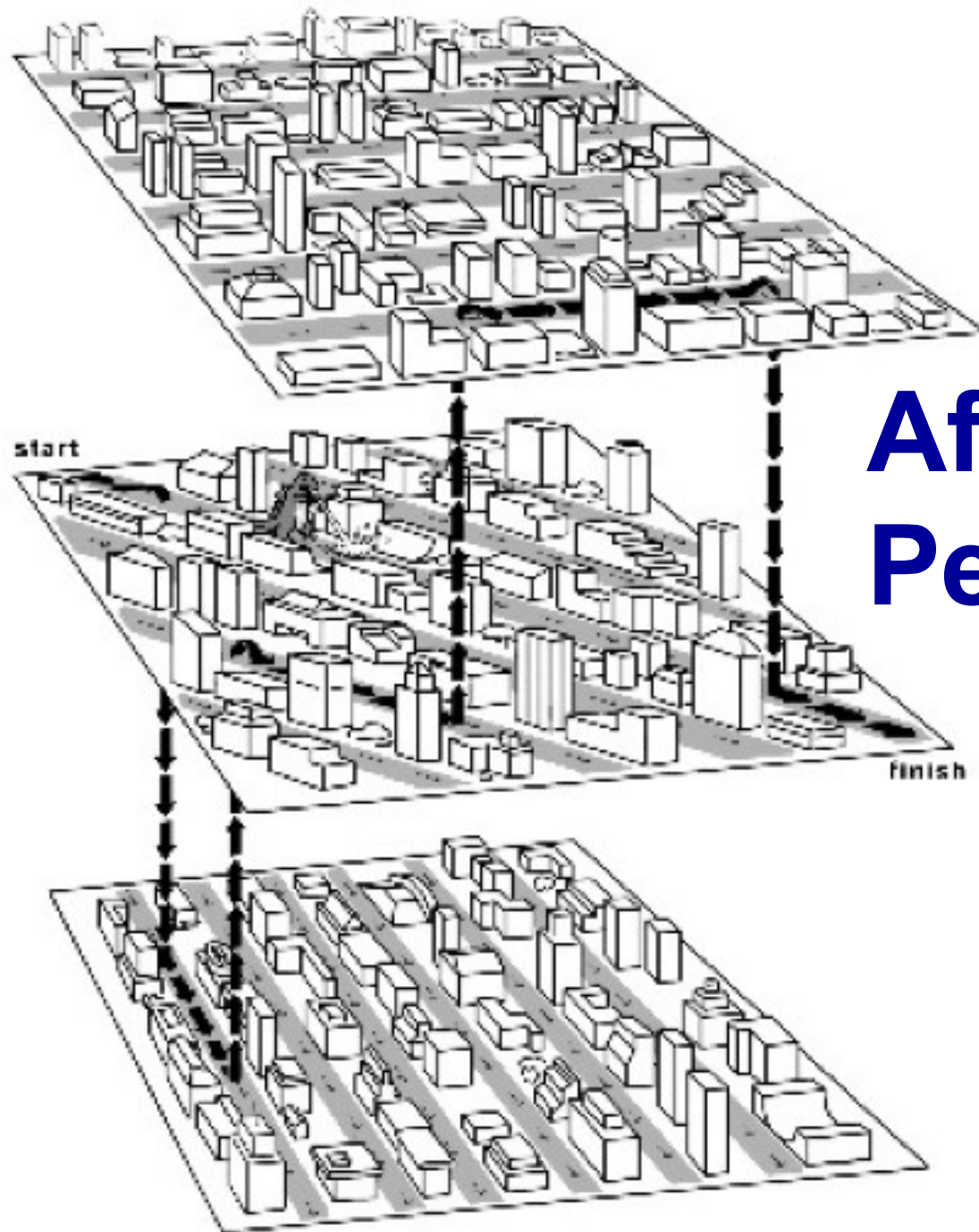
$$d(a, -) =$$

$$d(-, a) = -2$$



G T T - A C

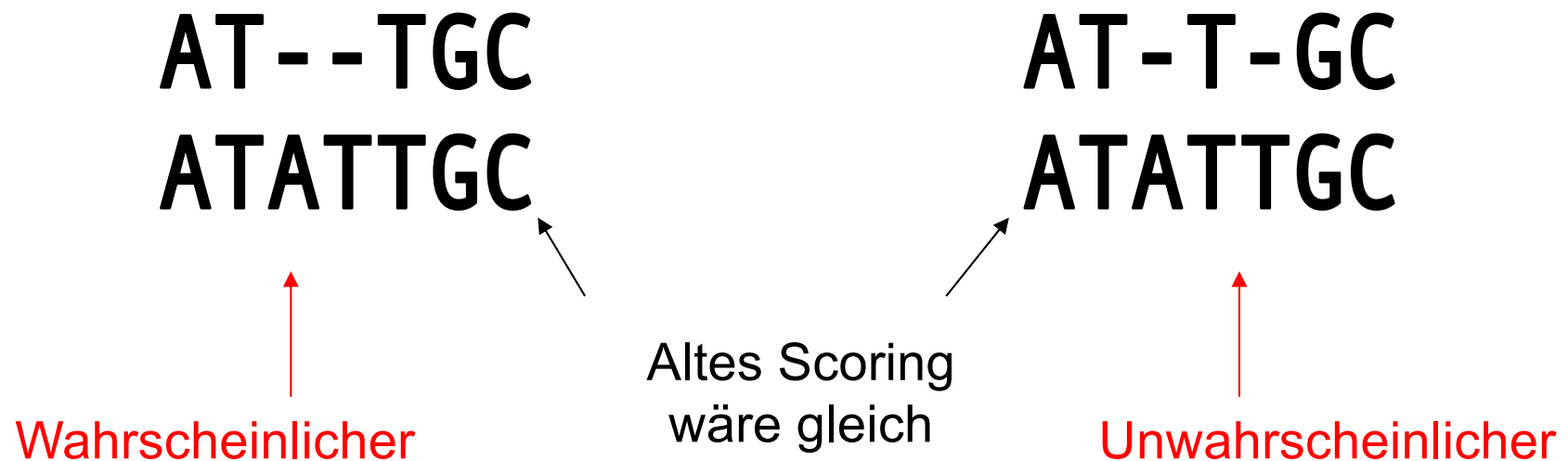
G T T G A C



Affine Lücken- Penalties

Affine Lücken-Penalties

- Bisher: Strafkosten σ für jedes Indel
- Biologie: eine Mutation mit k Indels ist wahrscheinlicher als k einzelne Mutationen:



Affine Lücken-Penalties

Score für eine Lücke der Länge k :

$$-(\sigma + \varepsilon(k-1))$$

mit $\sigma > 0$ Strafe für das Öffnen der Lücke:

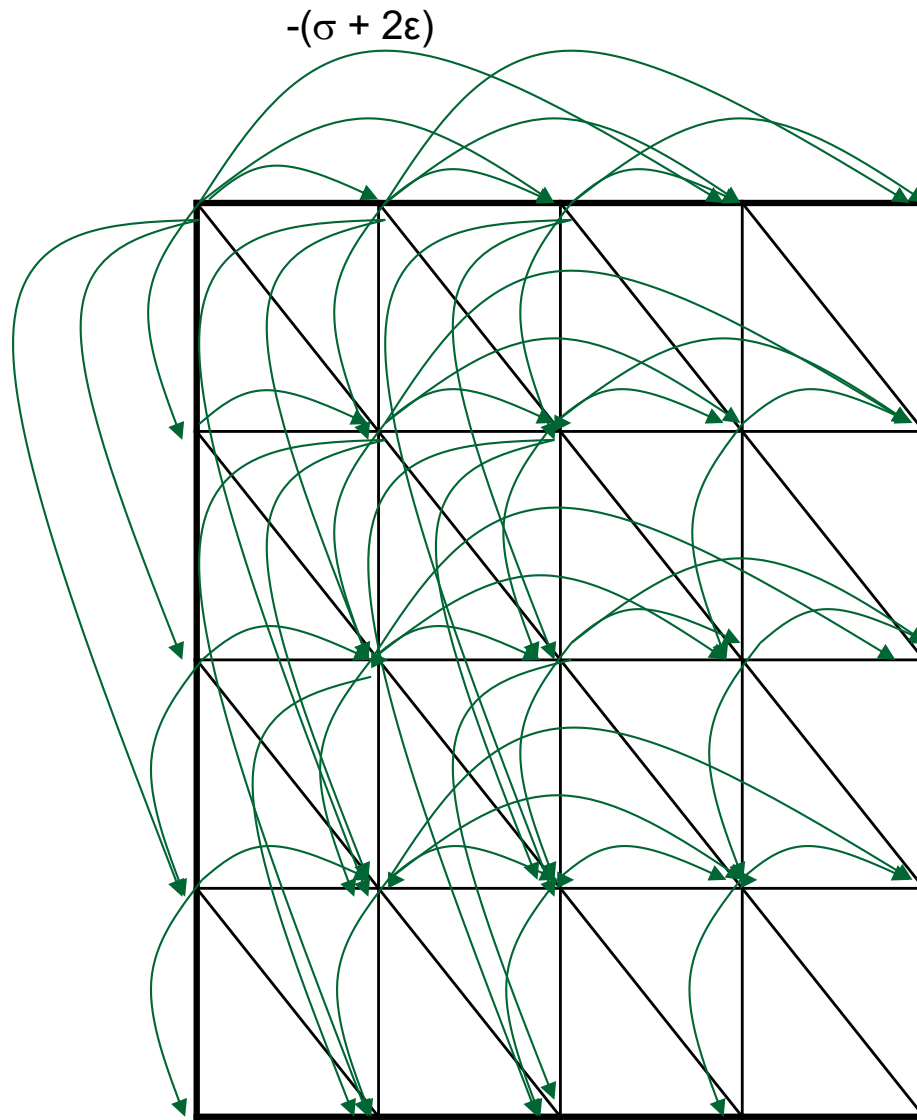
Gap opening penalty

und $\varepsilon > 0$ Strafe für das weitere Vergrößern der Lücke:

Gap extension penalty

In der Regel $\sigma > \varepsilon$: Vergrößern ist billiger (biologisch wahrscheinlicher)
als das Öffnen einer neuen Lücke

Affine Gaps und der Edit-Graph



Affine Gaps entsprechen langen (k)
horizontalen und vertikalen
Kanten mit Gewicht $-(\sigma + \varepsilon(k-1))$

Davon gibt es viele

→ Komplexität steigt von $O(n^2)$
auf $O(n^3)$

- Außer wir sind schlau!

Anzahl graue Kanten:

Betrachte Knoten (i, j) : ~~$(n-j)$~~ $(n-j)$ horiz. Kanten +
 $(n-i)$ vertikale Kanten.

$$\text{Insgesamt } \sum_{i=1}^n \sum_{j=1}^n (n-i) + \sum_{i=1}^n \sum_{j=1}^n (n-j)$$

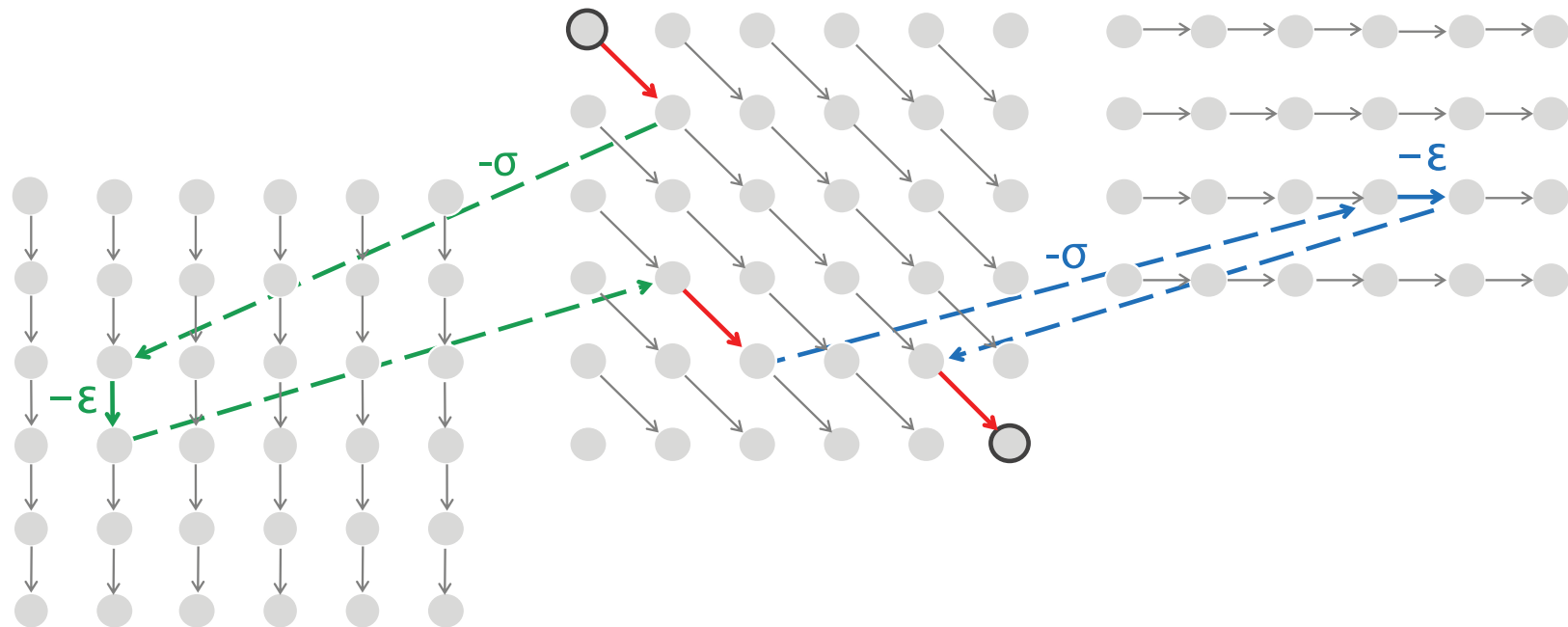
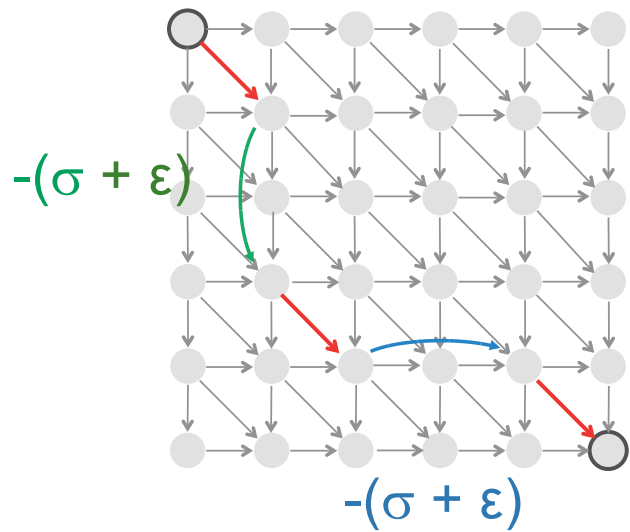
$$= 2 \sum_{i=1}^n \sum_{j=1}^n (n-i) = 2n^3 - \sum_{i=1}^n \sum_{j=1}^n i$$

$$= 2n^3 - \sum_{i=1}^n \frac{n(n+1)}{2}$$

$$= 2n^3 - \frac{n \cdot n \cdot (n+1)}{2}$$

$$= 2n^3 - \frac{n^3 + n^2}{2} = \mathcal{O}(n^3)$$

Beispiel und Idee



Affine Lücken: Manhattan in 3 Ebenen

$$\downarrow s_{i,j} = \max \begin{cases} \downarrow s_{i,j-1} - \varepsilon \\ s_{i,j-1} - \sigma \end{cases}$$

Setze Lücke in w fort

Öffne Lücke in w (Mitte \rightarrow oben)

$$s_{i,j} = \max \begin{cases} s_{i-1,j-1} + \delta(v_i, w_j) \\ \downarrow s_{i,j} \\ \rightarrow s_{i,j} \end{cases}$$

Match oder Mismatch

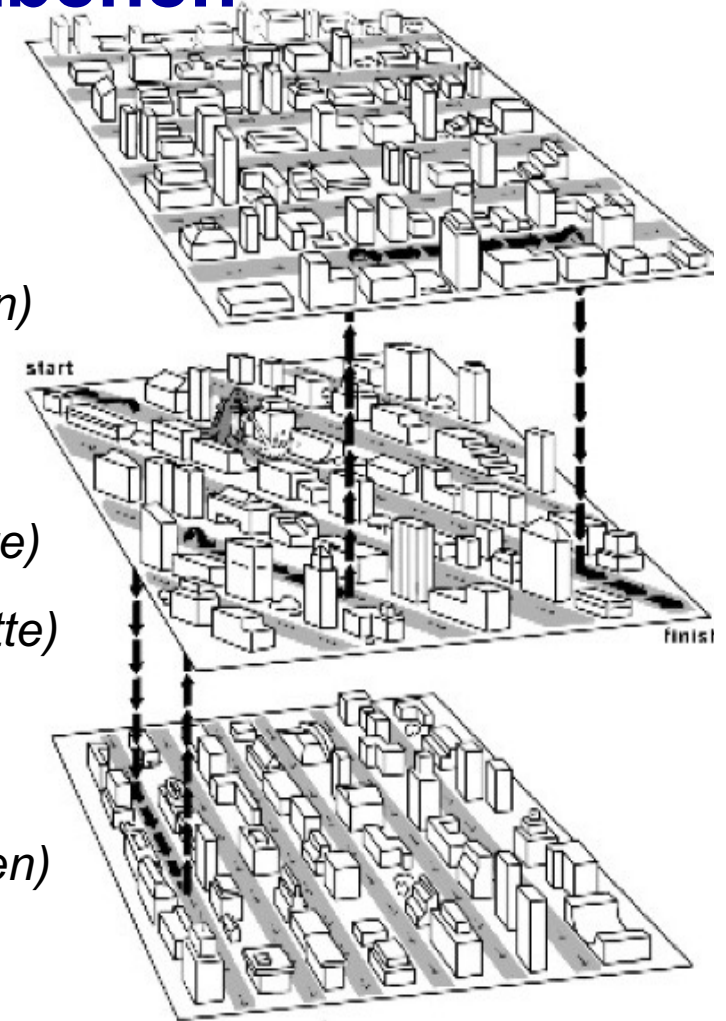
Ende Lücke in w (oben \rightarrow Mitte)

Ende Lücke in v (unten \rightarrow ?Mitte)

$$\rightarrow s_{i,j} = \max \begin{cases} \rightarrow s_{i-1,j} - \varepsilon \\ s_{i-1,j} - \sigma \end{cases}$$

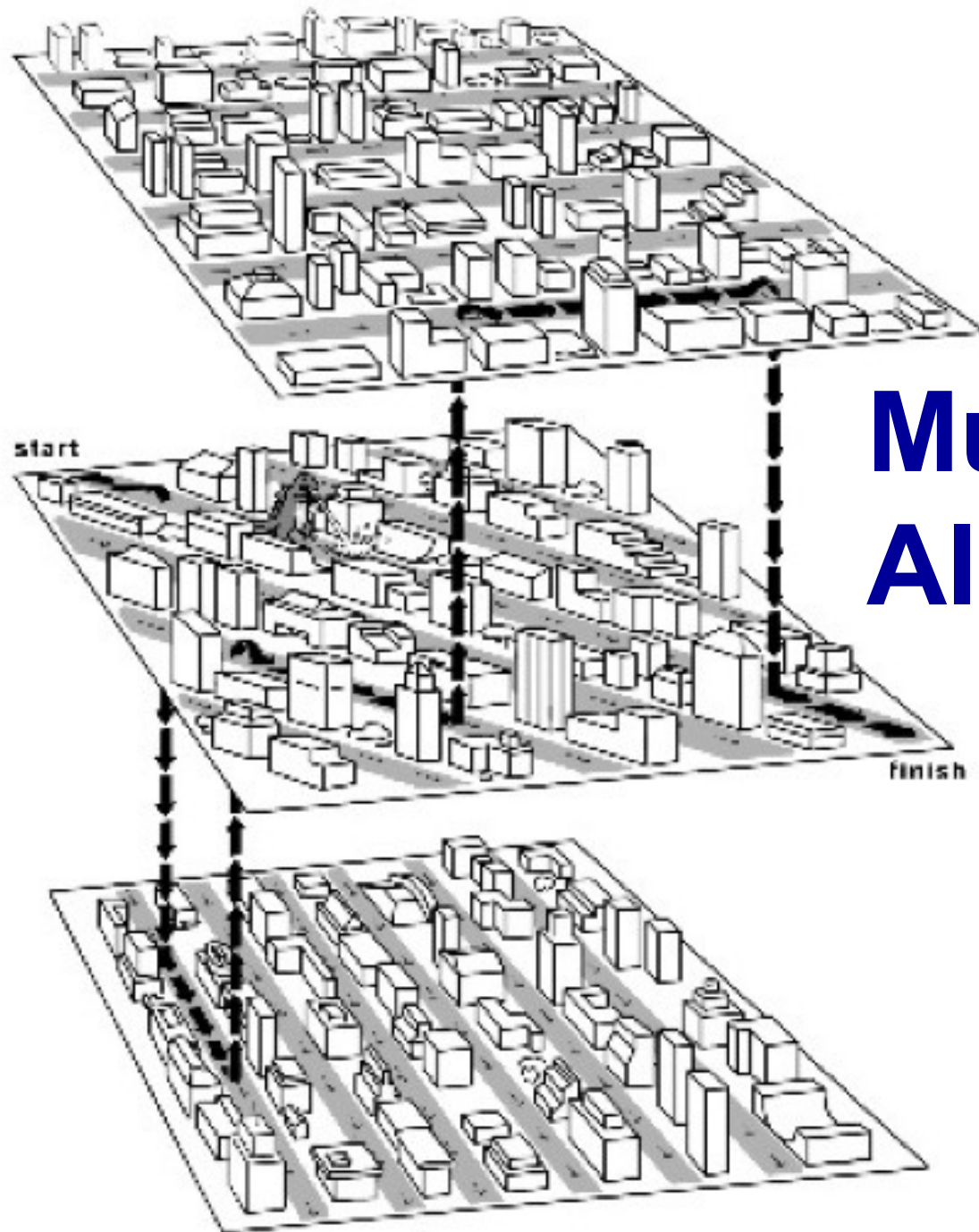
Setze Lücke in v fort

Öffne Lücke in v (Mitte \rightarrow ?unten)



- Sprung-Strafe $-\sigma$ für Bewegungen vom Hauptlevel nach oben oder unten
- Fortsetzungs-Strafe $-\varepsilon$ für Schritte im oberen oder unteren Level

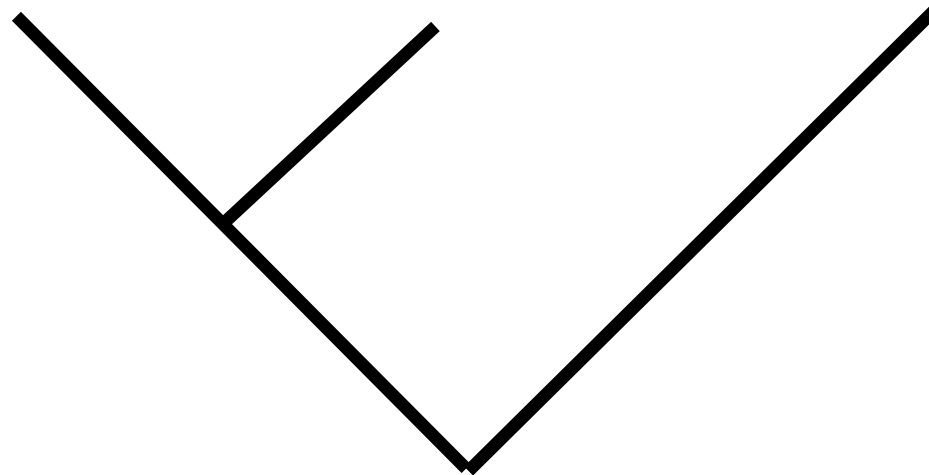
Noch weitere Änderungen für Initialisierung \rightarrow Übung.



Multiple Alignment

Multiples Alignment

Um Veränderungen zu verstehen, brauchen wir >2 Sequenzen



Verallgemeinerung von paarweisem Alignment

Alignment von 2 Sequenzen ~ 2-reihige Matrix

ATGCG-
ACGT-A

Alignment von 3 Sequenzen ~ 3-reihige Matrix (etc.)

AT-GCG-
A-CGT-A
ATCAC-A

...

Score: mehr Info (mehr konservierte Positionen)
→ besseres Alignment

Alignment-Pfade

0	1	1	2	3	4
	A	-	T	G	C
0	1	2	3	3	4
	A	A	T	-	C
0	0	1	2	3	4
	-	A	T	G	C

x-Koordinate

y-Koordinate

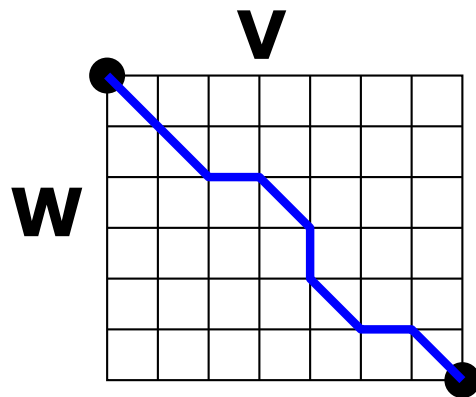
z-Koordinate

Pfad in 3-D:

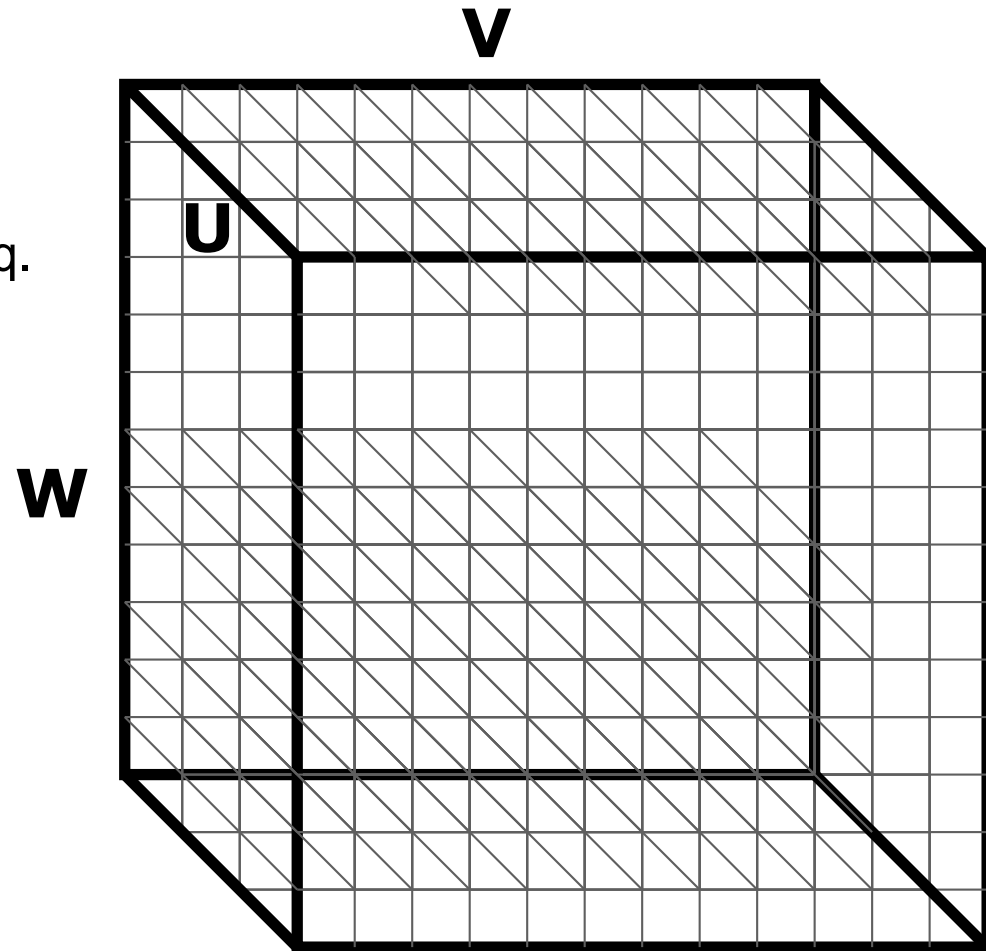
$(0,0,0) \rightarrow (1,1,0) \rightarrow (1,2,1) \rightarrow (2,3,2) \rightarrow (3,3,3) \rightarrow (4,4,4)$

2-D vs. 3-D Alignment-Gitter

- Selbe Strategie wie in 2-D
- 3-D “Manhattan-Würfel”, jede Achse repräsentiert 1 Seq.
- Gehe von Quelle zu Senke (für globales Alignment)

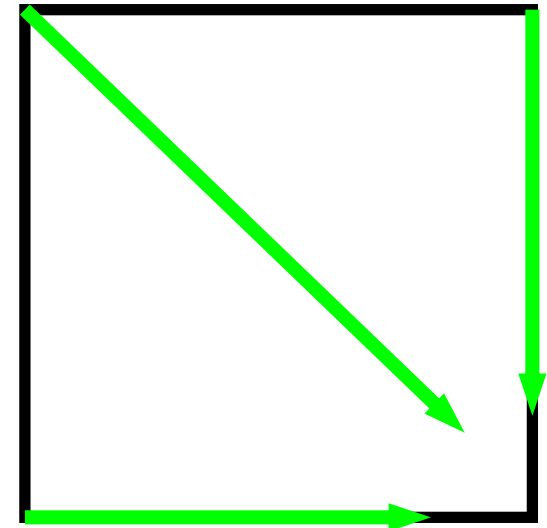
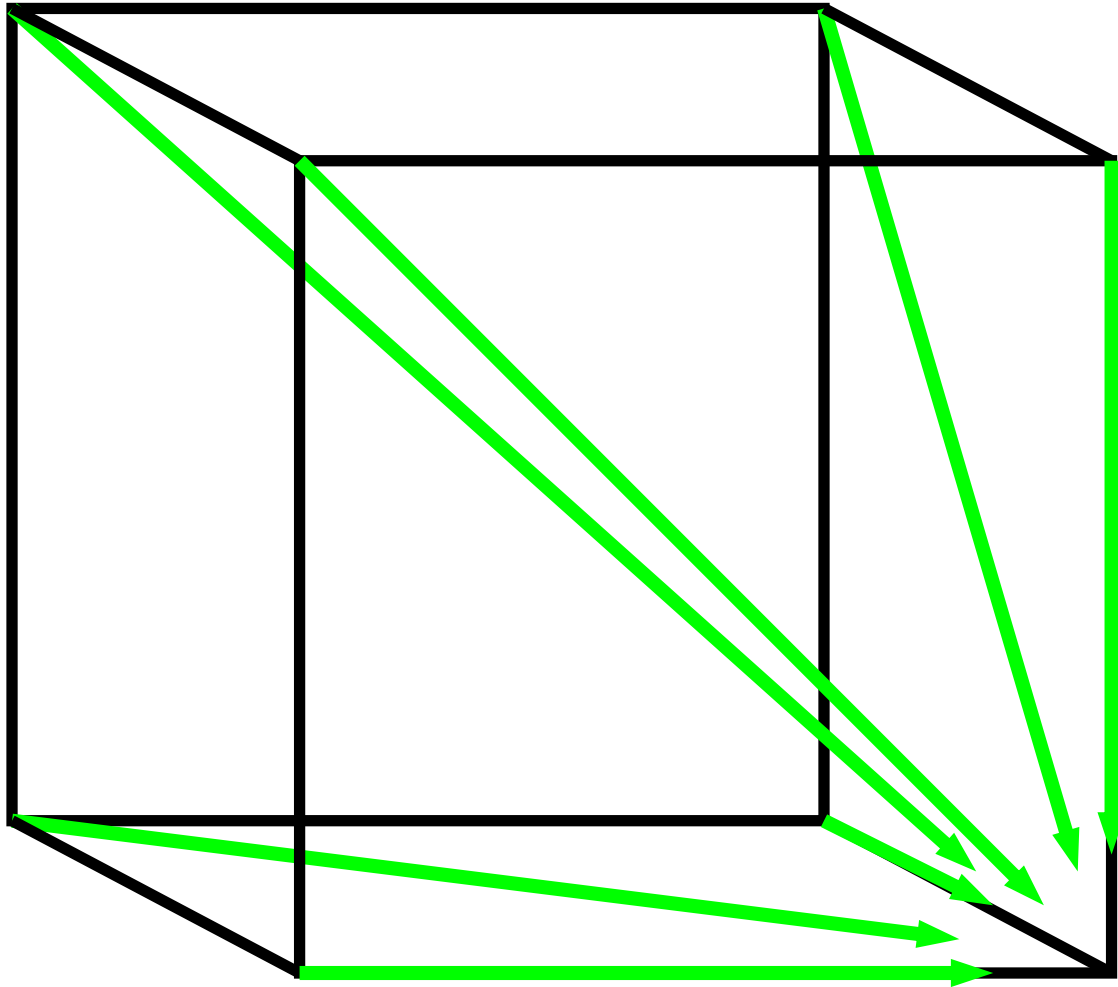


2-D edit graph



3-D edit graph

3-D vs. 2-D Alignment-Zellen



In **2-D**, 3 Kanten
pro Quadrat

In **3-D**, 7 Kanten
pro Würfel

Multiples Alignment: Rekurrenz

$$s_{i,j,k} = \max$$

$$s_{i-1,j-1,k-1} + \delta(v_i, w_j, u_k)$$

$$s_{i-1,j-1,k} + \delta(v_i, w_j, _)$$

$$s_{i-1,j,k-1} + \delta(v_i, _, u_k)$$

$$s_{i,j-1,k-1} + \delta(_, w_j, u_k)$$

$$s_{i-1,j,k} + \delta(v_i, _, _)$$

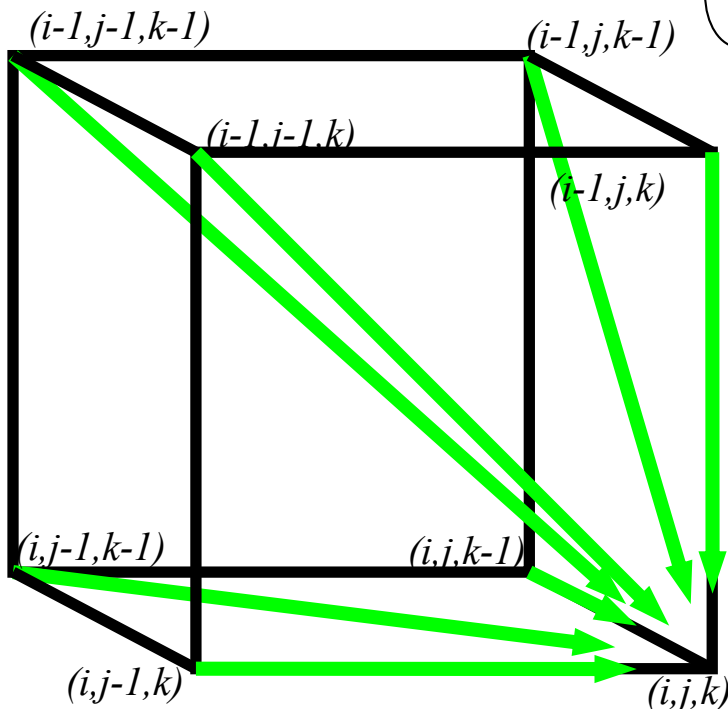
$$s_{i,j-1,k} + \delta(_, w_j, _)$$

$$s_{i,j,k-1} + \delta(_, _, u_k)$$

Würfel-Diagonale:
0 Indels

Flächen-Diagonale:
1 Indel

Würfelkante:
2 Indels



$\delta(x, y, z)$ ist ein Eintrag der 3-D Scoring-Matrix

Multiples Alignment: Laufzeit

Für 3 Sequenzen der Länge n :

$$7n^3 \rightarrow O(n^3)$$

Für k Sequenzen: k -dimensionales Manhattan:

$$(2^k - 1)(n^k) \rightarrow O(2^k n^k)$$

Dynamische Programmierung kann leicht auf k Sequenzen verallgemeinert werden, aber Laufzeit ist inpraktikabel

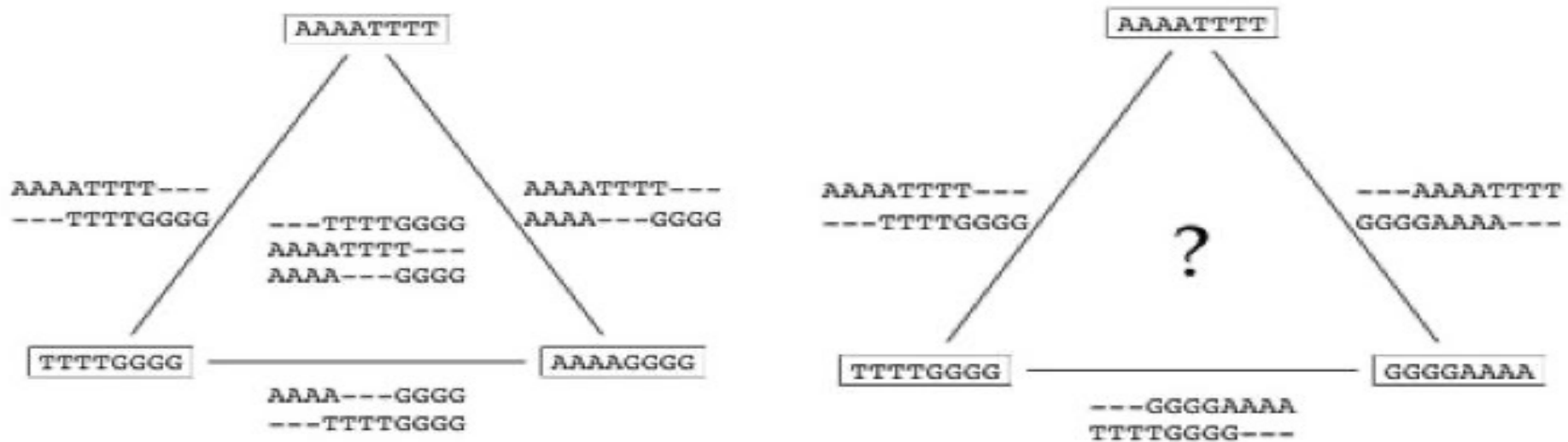
Multiples Alignment \leftrightarrow paarweises Alignment

Multiples Alignment induziert 3 paarweise Alignments:

x: AC-GCGG-C
y: AC-GC-GAG
z: GCCGC-GAG

→ **x:** ACGCGG-C; **x:** AC-GCGG-C; **y:** AC-GCGAG
y: ACGC-GAC; **z:** GCCGC-GAG; **z:** GCCGCGAG

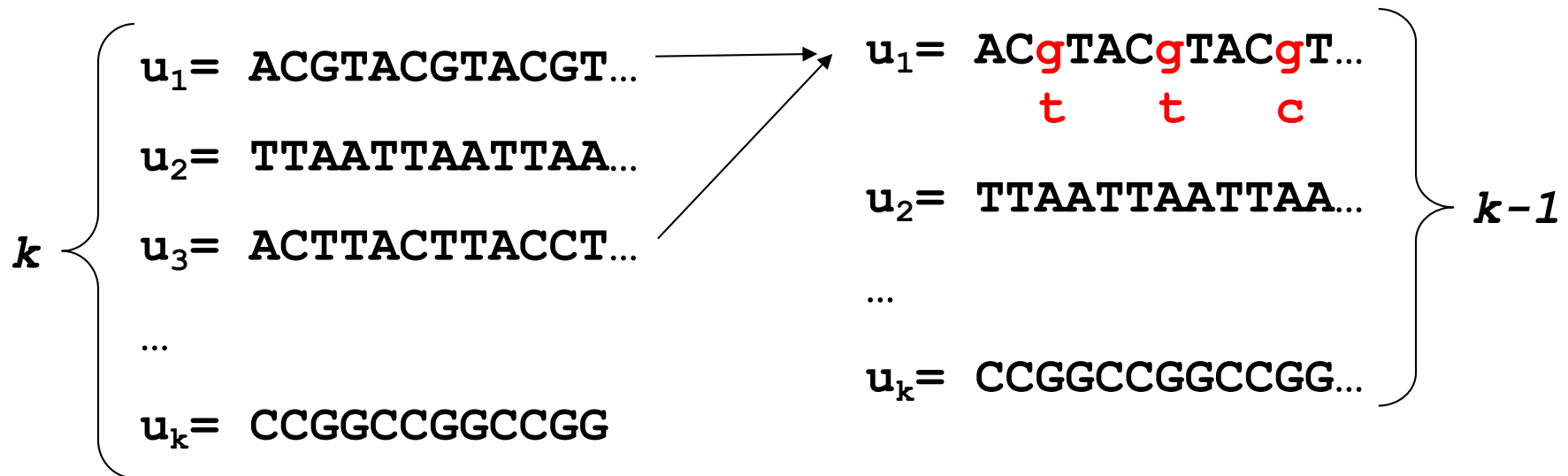
Induzieren 3 paarweise Alignments ein multiples Alignment?



Multiples Alignment: Gieriger Algorithmus

Wähle die 2 ähnlichsten* Sequenzen und kombiniere sie in ein Alignment:
 $k \rightarrow k - 1$ Sequenzen/Profile

* besser: die am nächsten verwandten



Profil-Repräsentation eines multiplen Alignments

Alignment:

-	A	G	G	C	T	A	T	C	A	C	C	T	G
T	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	T	C	A	C	-	G	G
C	A	G	-	C	T	A	T	C	G	C	-	G	G

Profil:

A		1				1			.8				
C	.6			1			.4	1		.6	.2		
G			1	.2					.2			.4	1
T	.2				1		.6					.2	
-	.2		.8							.4	.8	.4	

Progressives Alignment: Integriere die Sequenzen schrittweise in das Profil, so dass das Profil nach i Schritten $i+1$ Sequenzen repräsentiert

ClustalW

1. Paarweise Alignments
2. → Führungs-Baum (für Reihenfolge)
3. Progressives Alignment (Gieriger Alg. mit Profilen)

Probleme:

- Gaps im Profil können nicht korrigiert werden
- Profile berücksichtigen Verwandtschaft nicht mehr