

# Multiple alignment by sequence annealing

## Alignment Poset

$\sigma_i^a$  ...  $i$ -te Element einer Sequenz  $\sigma^a = \sigma_1^a, \dots, \sigma_n^a$  der Länge  $n$

$S = \{\sigma^1, \dots, \sigma^k\}$  ... Menge der  $n_1 + n_2 + \dots + n_k$  Sequenzbuchstaben der Sequenzen  $\sigma^1, \sigma^2, \dots, \sigma^k$  der Längen  $n_1, n_2, \dots, n_k$

$P = \{c_1, \dots, c_m\}$  ... partiell geordnete Menge der Sequenzbuchstaben  $S$  zusammen mit der surjektiven Funktion  $\varphi: S \rightarrow P$ , so dass  $\varphi(\sigma_i^a) < \varphi(\sigma_j^a)$ , falls  $i < j$   
 $\rightarrow$  *partielles globales multiples Alignment*

Die Elemente von  $P$  korrespondieren also zu den Spalten des multiplen Alignments und die partielle Ordnung gibt die Reihenfolge an, in der die Spalten auftreten müssen.

Beispiel:

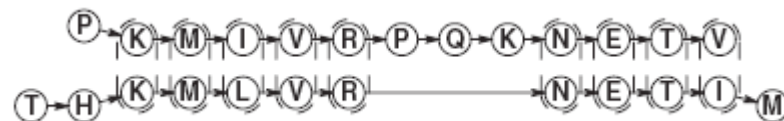
```

. . 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
    
```

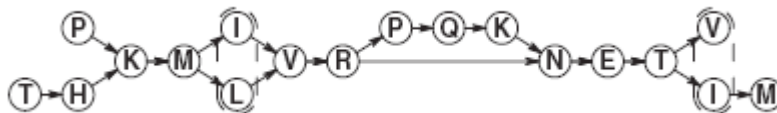
2 Sequenzen



Sequenz 1 als gerichteter azyklischer Graph



Markierung der Alignment-Knoten



PO-MSA Repräsentation eines paarweisen Protein-Sequenz-Alignments

## Scoring-Funktion für multiple Alignments

$M$  ... Menge aller partiellen multiplen Alignments einer Menge Sequenzbuchstaben  $S$

$f: M \rightarrow \mathbb{R}$  ... Funktion, die jedem multiplen Alignment einen „score“ zuordnet

$f_D$  (Developer score) – equivalent zum sum-of-pairs score

Verhältnis der Anzahl der korrekt aligneden Paare im Alignment zur Anzahl der aligneden Paare im Referenzalignment

*alignment metric accuracy (AMA)* ... Verhältnis von Resten, die korrekt zu einem anderen Rest oder zu einer Gap aligned sind (Schwartz et al, 2006)

$h^p$  – predicted Alignment

$h^r$  – reference Alignment

$h^i, h^j \in A_{n_1, n_2, \dots, n_k}$  - MSAs mit  $k$  Sequenzen der Längen  $n_1$  bis  $n_k$

$$g(h^p, h^r) = 1 - \frac{d(h^p, h^r)}{(k-1) \sum_{i=1}^k n_i} \quad \text{mit} \quad d(h^i, h^j) = \sum_{s^1=1}^{k-1} \sum_{s^2>s^1}^k d(h_{s^1, s^2}^i, h_{s^1, s^2}^j)$$

$$\begin{aligned}
d(h^i, h^j) &= 2|h_H^i| + |h_I^i| + |h_D^i| - 2|h_H^i \cap h_H^j| - |h_I^i \cap h_I^j| - |h_D^i \cap h_D^j| \\
&= 2|h_H^i| + |h_I^i| + |h_D^i| - 2|h_H^i \cap h_H^j| - |h_I^i \cap h_I^j| - |h_D^i \cap h_D^j| \\
&= n + m - 2|h_H^i \cap h_H^j| - |h_I^i \cap h_I^j| - |h_D^i \cap h_D^j| \\
h_H &= \{(i, j) : (\sigma_i^1 \diamond \sigma_j^2) \in h\} \text{ (Align)} \\
h_D &= \{i : (\sigma_i^1 \diamond -) \in h\} \text{ (Deletion)} \\
h_I &= \{j : (\text{sigma}_j^2 \diamond -) \in h\} \text{ (Insertion)}
\end{aligned}$$

## Sequence Annealing

Ziel:  $\text{argmax}_{M \in \mathcal{M}} (f(M))$  finden (wobei  $M$  über alle partiellen multiplen Alignments geht)

$$L = \sum_i n_i \quad \dots \text{ Länge aller Sequenzen}$$

$M_{\text{null}}$  ... Null globales multiples Alignment von  $k$  Sequenzen der Längen  $n_1, \dots, n_k$   
(disjunkte Vereinigung von  $k$  Ketten)

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

Wahrscheinlichkeiten:

- probabilistisches Modell muss gegeben sein
- Match posterior probability:  $P(\sigma_i^1 \diamond \sigma_j^2 | \sigma^1, \sigma^2, \theta)$
- Gap posterior probability:  $P(\sigma_i^1 \diamond - | \sigma^1, \sigma^2, \theta)$   
 $P(- \diamond \sigma_j^2 | \sigma^1, \sigma^2, \theta)$

- Gap-Faktor:  $G_f$

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

- $G_f=0$ :  $f^0$  berechnet  $f_D$ -score,  $G_f=0.5$ :  $f^{0.5}$  berechnet AMA-score
- allgemein: je größer  $G_f$ , umso höher die Spezifität; je kleiner  $G_f$ , umso höher die Sensitivität
- Gewichtsfunktion:  $w(p)$

- jedem Paar  $p=(c_k, c_l)$  soll ein Gewicht zugeordnet werden
- bei positiven Gewichten soll im Falle des Alignments dieses Paares das entstehende multiple Alignment einen höheren score haben als vorher:  $f(M_{i-1}) \geq f(M_i)$
- Gewichtsfunktionen werden aus  $f$  berechnet

- statisch: einmal am Anfang des Algorithmus' berechnet
- dynamisch: für jedes Kandidatenpaar  $p$  neu berechnet

- das Paar mit dem höchsten Gewicht wird jeweils gemerged
- $P_{\text{match}}(c_k, c_l)$  - Wahrscheinlichkeit, dass  $c_k$  und  $c_l$  matchen

- $$P_{\text{match}} = \sum_{\{\sigma_i^a \in \varphi^{-1}(c_k)\}} \sum_{\{\sigma_j^b \in \varphi^{-1}(c_l)\}} P(\sigma_i^a \diamond \sigma_j^b | \sigma^a, \sigma^b, \theta)$$

- $P_{\text{gap}}(c_k, c_l)$  - Wahrscheinlichkeit, dass  $c_k$  und  $c_l$  eine Lücke bilden

- $$P_{\text{gap}} = \sum_{\{\sigma_i^a \in \varphi^{-1}(c_k)\}} P(\sigma_i^a \diamond - | \sigma^a, \sigma^b, \theta) + \sum_{\{\sigma_j^b \in \varphi^{-1}(c_l)\}} P(- \diamond \sigma_j^b | \sigma^a, \sigma^b, \theta)$$

- Gewichtungsfunktion: *maxstep*
  - $w_{maxstep}^{G_f}(c_k, c_l) = \frac{P_{match} - G_f \cdot P_{gap}}{|\Phi^{-1}(c_k) \Phi^{-1}(c_l)|}$ , falls  $c_k \neq c_l$ , sonst  $-\infty$
  - Anstieg in der Scoring-Funktion  $f^{G_f}$  soll maximal sein
  
- Gewichtungsfunktion: *tgf*
  - $w_{tgf}^{G_f}(c_k, c_l) = \frac{P_{match}}{P_{gap}} - G_f$ , falls  $c_k \neq c_l$ , sonst  $-\infty$
  - $temp = w_{tgf}^{G_f}(p_{max}) + G_f$
  - scoring-Funktion  $f^{temp}(M)$  kann neu berechnet werden
  - Gewichte werden somit dynamisch neu berechnet

### Algorithmus:

1.  $M_L \leftarrow M_{null}$
2.  $i \leftarrow L$
3. **while**  $\exists c_k^{M_i}, c_l^{M_i}$ , so dass  $c_k^{M_i}$  und  $c_l^{M_i}$  zu  $M'$  gemerged wird und  $f(M') \geq f(M_i)$  **do**
4.  $M_{i-1} \leftarrow M'$
5.  $i \leftarrow i-1$
6. **end while**

vereinfacht:

„Solange es in dem aktuellen partiellen multiplen Alignment noch ein Buchstabenpaar gibt, dass den score des Alignment erhöht würde, merge diese.“

**Programm:** <http://baboon.math.berkeley.edu/amap/>

### Quellen

A.S. Schwartz and L. Pachter, *Multiple Alignment by Sequence Annealing*, Bioinformatics 23 (2007), e24–e29.

A. S. Schwartz, E. W. Myers and L. Pachter, *Alignment Metric Accuracy*, arXiv:q-bio/0510052