

A Variation on Algorithms for Pairwise Global Alignments

Ulrike Mueckstein

Institute for Theoretical Chemistry and Structural Biology

University Vienna

<http://www.tbi.univie.ac.at/~ulim/>

Bled, 2002

Various Alignments

```

CE      -----QAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSEVPQNN-PELQAHAGKVFKLVYE
-----DKTNVKAAWGKVGGAHAGEYGAELERMFLSFPTTKTYF--PHFDLSHGSAQVKGHGKKVADALTNVA

TOP 6.7 --LTE-QAALVKSSWEEFN-----HTHRFFILVLE-APAAK-----HAGK-----
--LSP-DKTNVKAAWGKVG-----YGAEALERMFL-FPTTK-----KVAD-----

SARF2  --LTSQAALVKSSWEEFNANI-KHTRFFILVLEIAPAAKDLF----KGTSEVP--NPELQAHAGKVFKLVYE
--LSPADKTNVKAAWGKVGAAH-EYGAEALERMFLSFPTTKTYF----FDLSHGS--K-GHGKKVADALTNVA

MATRAS -ALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSEVPQNNPELQAHAGKVFKLVYEA
-VLSPADKTNVKAAWGKVGGAHAGEYGAELERMFLSFPTTKTYFPHFD----LSHGSAQVKGHGKKVADALTNA

1GDJ   GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSEVPQNNPELQAHAGKVFKLVYEA
2HHB A -VLSPADKTNVKAAWGKVGGAHAGEYGAELERMFLSFPTTKTYFPHF----DLSHGSAQVKGHGKKVADALTNA
  
```

alignment reliability

| | |
|------|-----|
| high | low |
|------|-----|

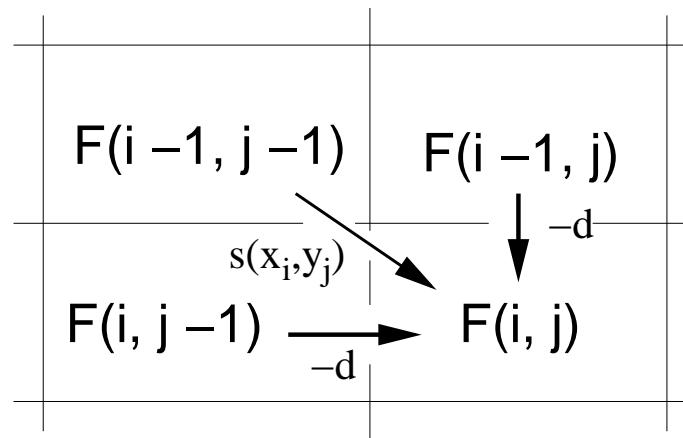
alignment between *Lupinus luteus* leghaemoglobin (1GDJ)
and human chain A deoxyhaemoglobin (2HHB A)

Needleman-Wunsch Algorithm

- ★ dynamic programming algorithm: an optimal alignment is build recursively, using previous solutions for optimal alignments for smaller subsequences.

| | | |
|-------------------|--------------------|--------------------|
| XXXx _i | XXXx _i | XXx _i - |
| YYYy _j | YYy _j - | YYYy _j |

- ★ The score of the best alignment between the initial segments x_1, x_2, \dots, x_i and y_1, y_2, \dots, y_j is stored in matrix F at position $F(i, j)$.



Recursion

- match (x_i, y_j) :

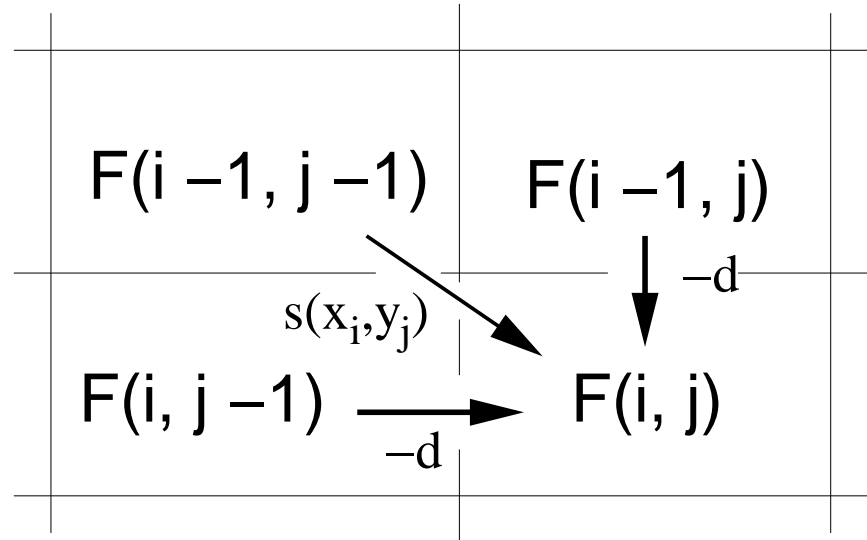
$$F(i, j) = F(i - 1, j - 1) + s(x_i, y_i)$$

- gap in sequence \mathbf{x} $(-, y_j)$:

$$F(i, j) = F(i, j - 1) - d$$

- gap in sequence \mathbf{y} $(x_i, -)$:

$$F(i, j) = F(i - 1, j) - d$$



$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_i) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

★ Initialisation

$$F(0, 0) = 0; \quad F(i, 0) = -id; \quad F(0, j) = -jd;$$

★ Extension

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_i) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

★ Backtracking

$F(n, m)$, is by **definition** the best score for an alignment of x_1, x_2, \dots, x_n and y_1, y_2, \dots, y_m .

$$s(x, y) = 2 \text{ if } x = y;$$

$$s(x, y) = 0 \text{ if } x \neq y; d = 2$$

| | | | | | |
|---|----|----|----|----|----|
| | - | A | U | G | G |
| - | 0 | -2 | -4 | -6 | -8 |
| A | -2 | 2 | 0 | -2 | -4 |
| G | -4 | 0 | 2 | 2 | 0 |

Gap penalties

- linear: $\lambda(l_g) = -l_g d$

l_g : length of the gap; d : gap penalty;

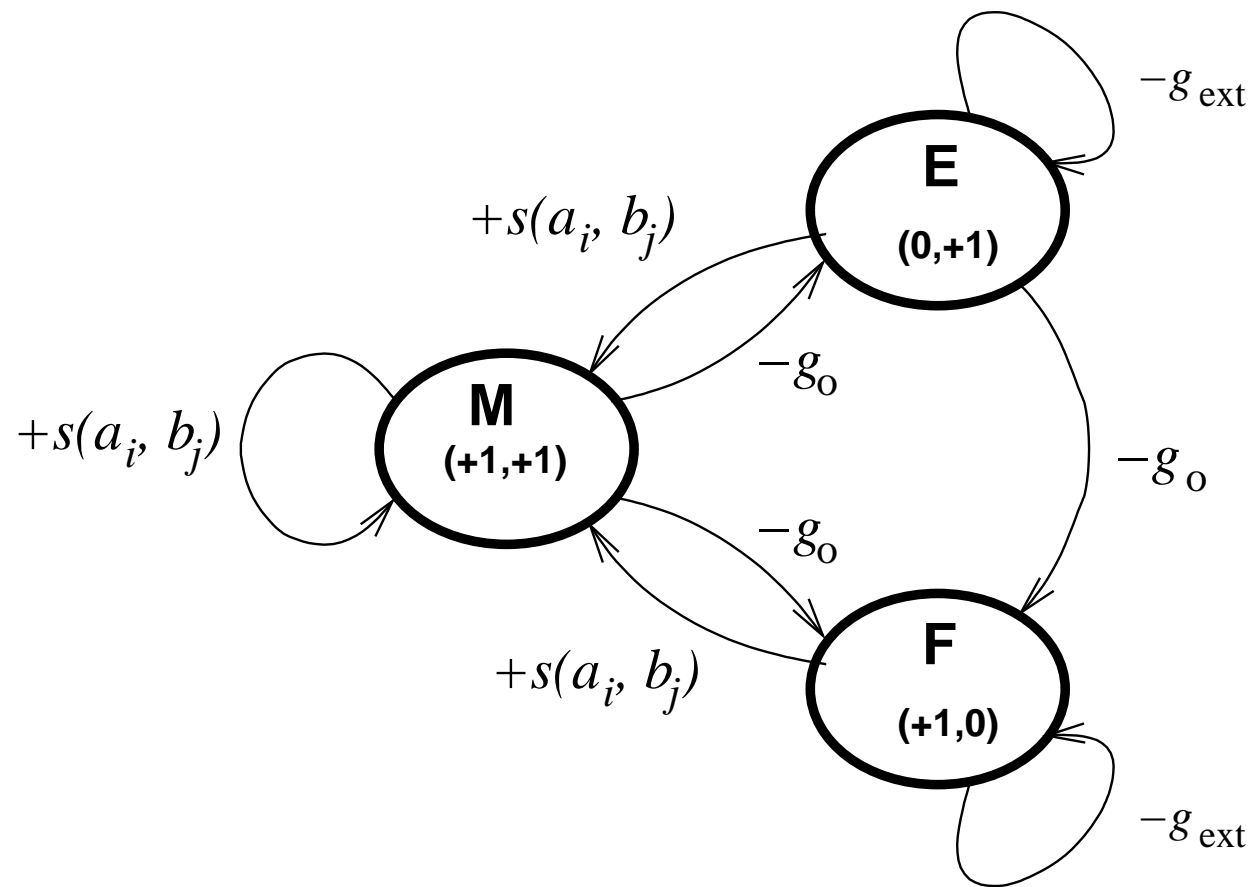
- affine: $\lambda(l_g) = -(g_o + g_{ext}(l_g - 1))$

g_o : gap open and g_{ext} : gap extension penalty;

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + s(a_i, b_j) \\ E(i-1, j-1) + s(a_i, b_j) \\ F(i-1, j-1) + s(a_i, b_j) \end{cases}$$

$$E(i, j) = \max \begin{cases} M(i, j-1) & - g_o \\ E(i, j-1) & - g_{ext} \end{cases}$$

$$F(i, j) = \max \begin{cases} M(i-1, j) & - g_o \\ E(i-1, j) & - g_o \\ F(i-1, j) & - g_{ext} \end{cases}$$



Alignment Scores: Substitutions

- Random Model: residues occur independently with some frequency q_a :

$$P(a, b | R) = q_a q_b$$

- Match Model: Aligned residue pairs occur with a joint probability p_{ab} :

$$P(a, b | M) = p_{ab}$$

- ratio of these two likelihoods: *odds ratio*:

$$\frac{P(a, b | M)}{P(a, b | R)} = \frac{p_{ab}}{q_a q_b} = p(a, b)$$

- additive scoring system: *log-odds ratio*

$$s(a, b) = \log \left(\frac{p_{ab}}{q_a q_b} \right)$$

affine transformation of log likelihood ratio

$$s(a, b) = \omega + k \log \left(\frac{p_{ab}}{q_a q_b} \right)$$

Alignment Scores: Gaps (Indels)

- probability of a gap

$$p(\mathbf{g}) = f(l_g) \prod_{i \text{ in } \mathbf{g}} q_{a_i}$$

- odds ratio

$$\frac{p(\mathbf{g})}{\prod_{i \text{ in } \mathbf{g}} q_{a_i}} = \frac{f(l_g) \prod_{i \text{ in } \mathbf{g}} q_{a_i}}{\prod_{i \text{ in } \mathbf{g}} q_{a_i}} = f(l_g)$$

- additive scoring system: *log-odds ratio*

$$\gamma(l_g) = \log f(l_g)$$

affine transformation

$$\gamma(l_g) = \omega + k \log f(l_g)$$

Probability of an alignment

Probability of a match

$$p_{ab} = p(a, b)q_aq_b$$

Probability of a gap

$$p(\mathbf{g}) = p_{-b_j} = p_{a_i-} = f(l_g) \prod_{i \text{ in } \mathbf{g}} q_{a_i}$$

Probability of the whole alignment

$$\begin{aligned} \text{Prob}(\mathcal{A}) &= \prod_{i \in \text{sub}} p_{a_i^* b_i^*} \prod_{i \in \text{ins}} p_{-b_i^*} \prod_{i \in \text{del}} p_{a_i^* -} = \\ &= p(\mathbf{a}) p(\mathbf{b}) \prod_{i \in \text{indel}} f(l_{\text{indel}}) \prod_{i \in \text{sub}} p(a_i^*, b_i^*) \end{aligned}$$

$$\frac{\text{Prob}(\mathcal{A})}{p(\mathbf{a}) p(\mathbf{b})} = \prod_{i \in \text{indel}} f(l_{\text{indel}}) \prod_{i \in \text{sub}} p(a_i^*, b_i^*)$$

$$S(\mathcal{A}) = \omega + k \left\{ \sum_{l_g} \log f(l_g) + \sum_{(i,j) \in \mathcal{A}} \log p(a_i, b_j) \right\}$$

$$\begin{aligned} e^{S(\mathcal{A})} &= e^\omega \prod_{l_g} e^{k \log f(l_g)} \prod_{(i,j) \in \mathcal{A}} e^{k \log p(a_i, b_j)} \\ &= e^\omega \left\{ \prod_{l_g} f(l_g) \prod_{(i,j) \in \mathcal{A}} p(a_i, b_j) \right\}^k \\ &= e^\omega \left\{ \frac{\text{Prob}(\mathcal{A})}{p(\mathbf{a})p(\mathbf{b})} \right\}^k \end{aligned}$$

$$e^{S(\mathcal{A})/k} = \frac{e^{(\omega/k)}}{p(\mathbf{a})p(\mathbf{b})} \text{Prob}(\mathcal{A})$$

Partition Function Z

$$c = \frac{p(\mathbf{a}) p(\mathbf{b})}{e^{(\omega/k)}}$$

$$\text{Prob}(\mathcal{A}) = c e^{\frac{S(\mathcal{A})}{k}}$$

$$\sum_{\mathcal{A}} \text{Prob}(\mathcal{A}) = 1 = c \sum_{\mathcal{A}} e^{\frac{S(\mathcal{A})}{k}}$$

partition function Z : sum of the Boltzmann factors for all possible states

$$Z = \sum_{\mathcal{A}} e^{\frac{S(\mathcal{A})}{k}}$$

$c = \frac{1}{Z}$: relationship between “energies” and probabilities of states

$$\text{Prob}(\mathcal{A}) = \frac{1}{Z} e^{\frac{S(\mathcal{A})}{k}}$$

Partition Function Z

- e.g. Dayhoffs PAM Matrices

$$s(a, b) = 10 \log_{10} \left(\frac{p_{ab}}{q_a q_b} \right)$$
$$k = \frac{10}{\log_e 10} \approx 4.3429$$

- temperature dependent partition function

$$Z(T) = \sum_{\mathcal{A}} e^{\frac{S(\mathcal{A})}{kT}} = \sum_{\mathcal{A}} e^{\beta S(\mathcal{A})}$$

Match Probabilities

$$\begin{aligned}
 \Omega_{i,j} &= \{\mathcal{A} \mid (i,j) \in \mathcal{A}\} \\
 \text{Prob}(\Omega_{i,j}) &= \frac{1}{Z} \sum_{\mathcal{A} \in \Omega_{i,j}} e^{S_{\mathcal{A}}} = \frac{Z(\Omega_{i,j})}{Z} \\
 S(\mathcal{A} \in \Omega_{i,j}) &= S(\mathcal{A}_{1,1}^{i,j}) + S(\mathcal{A}_{i,j}^{m,n}) - s(a_i, b_j) \\
 \\
 Z(\Omega_{i,j}) &= \sum_{\mathcal{A} \in \Omega_{i,j}} e^{\beta S(\mathcal{A}_{1,1}^{i,j}) + \beta S(\mathcal{A}_{i,j}^{m,n}) - \beta s(a_i, b_j)} \\
 &= \underbrace{\sum_{\mathcal{A} \in \mathcal{A}_{1,1}^{i,j}} e^{\beta S(\mathcal{A}_{1,1}^{i,j})}}_{Z_{ij}^M} \times \underbrace{\sum_{\mathcal{A} \in \mathcal{A}_{i,j}^{m,n}} e^{\beta S(\mathcal{A}_{i,j}^{m,n})}}_{\widehat{Z}_{ij}^M} \times e^{-\beta s(a_i, b_j)} \\
 &= Z_{ij}^M \widehat{Z}_{ij}^M e^{-\beta s(a_i, b_j)} \\
 \\
 \text{prob}(i,j) &= \frac{Z_{i,j} \widehat{Z}_{i,j}}{e^{s(a_i, b_j)} Z}
 \end{aligned}$$

Algorithm for the calculation of the partition function

$$Z_{i,j}^M = \left(Z_{i-1,j-1}^M + Z_{i-1,j-1}^E + Z_{i-1,j-1}^F \right) e^{\beta s(a_i, b_i)}$$

$$Z_{i,j}^E = Z_{i,j-1}^M e^{\beta g_o} + Z_{i,j-1}^E e^{\beta g_{\text{ext}}}$$

$$Z_{i,j}^F = \left(Z_{i-1,j}^M + Z_{i-1,j}^E \right) e^{\beta g_o} + Z_{i-1,j}^F e^{\beta g_{\text{ext}}}$$

$$Z_{i,j} = Z_{i,j}^M + Z_{i,j}^E + Z_{i,j}^F$$

Stochastic Backtracking

Probability of each state depends on the previous state
match

$$p(a, b) = \frac{Z_{i-1,j-1}^M e^{\beta s(a_i, b_j)}}{Z_{i,j}^M}$$

$$p(-, b) = \frac{Z_{i-1,j-1}^E e^{\beta s(a_i, b_j)}}{Z_{i,j}^M}$$

$$p(a, -) = \frac{Z_{i-1,j-1}^F e^{\beta s(a_i, b_j)}}{Z_{i,j}^M}$$

gap in **a**

$$p(a, b) = \frac{Z_{i-1,j}^M e^{\beta g_0}}{Z_{i,j}^F}$$

$$p(-, b) = \frac{Z_{i-1,j}^E e^{\beta g_0}}{Z_{i,j}^F}$$

$$p(a, -) = \frac{Z_{i,j-1}^F e^{\beta g_{\text{ext}}}}{Z_{i,j}^F}$$

where $Z_{i,j}^M = (Z_{i-1,j-1}^M + Z_{i-1,j-1}^E + Z_{i-1,j-1}^F) e^{\beta s(a_i, b_j)}$

where $Z_{i,j}^F = (Z_{i-1,j}^M + Z_{i-1,j}^E) e^{\beta g_0} + Z_{i-1,j}^F e^{\beta g_{\text{ext}}}$

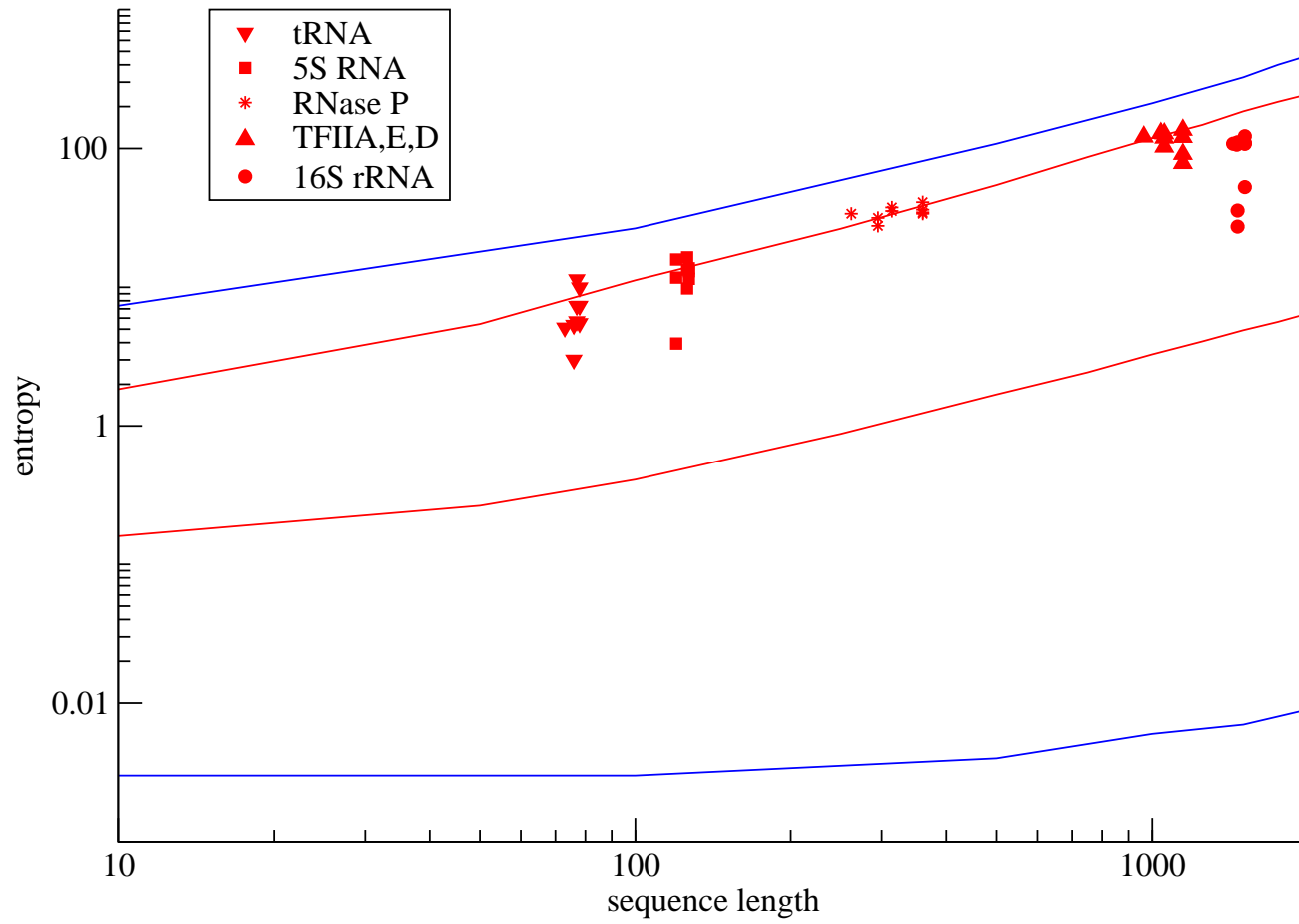
gap in **b**

$$p(a, b) = \frac{Z_{i,j-1}^M e^{\beta g_0}}{Z_{i,j}^E}$$

$$p(-, b) = \frac{Z_{i,j-1}^E e^{\beta g_{\text{ext}}}}{Z_{i,j}^E}$$

where $Z_{i,j}^E = Z_{i,j-1}^M e^{\beta g_0} + Z_{i,j-1}^E e^{\beta g_{\text{ext}}}$

$$\Delta S^{\text{ensemble}} = S(\mathcal{A}_{\text{opt}}) - kT \ln Z$$



CE
TOP 6.7
SARF2
MATRAS
COMPARER
consensus

```

-----QAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSEVPQNN-PELQAHAGKVFKLVEAAIQLE
-----DKTNVKAAWGKVGGAHAGEYGAEALERMFLSFPTTKTYF--PHFDLSHGSAQVKGHGKKVADALTNAVAHV----
--LTE-QAALVKSSWEEFN-----HTRFFILVLE-APAAK-----HAGK-----
--LSP-DKTNVKAAWGKVG-----YGAEALERMFL-FPTTK-----KVAD-----
--LTESQAALVKSSWEEFNANI-KHTRFFILVLEIAPAAKDLF----KGTSEVP--NPQLAHAGKVFKLVEE-----
--LSPADKTNVKAAWGKVGGAH-EGAEALERMFLSFPTTKTYF----FDLSHGSS--K-GHGKKVADALTNAVA-----
-ALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSE-VPQNNPELQAHAGKVFKLVEAAIQLE
-VLSPADKTNVKAAWGKVGGAHAGEYGAEALERMFLSFPTTKTYFPHFD-----LSHGSAQVKGHGKKVADALTNAVAHV-
GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSE-VPQNNPELQAHAGKVFKLVEAAIQLE
-VLSPADKTNVKAAWGKVGGAHAGEYGAEALERMFLSFPTTKTYFPHFD-----LSHGSAQVKGHGKKVADALTNAVAHVVD
GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSKLGKGTSEVPQNN-PELQAHAGKVFKLVEAAIQLE
-VLSPADKTNVKAAWGKVGGAHAGEYGAEALERMFLSFPTTKTYF--PHFDLSHGSAQVKGHGKKVADALTNAVAHV----

```

CE
TOP 6.7
SARF2
MATRAS
COMPARER
consensus

```

VTGVVVTDATLKNLGSVHV-SKGVADAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMD---
----DDMPNALALSADLHAHKLRVDPVNFKLLSHCLLVTLAAHLP AEFTPAVHASLDKFLASVSTVLT SKYR---
-----NLG-----VADAHFPVVKEAILKTIKEVVG-KWSEELNSAWTIAYDELAIVIKKEM---
-----ALS-----VDPVNFKLLSHCLLVTLAAHLP-EFTPAVHASLDKFLASVSTVLT SKY---
-----LKNLGSVHV-SKGVADAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMD---
-----LSALSADLHAHKLRVDPVNFKLLSHCLLVTLAAHLP AEFTPAVHASLDKFLASVSTVLT SKYR---
VTGVVVTDATLKNLGSVHVS-KGVADAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMD---
--DDMPNA--LSALSADLHAHKLRVDPVNFKLLSHCLLVTLAAHLP AEFTPAVHASLDKFLASVSTVLT SKYR---
VTGVVVTDATLKNLGSVHVSK-GVADAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMD DAA
D-----MPNALALSADLHAHKLRVDPVNFKLLSHCLLVTLAAHLP AEFTPAVHASLDKFLASVSTVLT SKYR---
VTGVVVTDATLKNLGSVHV-SKGVADAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMD DAA
----DDMPNALALSADLHAHKLRVDPVNFKLLSHCLLVTLAAHLP AEFTPAVHASLDKFLASVSTVLT SKYR---

```

struct. aln.
 GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSKLGKTSEVPQNN-PELQAHAGKVFKLVYEAAIQLE
 -VLSPADKTNVKAAWGKVGGAHAGEYGAEALERMFLSFPTTKTYF--PHFDLSHGSAQVKGHGKKVADALTNVAHV----

opt. aln. GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSEVPQNNPELQAHAGKVFKLVYEAAIQLEV
 -VLSPADKTNVKAAWGKVGGAHAGEYGAEALERMFLSFPTTKTYFPHE----DLSHGSAQVKGHGKKVADALTNVAHVDD

struct. aln.
 VTGVVVTDATLKNLGSVHV-SKGVADAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMDDAA
 ----DDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR----

opt. aln. TGVVVTDATLKNLGSVHVSKGVAD--AHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMDDAA
 MPNALSALSDLHAHKLRVDP-----VNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR----

