

# Problemseminar

Multiple alignment by sequence annealing

Ariel S. Schwartz and Lior Pachter

## 0.Gliederung

### 1.Alignment posets

# Gliederung

1. Einführung
2. Alignment Posets
3. Sequence Annealing
4. Results

# 1. Einführung

## 2. Partial order graphs

# Progressives Alignment

- Multiples Alignment durch paarweises Alignment



1. Einführung
- 2. Partial order graphs**
3. Sequence Annealing

# Partial order graphs

1. Einführung

2. Partial order graphs

3. Sequence Annealing

## Alignment-Problem:

$S_1$ :                    . . . . . ACATGTCGAT . . . . . AGGTG  
 $S_2$ :                    TGCAC . . . . . TCGATACATAAGGTG

Consensus1: TGCACACATGTCGATACATAAGGTG

Consensus2: ACATGTGCACTCGATACATAAGGTG

. . .

1. Einführung

2. Partial order graphs

3. Sequence Annealing

## Neue Forderungen:

- Alignment durch paarweises dynamisches Alignment möglich (wie bisher auch bei Konsensussequenz)
- Kein Informationsverlust gegenüber dem MSA (Multiple Sequence Alignment)

1. Einführung

**2. Partial order graphs**

3. Sequence Annealing

- Welche Sequenzpositionen sind zueinander aligned?
- Ordnung dieser Positionen in den Sequenzen?

→ Partial Order Graph

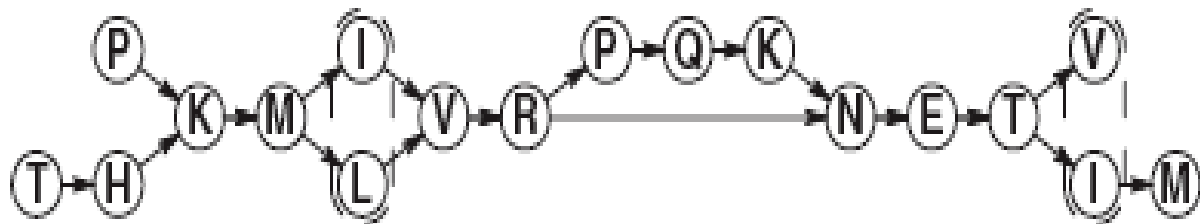
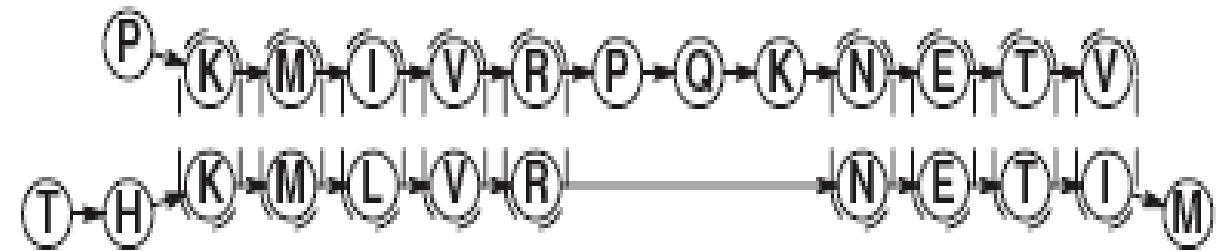


1. Einführung

## 2. Partial order graphs

3. Sequence Annealing

. . P K M I V R P Q K N E T V .  
T H . K M L V R . . . N E T I M



1. Einführung

2. Partial order graphs

3. Sequence Annealing

- Welche Sequenzpositionen sind zueinander aligned? ✓
- Ordnung dieser Positionen in den Sequenzen? ✓

2. Partial order graphs

**3. Sequence Annealing**

4. Results

# Sequence annealing

2. Partial order graphs

3. Sequence Annealing

4. Results

Ausgangsposition:

- $k$  Sequenzen der Länge  $n_1, n_2, \dots, n_k$
- $\sigma_i^a$  -  $i$ -te Position in der Sequenz  $a$
- $L = \sum_{(i)} n_i$

2. Partial order graphs

### 3. Sequence Annealing

4. Results

- $M_{\text{null}}$  – Null-Alignment

- Länge  $L$

- Bsp: 

ABCD	-----	-----
-----	EFGH	-----
-----	-----	IJKL

2. Partial order graphs

**3. Sequence Annealing**

4. Results

## Idee

- Alle möglichen Paare bekommen Gewichte (posterior Wahrscheinlichkeiten für matches und gaps)
- Positionen, die am wahrscheinlichsten homolog sind, werden zuerst aligned

2. Partial order graphs

3. Sequence Annealing

4. Results

## Zielfunktion

- Exaktheit des multiplen Sequence Alignment muss bestimmt werden
- *Developer score ( $f_D$ )* (*~sum-of-pair-score*)
- *Alignment metric accuracy (AMA)*
- Zielfunktion, um sowohl  $f_D$  als auch *AMA* zu optimieren

2. Partial order graphs

3. Sequence Annealing

4. Results

- probabilistisches Modell muss gegeben sein

- *Match posterior probability*

$$P(\sigma_i^1 \diamond \sigma_j^2 | \sigma^1, \sigma^2, \theta)$$

- *Gap posterior probability*

$$P(\sigma_i^1 \diamond \text{---} | \sigma^1, \sigma^2, \theta)$$

$$P(\text{---} \diamond \sigma_j^2 | \sigma^1, \sigma^2, \theta)$$

- Gap-Faktor:  $G_f$



2. Partial order graphs

3. Sequence Annealing

4. Results

Resultierende Zielfunktion  $f$

$$\begin{aligned} f^{G_f}(M) = & \sum_{\sigma^a, \sigma^b \text{ mit } a \neq b} \left( \sum_{\{(j,k) | \varphi^M(\sigma_j^a) = \varphi^M(\sigma_k^b)\}} P(\sigma_j^a \diamond \sigma_k^b | \sigma^a, \sigma^b, \theta) \right. \\ & + G_f \cdot \sum_{\{j | \forall \sigma_k^b \varphi^M(\sigma_j^a) \neq \varphi^M(\sigma_k^b)\}} P(\sigma_j^a \diamond \neg | \sigma^a, \sigma^b, \theta) \\ & \left. + G_f \cdot \sum_{\{k | \forall \sigma_j^a \varphi^M(\sigma_j^a) \neq \varphi^M(\sigma_k^b)\}} P(\neg \diamond \sigma_k^b | \sigma^a, \sigma^b, \theta) \right) \end{aligned}$$

2. Partial order graphs

3. Sequence Annealing

4. Results

$G_f=0$ :  $f^0$  berechnet  $f_D$ -score

$G_f=0.5$ :  $f^{0.5}$  berechnet  $AMA$ -score

je größer  $G_f$  umso höher die Spezifität

je kleiner  $G_f$  umso höher die Sensitivität

2. Partial order graphs

3. Sequence Annealing

4. Results

Gewichtsfunktion –  $w(p)$

- jedem Paar  $p=(c_k, c_i)$  soll ein Gewicht zugeordnet werden
- $f(M_{i-1}) \geq f(M_i)$  bei positiven Gewichten

2. Partial order graphs

**3. Sequence Annealing**

4. Results

- Aus  $f$  werden die Gewichtsfunktionen berechnet
  - statische vs. dynamische Version

2. Partial order graphs

3. Sequence Annealing

4. Results

- $P_{match}(c_k, c_l)$  - Wahrscheinlichkeit, dass  $c_k$  und  $c_l$  matchen
- $P_{gap}(c_k, c_l)$  – Wahrscheinlichkeit, dass  $c_k$  und  $c_l$  eine Lücke bilden
- Gap-Faktor  $G_f$

2. Partial order graphs

3. Sequence Annealing

4. Results

1. Gewichtungsfunktion: maxstep:

$$W_{\maxstep}^{G_f}(c_k, c_l) = \frac{P_{\text{match}} - G_f \cdot P_{\text{gap}}}{|\varphi^{-1}(c_k) \varphi^{-1}(c_l)|}, \text{ falls } c_k \neq c_l, \text{ sonst } -\infty$$

2. Partial order graphs

3. Sequence Annealing

4. Results

## 2. Gewichtungsfunktion: $tgf$

$$w_{tgf}^{G_f}(c_k, c_l) = \frac{P_{match}}{P_{gap}} - G_f, \text{ falls } c_k \neq c_l, \text{ sonst } -\infty$$

$$temp = w_{tgf}^{G_f}(p_{max}) + G_f$$

$$f^{temp}(M)$$

→dynamische Gewichte

2. Partial order graphs

3. Sequence Annealing

4. Results

## Algorithmus

1)  $M_i = M_{Null}$

2) **while**  $\exists$  Paar  $p$  mit  $w(p) > 0$   
    **do** Merge  $p_{max}$  zu  $M_{i-1}$   
    *i--*

3) **end while**



### 3. Sequence Annealing

## 4. Results

- AMAP 2.0 - <http://bio.math.berkeley.edu/amap/>
- Laufzeit vergleichbar mit anderen Programmen (Dialign-T, MUSCLE, T-Coffee)
- Methoden vermutlich auch für DNA anwendbar

### 3. Sequence Annealing

## 4. Results

- Benchmark-Test auf Proteinsequenzen
- 3 verschiedene scores wurden berechnet (Spezifität, Sensibilität, Mix)
- In fast allen Tests bestes Programm
- Zu jedem Zeitpunkt des Algorithmus' kann ein Alignment ausgegeben werden(angefangen von sehr spezifischen Alignments mit geringer Sensibilität bis hin zu sehr sensiblen Alignments mit geringerer Spezifität)

# Ende

