

# Algorithmen der Sequenzanalyse: Sequenzalignments

AlgSeq – 09/12/2024

Prof. M. Sammeth

# Multiples Sequenz-Alignement (MSA)

Multipel = mehr als 2 Sequenzen

Asp:

YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTE-FINHYGPTTEATIGA

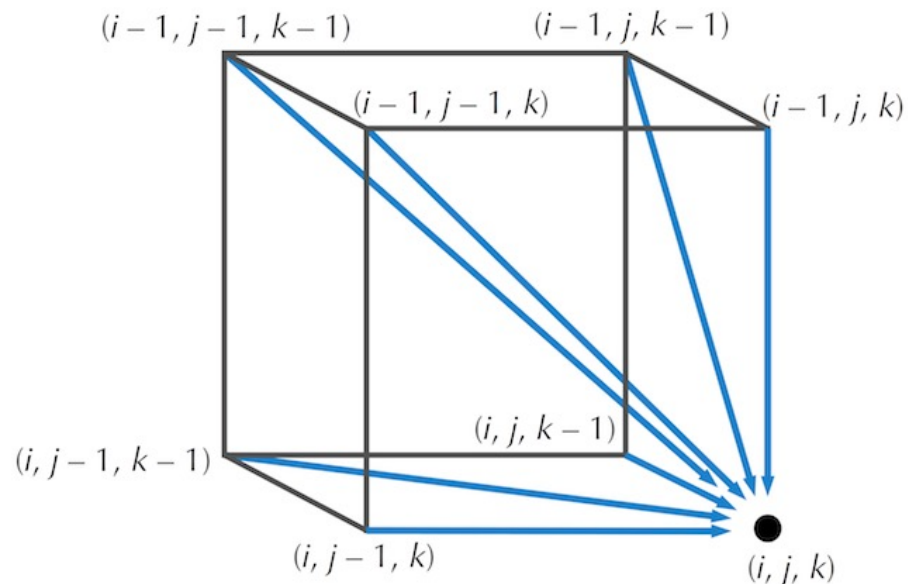
Orn:

-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYI-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

Val:

IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSA-----PTMISSLEILFAAGDRLSSQDAILARRAVSGV-Y-NAYGPTENTVLS

Größe des Alignment Graphen



- 2 Sequenzen mit längen  $m \sim n$ :  
 $n^2$
- 3 Sequenzen:  $n^3$
- k Sequenzen:  $n^k$

Bsp: Proteinfamilie mit 10 Mitgliedern,  
je 500 Aminosäuren.

➔  $500^{10} \sim 10^{27}$  Knoten (!)

# Heuristiken: Progressives Alignment

## Heuristik

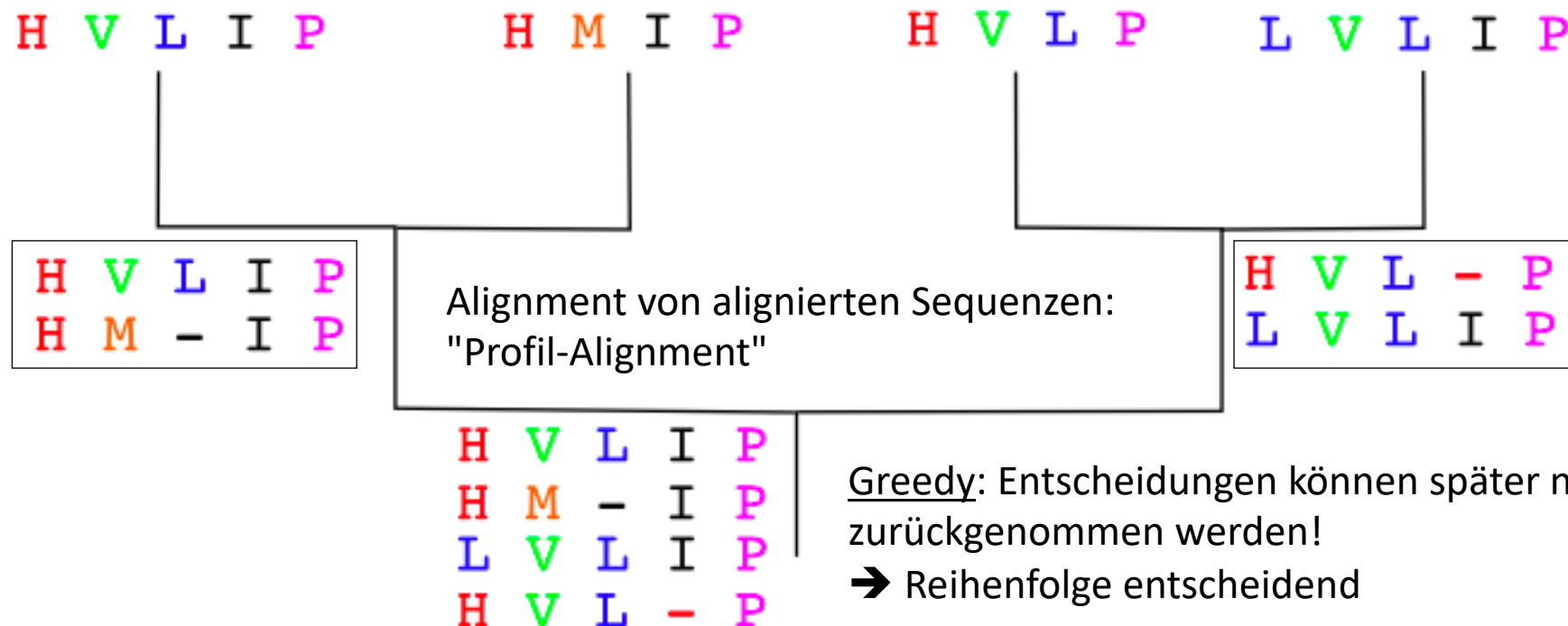
Abkürzung der Berechnung, effizienterer Algorithmus.

Heuristiken sind meist inexakt, es gibt jedoch auch exakte.

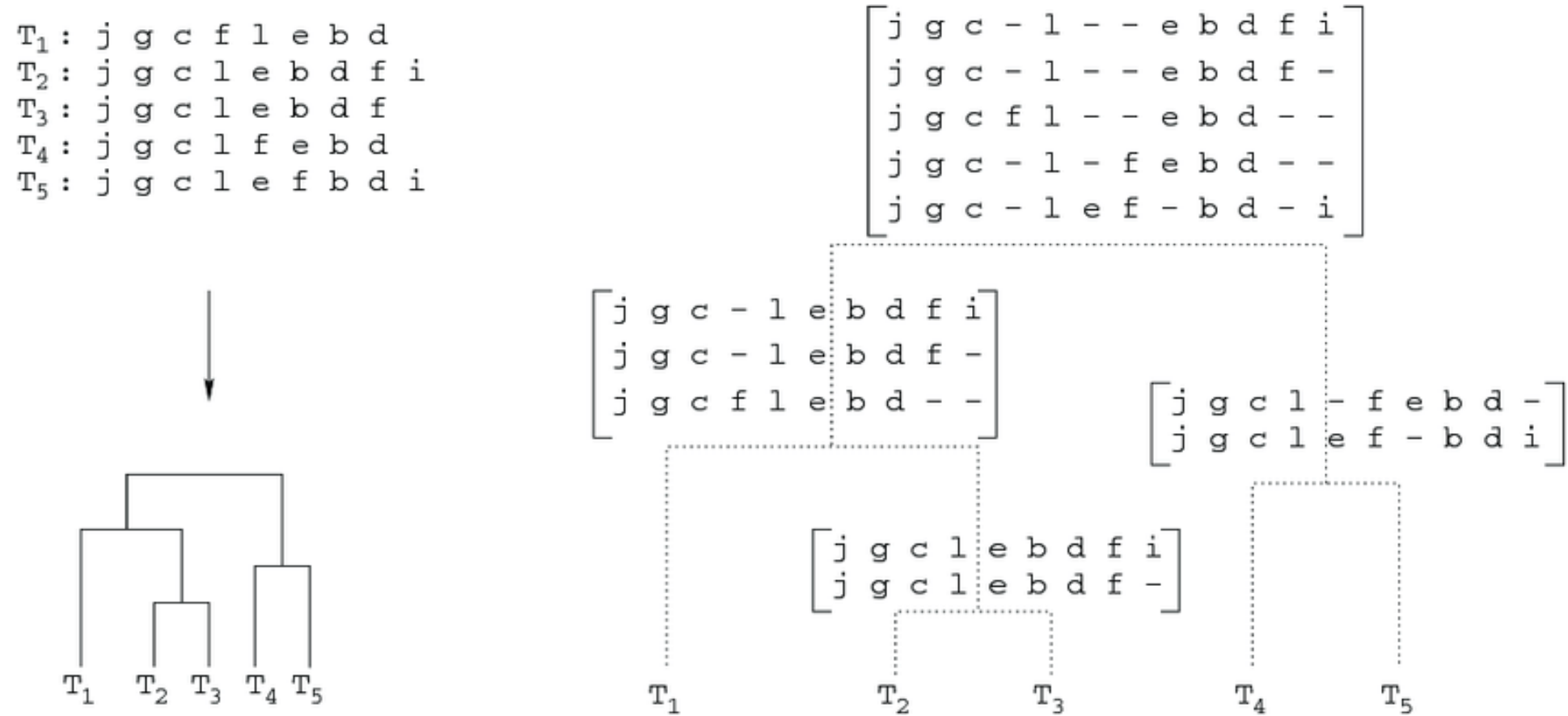
**Progressives Alignment: CLUSTALV, CLUSTALW, CLUSTALX**



*Des Higgins*



# Inexakte Heuristiken: Progressives MSA



Reihenfolge nach Sequenzähnlichkeit (paarweises Alignment)

# Alignment Profile und Profil Alignments

Alignment

|   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|
| T | C | G | G | G | - | g | T | T | T | t | t |
| c | C | - | - | t | G | A | c | T | T | a | C |
| a | C | G | - | G | G | A | T | T | T | t | C |
| T | t | G | G | G | - | A | c | T | T | t | t |
| a | - | - | - | G | - | - | - | T | - | C | - |
| T | t | G | G | G | G | A | c | T | T | C | C |
| T | C | G | - | - | G | A | T | T | c | a | t |
| - | - | - | G | G | G | A | T | T | c | C | - |
| T | a | G | G | G | G | A | a | c | - | - | C |
| T | C | G | G | G | t | A | T | a | a | C | C |

Profil

|    |    |    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| A: | .2 | .1 | 0  | 0  | 0  | 0  | .8 | .1 | .1 | .1 | .2 | 0  |
| C: | .1 | .5 | 0  | 0  | 0  | 0  | 0  | .3 | .1 | .2 | .4 | .5 |
| G: | 0  | 0  | .7 | .6 | .8 | .6 | .1 | 0  | 0  | 0  | 0  | 0  |
| T: | .6 | .2 | 0  | 0  | .1 | .1 | 0  | .5 | .8 | .6 | .2 | .3 |



z.B.  $score(X, v_i) = \sum_{Y \in A} (Profil_i(Y) \frac{score(X, Y)}{|A|}) + (1 - \sum_{Y \in A} (Profil_i(Y))) \cdot score(X, -)$

# ClustalW

Valid format for input is: FASTA(Pearson)  
max number of sequences = 30  
max total length of sequences = 10000

[Help page](#)

More information on Clustal [home page](#)

Scoring matrix :

Opening gap penalty :

End gap penalty :

Output format :

Extending gap penalty :

Separation gap penalty :

Output order :

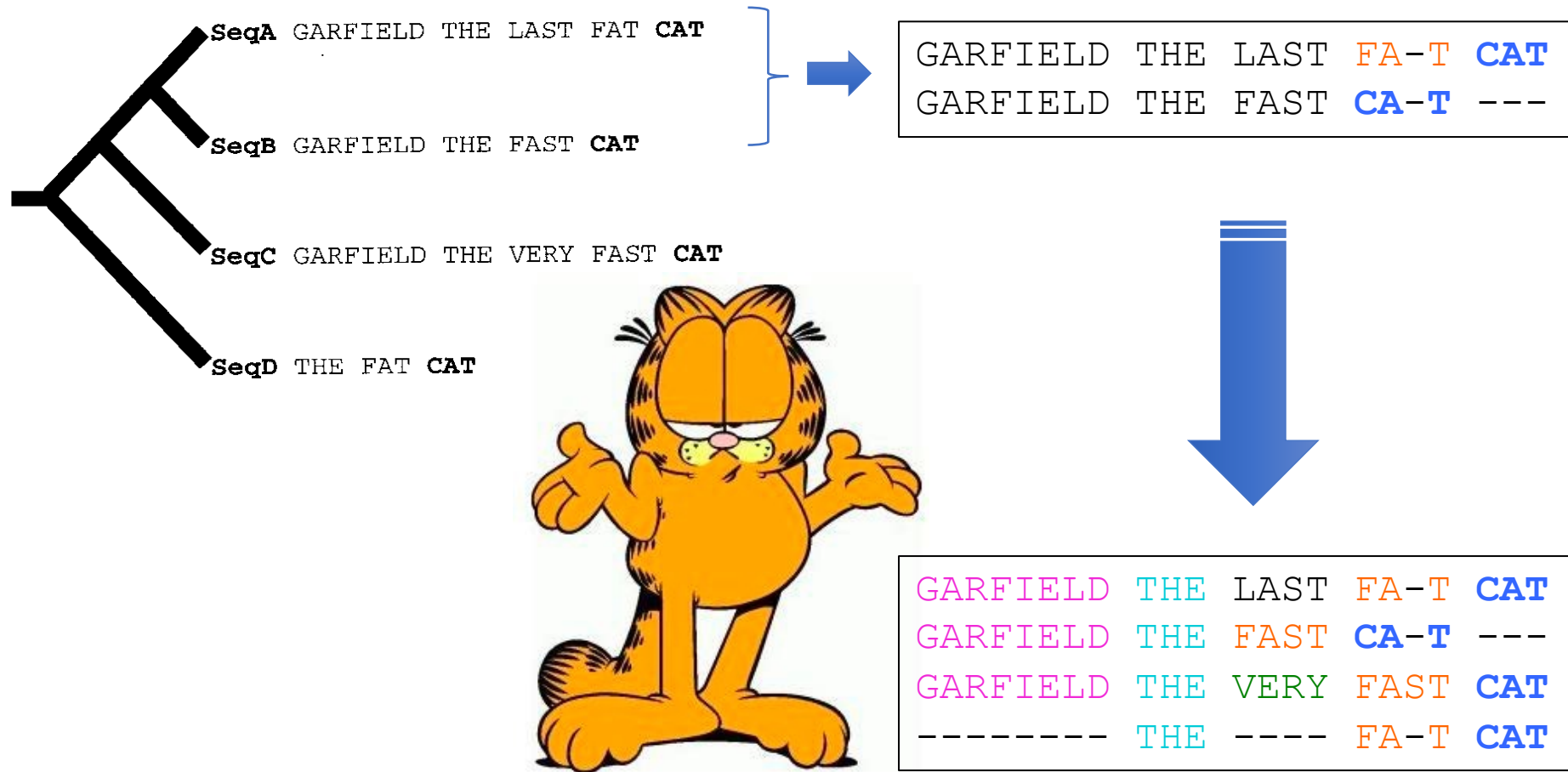
Input sequences:  
(see above for valid  
formats)

Run ClustalW

Clear Input

# Probleme mit dem Progressiven Alignment

*"Once a gap, always a gap"*



# Konsistenz-basiertes Alignment

## (1) Vorbereitung: berechne alle paarweisen Alignments.

SeqA GARFIELD THE LAST FAT CAT Prim. Weight = 88  
SeqB GARFIELD THE FAST CAT ---

SeqB GARFIELD THE ---- FAST CAT Prim Weight = 100  
SeqC GARFIELD THE VERY FAST CAT

SeqA GARFIELD THE LAST FA-T CAT Prim. Weight = 77  
SeqC GARFIELD THE VERY FAST CAT

SeqB GARFIELD THE FAST CAT Prim. Weight = 100  
SeqD ----- THE FA-T CAT

SeqA GARFIELD THE LAST FAT CAT Prim. Weight = 100  
SeqD ----- THE ---- FAT CAT

SeqC GARFIELD THE VERY FAST CAT Prim. Weight = 100  
SeqD ----- THE ---- FA-T CAT

## (2) Gewichtung: wie oft wurde jedes Positionspaar in den paarweisen Vergleichen aligniert?

SeqA GARFIELD THE LAST FAT CAT Weight = 88  
SeqB GARFIELD THE FAST CAT

SeqA GARFIELD THE LAST FAT CAT Weight = 77  
SeqC GARFIELD THE VERY FAST CAT  
SeqB GARFIELD THE FAST CAT

Seq1 GARFIELD THE LAST FAT CAT Weight = 100  
SeqD THE FAT CAT  
SeqB GARFIELD THE FAST CAT

### "Konsistenzen"

SeqA GARFIELD THE LAST FAT CAT  
SeqB GARFIELD THE FAST CAT

## (3) Dynamische Programmierung mit PSSM !

GARFIELD THE LAST FA-T CAT  
GARFIELD THE ---- FAST CAT

## (4) Progressives Multiples Alignment





# PSSM – Position-Specific Scoring Matrices

a) Alignment Matrix

|            |   |   |   |   |   |   |
|------------|---|---|---|---|---|---|
|            | A | A | T | T | G | A |
|            | A | G | G | T | C | C |
|            | A | G | G | A | T | G |
|            | A | G | G | C | G | T |
|            | 1 | 2 | 3 | 4 | 5 | 6 |
| A          | 4 | 1 | 0 | 1 | 0 | 1 |
| C          | 0 | 0 | 0 | 1 | 1 | 1 |
| G          | 0 | 3 | 3 | 0 | 2 | 1 |
| T          | 0 | 0 | 1 | 2 | 1 | 1 |
| consensus: | A | G | G | T | G | N |

$$\ln \frac{(n_{i,j} + p_i) / (N + 1)}{p_i} \approx \ln \frac{f_{i,j}}{p_i}$$

b) Weight Matrix

|                |      |      |      |      |      |   |
|----------------|------|------|------|------|------|---|
|                | 1    | 2    | 3    | 4    | 5    | 6 |
| A              | 1.2  | 0    | -1.6 | 0    | -1.6 | 0 |
| C              | -1.6 | -1.6 | -1.6 | 0    | 0    | 0 |
| G              | -1.6 | .96  | .96  | -1.6 | .59  | 0 |
| T              | -1.6 | -1.6 | 0    | .59  | 0    | 0 |
| test sequence: | A    | G    | G    | T    | G    | C |

# T-Coffee

(Tree-based **Consistency Objective Function**  
for alignment Evaluation)

"Konsistenz"-Scores können beliebige a priori bekannte Informationen in ein Progressives Alignment einfließen lassen.



Cedrik Notredame



## Der COFFEE – Mixer:

**M-Coffee** – kombiniert unterschiedliche multiple alignments

**R-Coffee** – RNA Sekundärstruktur Information

**3D-Coffee** – Protein Strukturinformation

...

# T-Coffee

*Aligns DNA, RNA or Proteins using the default T-Coffee*

## Sequences input

*Paste or upload your set of sequences in FASTA format*

### Sequences to align

[Click here to use the sample file](#)

- OR - [Click here to upload a file](#)

[Hide advanced options](#)

## Methods

*T-Coffee produces an alignment by combining the output of several alignment methods. Use this section to select the individual methods.*

### Pairwise Structural Methods

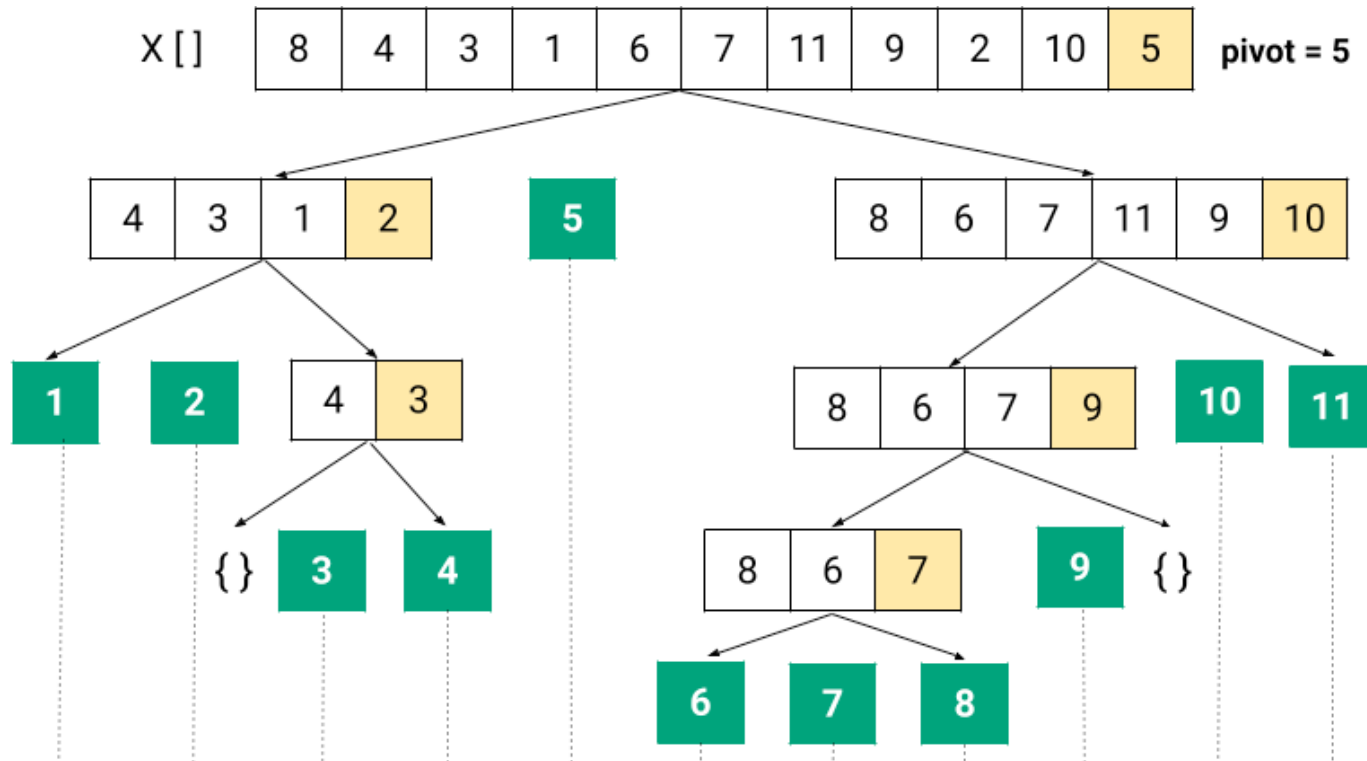
☐ sap\_pair ☐ TAlign\_pair ☐ mustang\_pair

### Multiple Methods

☐ pcma\_msa ☐ mafft\_msa ☐ clustalw\_msa ☐ dialigntx\_msa ☐ poa\_msa  
☐ muscle\_msa ☐ probcons\_msa ☐ t\_coffee\_msa ☐ amap\_msa ☐ kalign\_msa  
☐ fsa\_msa ☐ probconsRNA\_msa ☐ mus4\_msa

# Exakte Heuristiken: Divide-and-Conquer MSA

z.B.: Quicksort, Mergesort, binäre Suche...



### *Divide-and-Conquer Ansätze:*

- zerlegen (rekursiv) den Suchraum in kleinere Sub-Suchräume (*Divide*)
- finden für jeden (meist trivialen) Sub-Sub-Suchraum eine Lösung (*Conquer*)
- setzen die Lösung zur Gesamtlösung zusammen (*Merge*)

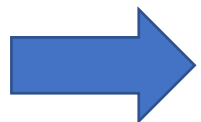
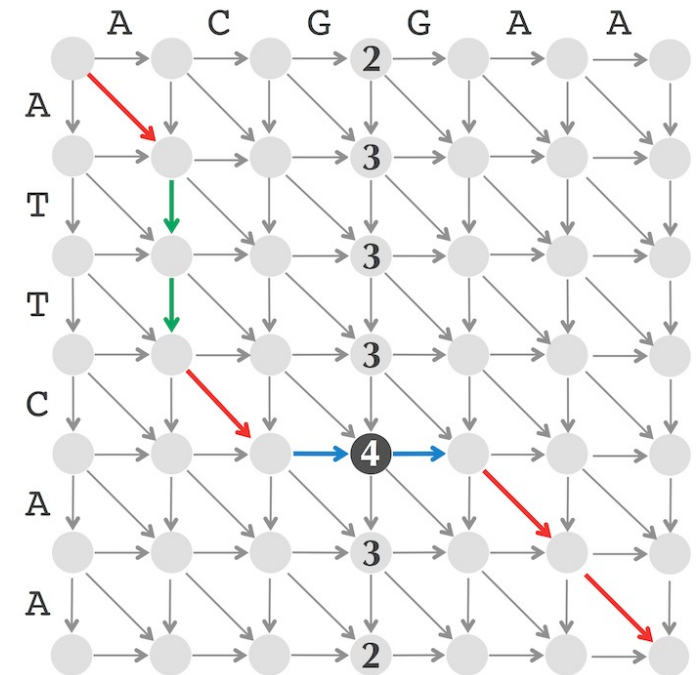
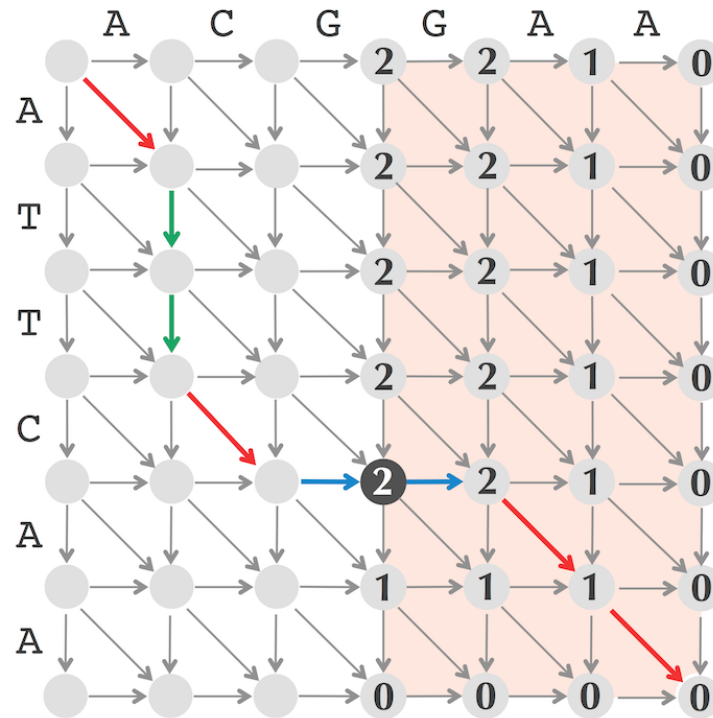
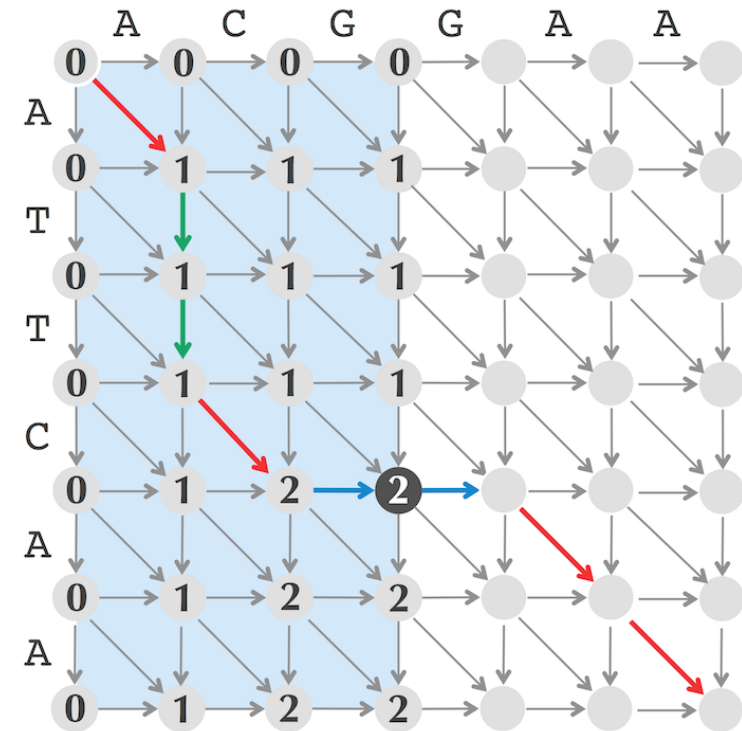


# Divide-and-Conquer für Paarweise Alignments

Optimales *Präfix*-, *Suffix*- und *Totales* Alignment Score

Bsp.: LCS-Scoring (Match= +1)

optimaler Präfix-Alignment Score  $(i,j)$  + optimales Suffix-Alignment Score  $(i,j)$  = Score des optimalen Alignments, d.h. längster Pfad durch  $(i,j)$

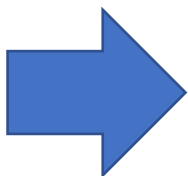
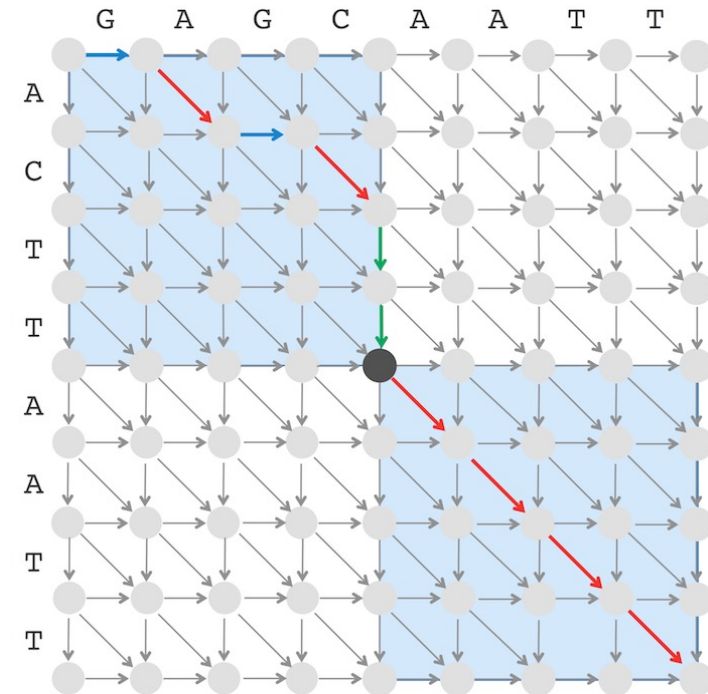
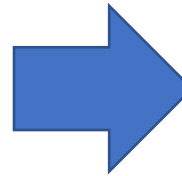
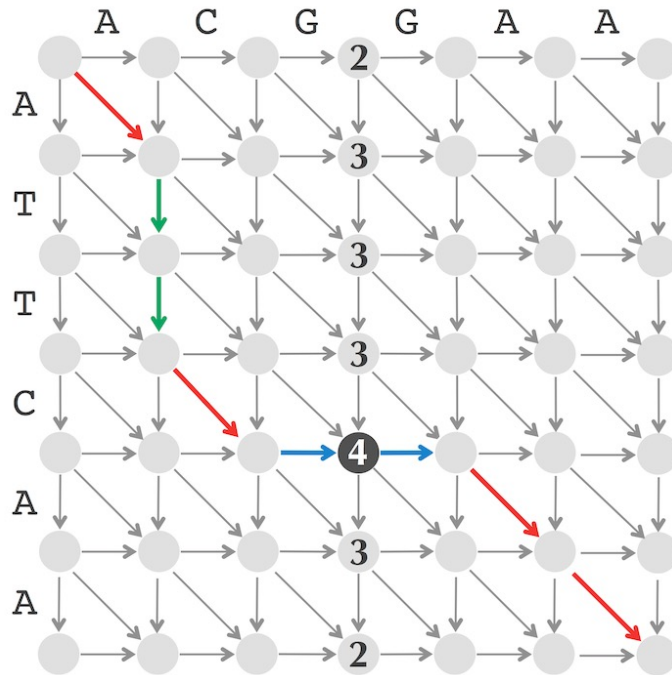


$$\text{Pfadlänge}(i,j) = \text{vonSource}(i,j) + \text{zurSink}(i,j)$$

# Der “Mittlere” Knoten

*Mittlere Spalte* := Spalte, die die Dynamische Programmier-Matrix halbiert

*Mittlerer Knoten* := (ein) Knoten in der mittleren Spalte des Alignments mit maximalem Score.



Ein Mittlerer Knoten  $(i,j)$  teilt die DP-Matrix in die zwei Submatrizen  $[0;i] \times [0;j]$  und  $[i+1;n] \times [j+1;m]$

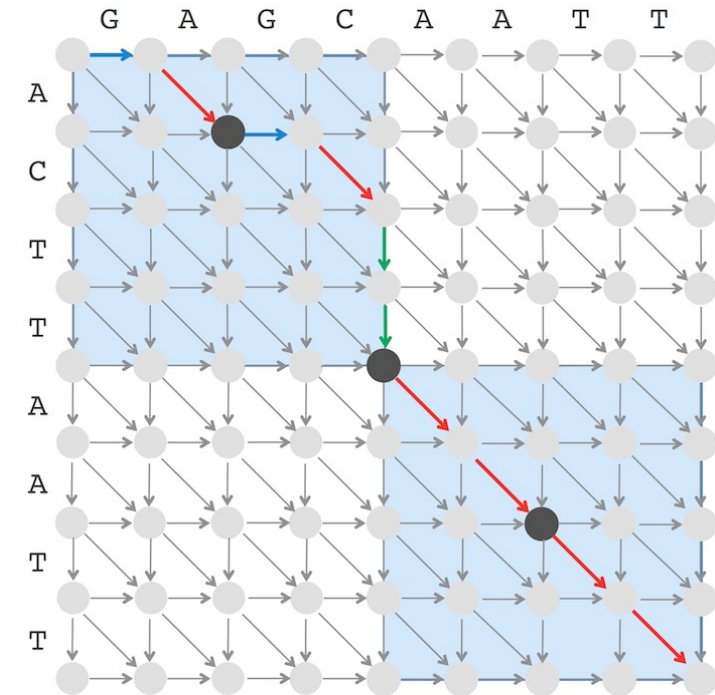
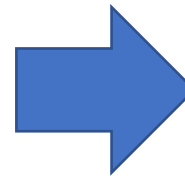
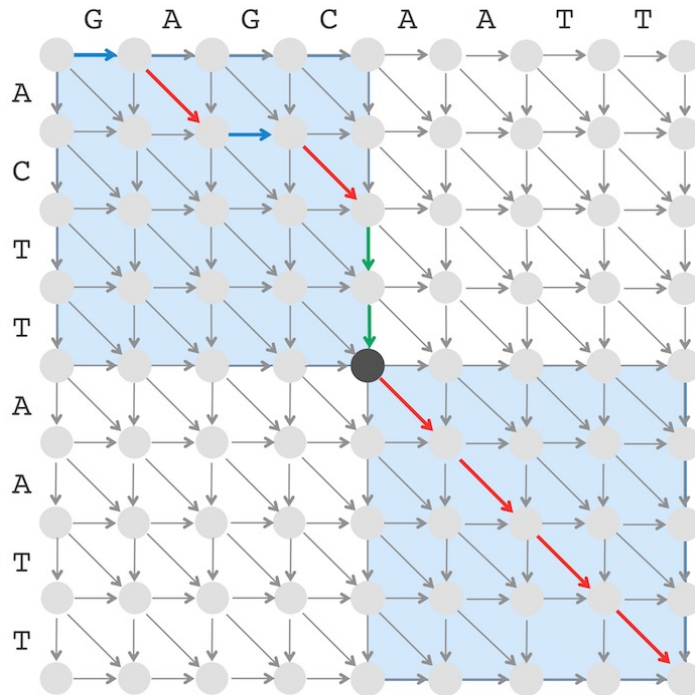


# Divide-and-Conquer mit mittleren Knoten

Idee:

- Der mittlere Knoten teilt das Problem in zwei +/- gleich große Subprobleme (*Divide*).
- Wenn links/oben und rechts/unten vom *mittleren Knoten* alle Knoten des längsten Pfades bekannt sind (*Conquer*),
- dann kann die Lösung einfach zusammengesetzt werden (*Merge*) als:

Knoten des längsten Pfades links/oben + mittlerer Knoten + Knoten des längsten Pfades rechts/unten

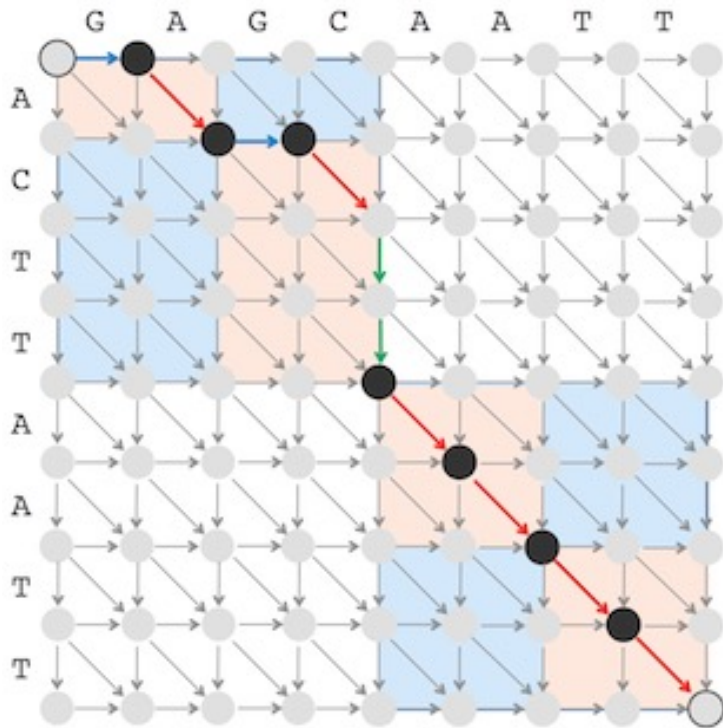


erster mittlerer Knoten: teilt DP Matrix in zwei Submatrizen

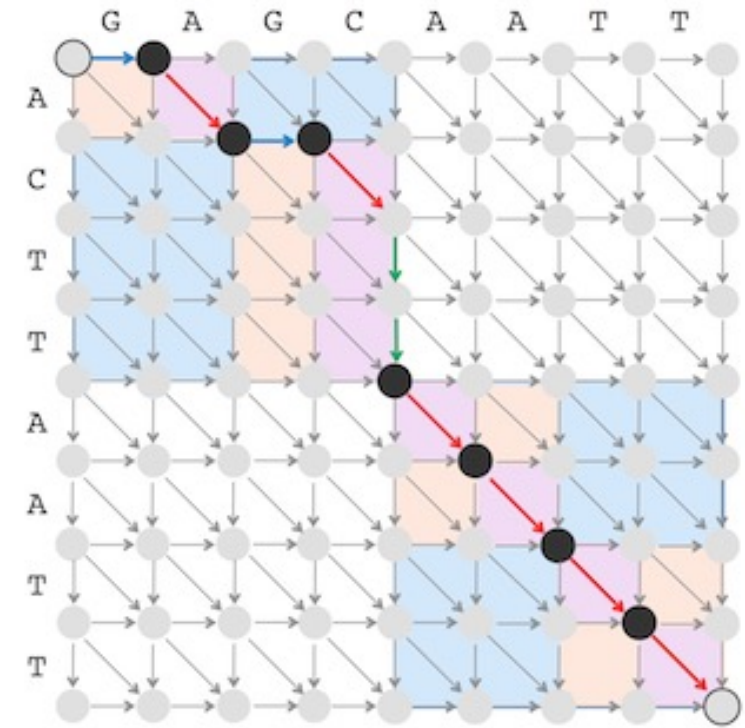
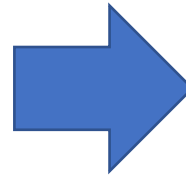
zweite mittlere Knoten, je einer pro Submatrix

# Paarweises Divide-and-Conquer Alignment

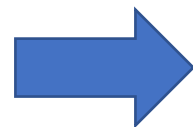
Rekursive Fortsetzung führt zu



dritte mittlere Knoten



(fast) alle Knoten des längsten Pfades bekannt

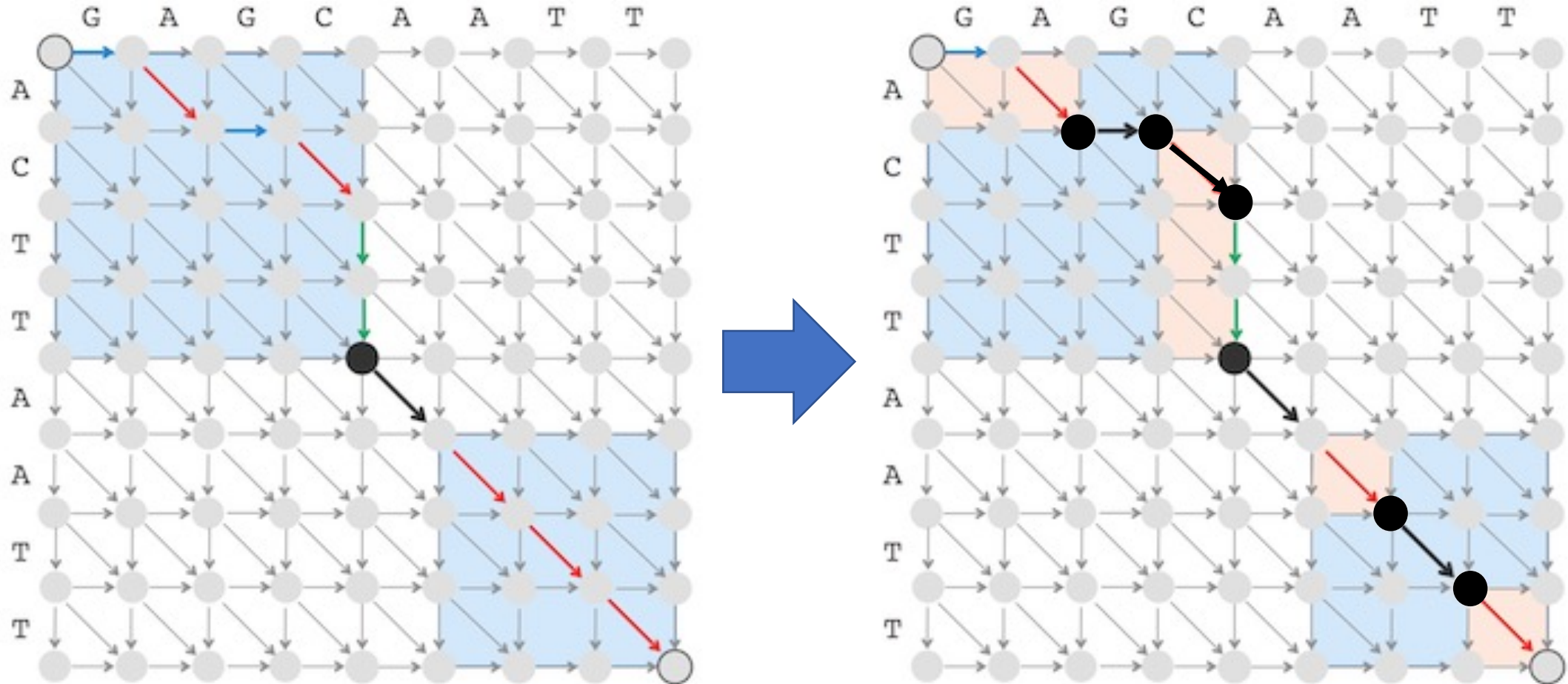


paarweises Alignment in  $O(n)$  Speicher möglich!



# Mittlere Kanten statt Mittlere Knoten

*Mittlere Kante* := ausgehende Kante des/eines mittleren Knotens, die Teil des/eines optimalen Alignments ist.



Rekursive Teilung über "mittlere Kanten" rekonstruiert intrinsisch den optimalen Alignement Pfad.

# Algorithmus Outline:

---

**Algorithm:** DCAALIGNMENT(*oben*, *unten*, *links*, *rechts*)

---

```
if links = rechts then
  | return Alignment aus (unten – oben) “↓” Kanten
if oben = unten then
  | return Alignment aus (rechts – links) “→” Kanten
mittlereSpalte =  $\lfloor (links + rechts)/2 \rfloor$ 
mittlereZeile = MITTLERERKNOTEN(oben, unten, links, rechts)
mittlereKante = MITTLEREKANTE(oben, unten, links, rechts)
output DCAALIGNMENT(oben, mittlereZeile, links, mittlereSpalte)
output mittlereKante
if mittlereKante = “→” or mittlereKante = “↘” then
  | mittlereSpalte = mittlereSpalte + 1
if mittlereKante = “↓” or mittlereKante = “↙” then
  | mittlereZeile = mittlereZeile + 1
output DCAALIGNMENT(mittlereZeile, unten, mittlereSpalte, rechts)
```

---

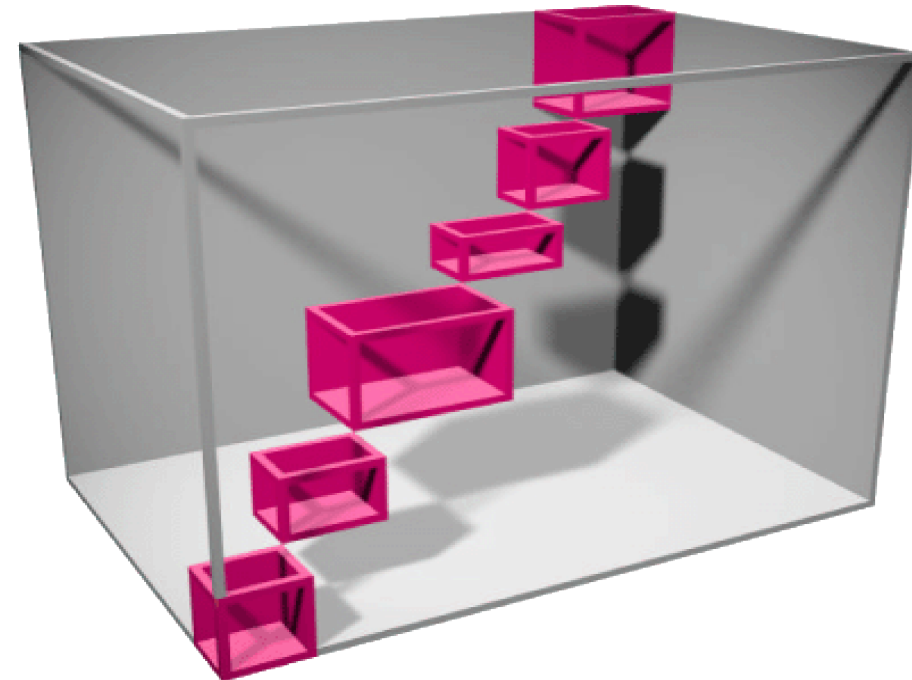
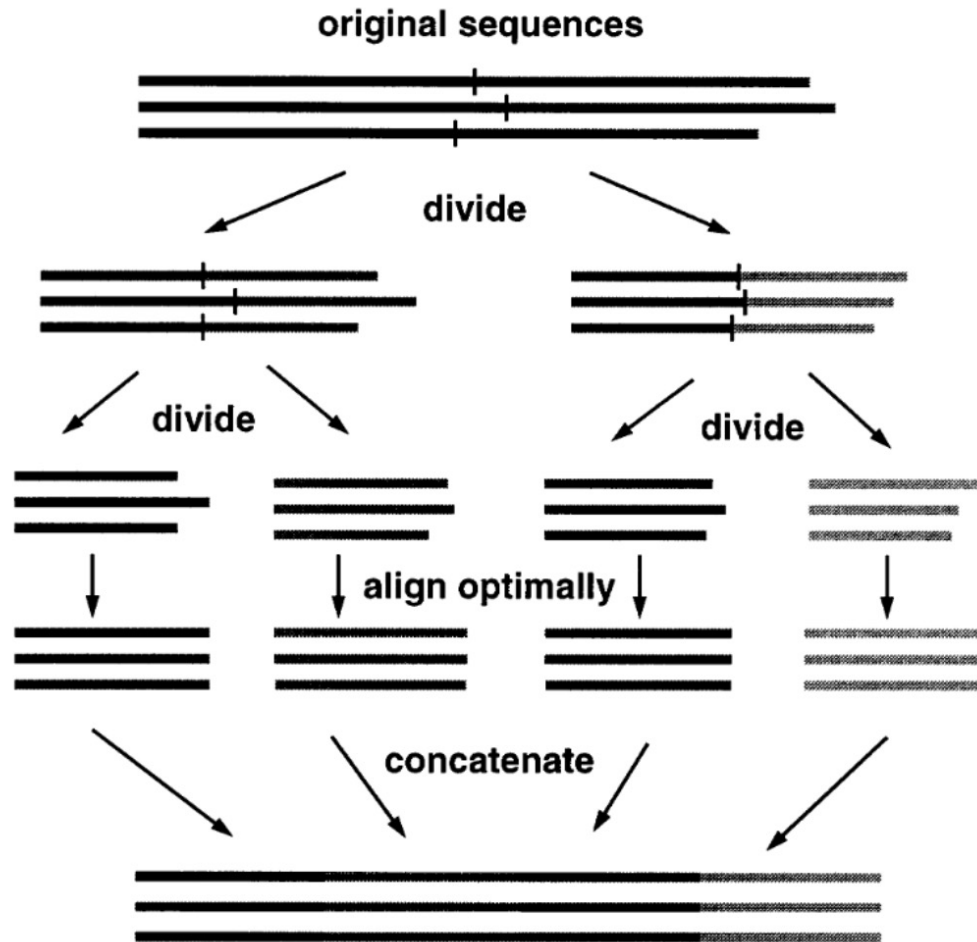
MITTLERERKNOTEN(...): gibt Koordinate *i* des mittleren Knotens der Substrings  $v_{oben+1} \dots v_{unten}$  und  $w_{links+1} \dots w_{rechts}$  zurück

MITTLEREKANTE(...): gibt die Richtung der ausgehenden Kante des mittleren Knotens (*i*, *mittlereSpalte*) zurück

# Divide-and-Conquer für MSA



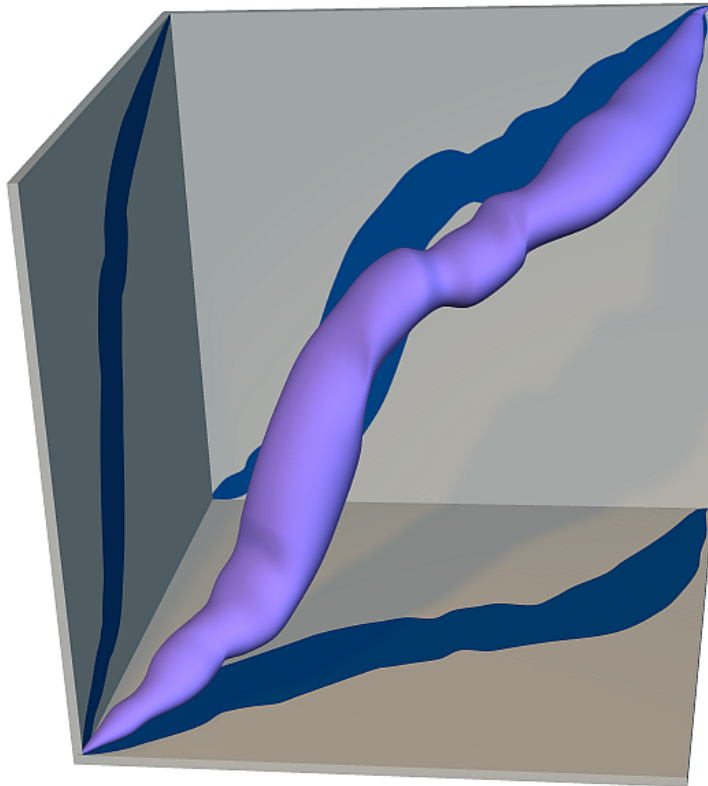
Jens Stoye



Stoye 1998, *Multiple sequence alignment with the divide-and-conquer method*.

[https://doi.org/10.1016/S0378-1119\(98\)00097-3](https://doi.org/10.1016/S0378-1119(98)00097-3) Bielefeld

# Exakte MSA Heuristiken: Branch-and-Bound



*Branch-and-Bound* Optimierung:

- exploriert systematisch den Ergebnisraum (*branching step*)
- berücksichtigt aber suboptimale Lösungen frühzeitig nicht weiter (*bounding step*)
- benötigt Annahme über das Score des optimalen Alignments und monotone Scoring Funktion.
- im *worst case* muss der ganze Datenraum exploriert werden.