

Approximating an optimal Sum-of-Pairs multiple alignment

Sum-of-Pairs Multiple Alignment

Problem

Let F be a set of k strings, each of length $\leq n$, we know how to an optimal SP-alignment M^* in time $O(n^k)$ using dynamic programming.

We will show how to compute an alignment M in time $O(k^2n^2)$ s.t.

$$SP(\mathcal{M}) < 2 \cdot SP(\mathcal{M}^*)$$

Notation

Let $d(x,y)$ be a metric between characters

Let $D(S,S')$ be the induced metric between strings as given by the optimal score of a global pairwise alignment (with linear gap cost)

Alignments consistent with a tree

$M:$	A	-	-	C	G	-	T	S_1
	A	T	T	C	-	-	T	
	C	T	-	C	G	-	A	
	A	-	-	C	G	G	T	
								S_4

$$\text{Score}(M(S_1, S_4)) = \text{Score}\left(\begin{array}{ccccccc} \text{A} & - & - & \text{C} & \text{G} & - & \text{T} \\ \text{A} & - & - & \text{C} & \text{G} & \text{G} & \text{T} \end{array}\right) \geq D(S_1, S_4)$$

Alignments consistent with a tree

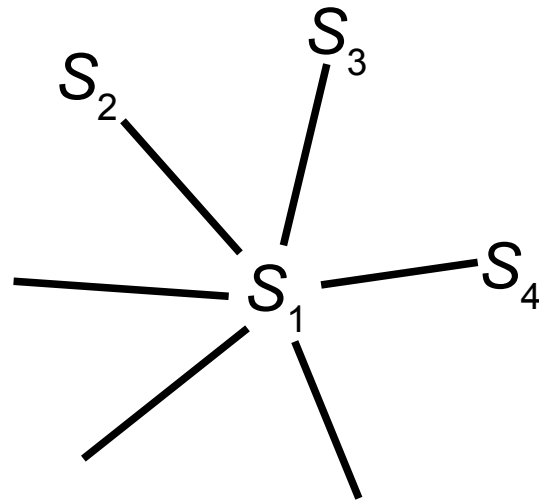
M :

A	-	-	C	G	-	T	S_1
A	T	T	C	-	-	T	
C	T	-	C	G	-	A	
A	-	-	C	G	G	T	S_4

$$\text{Score}(M(S_1, S_4)) = \text{Score}\left(\begin{array}{cccccc} \text{A} & - & - & \text{C} & \text{G} & - & \text{T} \\ \text{A} & - & - & \text{C} & \text{G} & \text{G} & \text{T} \end{array} \right) \geq D(S_1, S_4)$$

Definition (Gusfield, p. 347): Let F be a set of strings, and let T be a tree where each node is labeled with a distinct string from F . Then, a multiple alignment M of F is called *consistent* with T if the induced pairwise alignment of S_i and S_j has score $D(S_i, S_j)$ for each pair of strings (S_i, S_j) that label adjacent nodes in T .

The “guide” tree



s consistent with a tree

M :

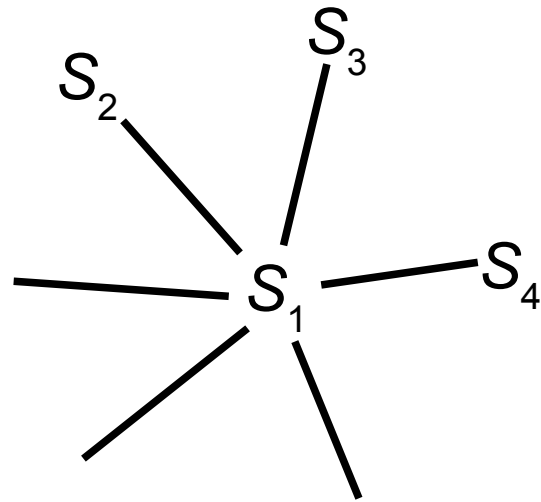
A	-	-	C	G	-	T	S_1
A	T	T	C	-	-	T	
C	T	-	C	G	-	A	
A	-	-	C	G	G	T	S_4

“=” if consistent
with “guide tree”

$$\text{Score}(M(S_1, S_4)) = \text{Score}\left(\begin{array}{ccccccc} \text{A} & - & - & \text{C} & \text{G} & - & \text{T} \\ \text{A} & - & - & \text{C} & \text{G} & \text{G} & \text{T} \end{array} \right) \geq D(S_1, S_4)$$

Definition (Gusfield, p. 347): Let F be a set of strings, and let T be a tree where each node is labeled with a distinct string from F . Then, a multiple alignment M of F is called *consistent* with T if the induced pairwise alignment of S_i and S_j has score $D(S_i, S_j)$ for each pair of strings (S_i, S_j) that label adjacent nodes in T .

The “guide” tree



s consistent with a tree

M :

A	-	-	C	G	-	T	S_1
A	T	T	C	-	-	T	
C	T	-	C	G	-	A	S_4
A	-	-	C	G	G	T	

“=” if consistent with “guide tree”

$$\text{Score}(M(S_1, S_4)) = \text{Score}\left(\begin{array}{cccccc} \text{A} & - & - & \text{C} & \text{G} & - & \text{T} \\ \text{A} & - & - & \text{C} & \text{G} & \text{G} & \text{T} \end{array} \right) \geq D(S_1, S_4)$$

Definition (Gusfield, p. 347): Let F be a set of strings, and let T be a tree where each node is labeled with a distinct string from F . Then, a multiple alignment M of F is called *consistent* with T if the induced pairwise alignment of S_i and S_j has score $D(S_i, S_j)$ for each pair of strings (S_i, S_j) that label adjacent nodes in T .

Lemma 14.6.1 (Gusfield, p. 347): For any set of strings F and for any tree T whose nodes are labeled by distinct strings of F , we can efficiently find a multiple alignment $M(T)$ of F that is consistent with T .

Algorithm

Input: A set F of k strings, each of length $\leq n$

Step 1 – Find the “center” string

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized.

Call the remaining strings S_2, S_3, \dots, S_k

Algorithm

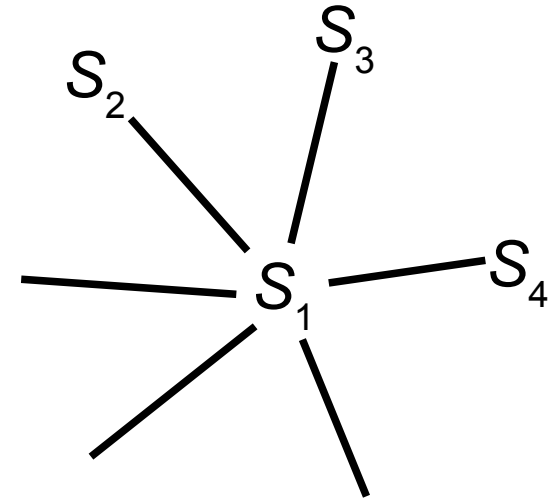
Input: A set F of k strings, each of

Step 1 – Find the “center” s

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized.

Call the remaining strings S_2, S_3, \dots, S_k

The “guide” tree



Algorithm

Input: A set F of k strings, each of

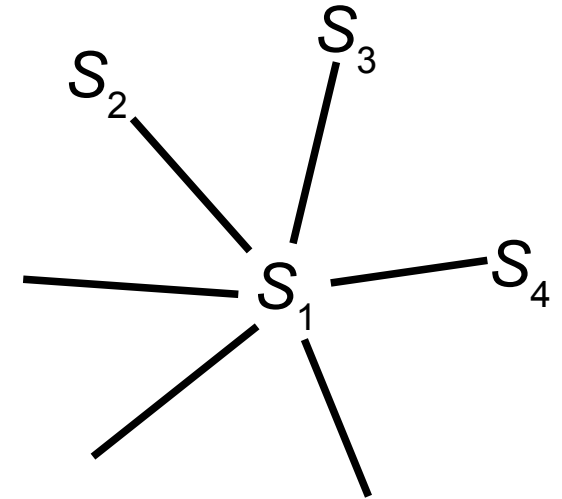
Step 1 – Find the “center” s

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized

Takes time $O(n^2)$ for each of the $k(k-1)$ pairs of strings

Call the remaining strings S_2, S_3, \dots, S_k

The “guide” tree



Algorithm

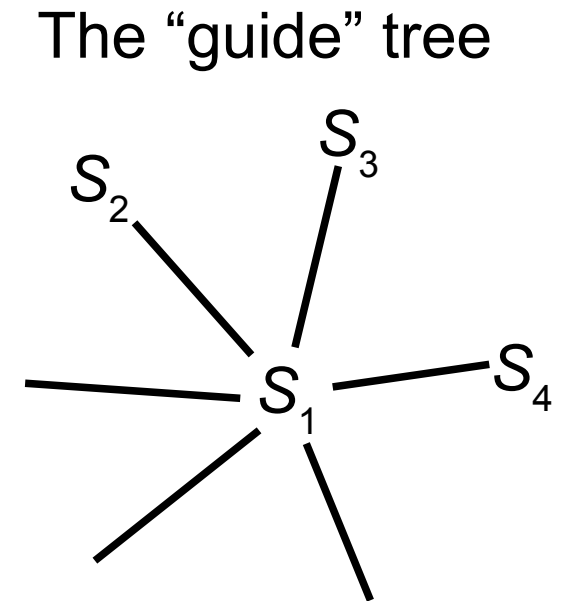
Input: A set F of k strings, each of

Step 1 – Find the “center” S_1

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized

Takes time $O(n^2)$ for each of the $k(k-1)$ pairs of strings

Call the remaining strings S_2, S_3, \dots, S_k



Step 2 – Construct alignment M cf. “guide tree”

$M_1 = [S_1]$

for $i = 2$ to k :

$A = \text{optalign}(S_1, S_i)$

$M_i = \text{"}M_{i-1} \text{ extended with } A\text{"}$

$M = M_k$

Example

Assume that $i=4$:

```

      A - - C G T
      A T T C - T
M3 = C T - C G A
  
```

$S_4 =$ A C G G T

```

      A C G - T
A = A C G G T
  
```

Extend M_3 with A gives:

```

      A - - C G - T
      A T T C - - T
      C T - C G - A
M4 = A - - C G G T
  
```

Note that the new column does not affect $\text{Score}(S_1, S_i)$ for $i < 4$

Algorithm

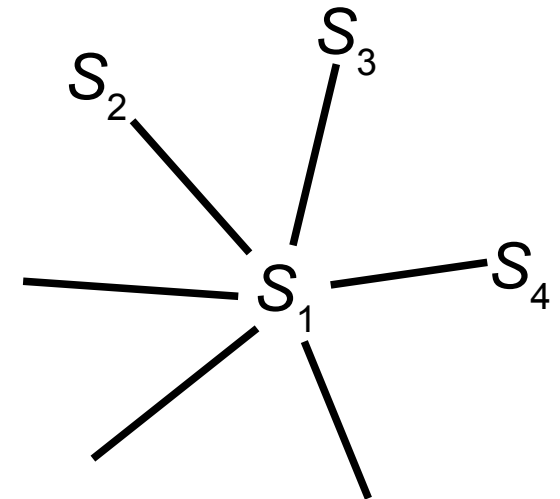
Let F of k strings, each of

– Find the “center” S_1

that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized

remaining strings S_2, S_3, \dots, S_k

The “guide” tree



Takes time $O(n^2)$ for each of the $k(k-1)$ pairs of strings

Construct alignment M cf. “guide tree”

```

M1 = [S1]
for i = 2 to k:
    A = optalign(S1, Si)
    Mi = "Mi-1 extended with A"
M = Mk
  
```

Example

Assume that $i=4$:

```

      A - - C G T
      A T T C - T
M3 = C T - C G A
  
```

```

S4 = A C G G T
  
```

```

      A C G - T
A = A C G G T
  
```

Extend M_3 with A gives:

```

      A - - C G - T
      A T T C - - T
      C T - C G - A
M4 = A - - C G G T
  
```

Note that the new column does not affect $\text{Score}(S_1, S_i)$ for $i < 4$

Algorithm

Let F of k strings, each of

– Find the “center” S_1

that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized

remaining strings S_2, S_3, \dots, S_k

– Construct alignment M of “guide tree”

```

M1 = [S1]
  
```

```

for i = 2 to k:
  
```

```

    A = optalign(S1, Si)
  
```

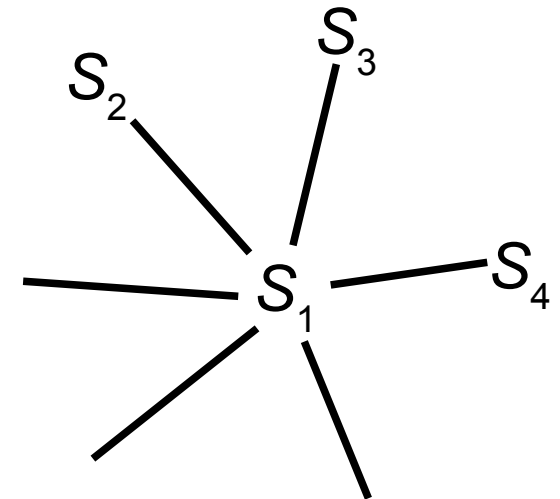
```

    Mi = "Mi-1 extended with A"
  
```

```

M = Mk
  
```

The “guide” tree



Takes time $O(n^2)$ for each of the $k(k-1)$ pairs of strings

Takes time $O(kn^2)$

Algorithm

Input: A set F of k strings, each of length $\leq n$

Step 1 – Find the “center” string

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized.

Call the remaining strings S_2, S_3, \dots, S_k

Step 2 – Construct alignment M cf. “guide tree”

```
M1 = [S1]  
for i = 2 to k:  
    A = optalign(S1, Si)  
    Mi = "Mi-1 extended with A"  
M = Mk
```

Algorithm

Input: A set F of k strings, each of length $\leq n$

Step 1 – Find the “center” string

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized.

Call the remaining strings S_2, S_3, \dots, S_k

Running time: $O(k^2n^2 + kn^2) = O(k^2n^2)$

Step 2 – Construct alignment M cf. “guide tree”

$M_1 = [S_1]$

for $i = 2$ to k :

$A = \text{optalign}(S_1, S_i)$

$M_i = \text{“}M_{i-1} \text{ extended with } A\text{”}$

$M = M_k$

How to extend M with A?

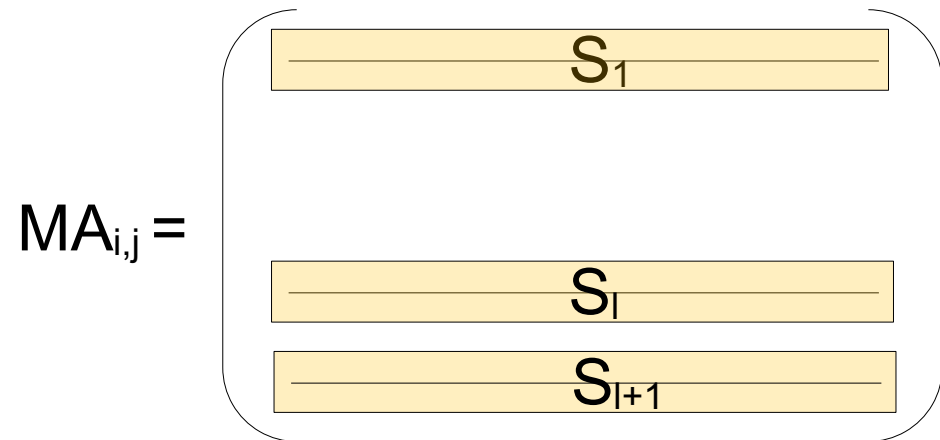
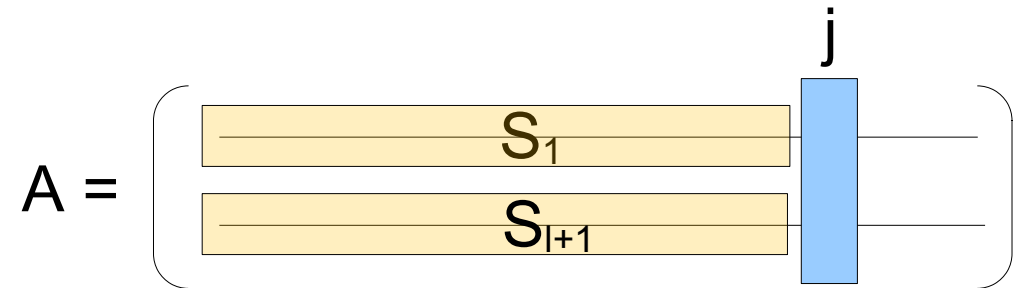
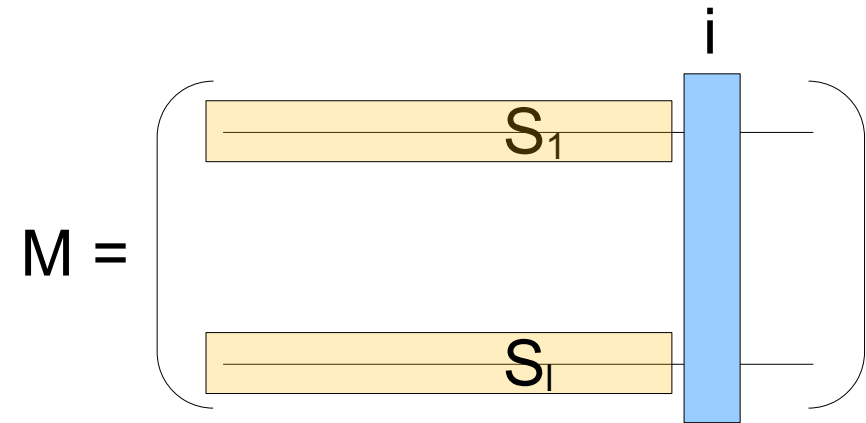
Let M be a multiple alignment, and let A be a pairwise alignment such that the first row of M and the first row of A is the same string if gaps are removed, i.e. like M_3 and A on the previous slides.

Let i be a column in M and j be column in A such the first row of M and A up to (but not including) column i and j respectively is the same string if gaps are removed.

Let MA be an extension of M with A that is consistent with the guide tree.

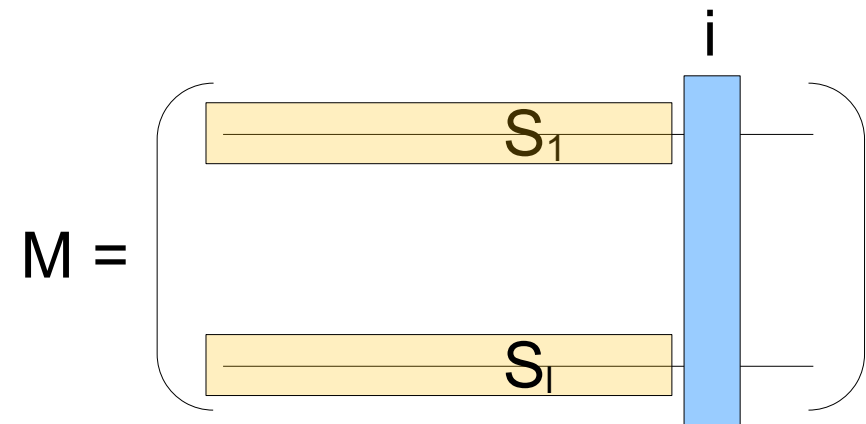
In particular, let $MA_{i,j}$ be an extension of all columns before column i in M and all columns before column j in A that is consistent with the guide tree.

Clearly $MA_{0,0}$ is the empty alignment, and $MA = MA_{\text{len}(M)+1, \text{len}(A)+1}$

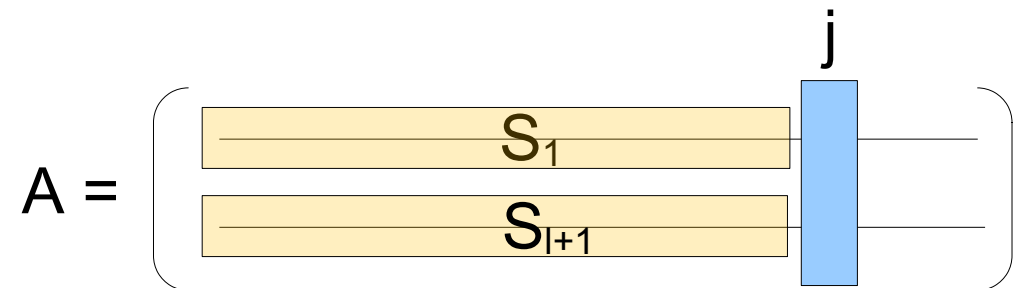


How to extend M with A?

Let M be a multiple alignment, and let A be a pairwise alignment such that the first row of M and the first row of A is the same string if gaps are removed, i.e. like M_3 and A on the previous slides.

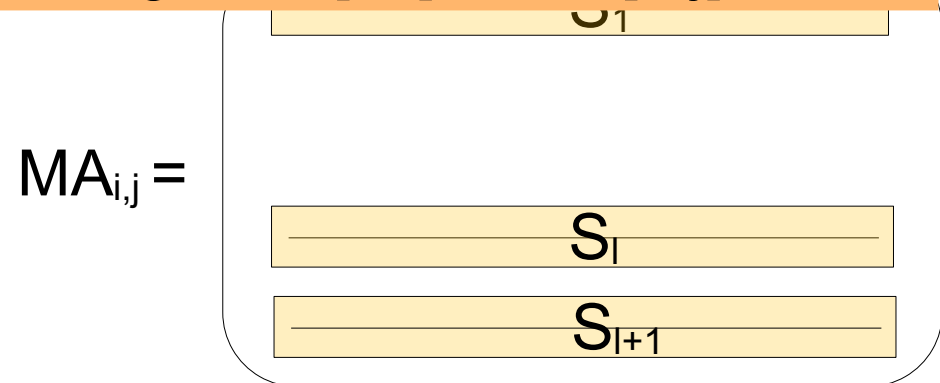


Let i be a column in M and j be column in A such the first row of M and A up to (but not including) column i and j respectively is the same string if gaps are removed.



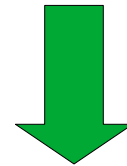
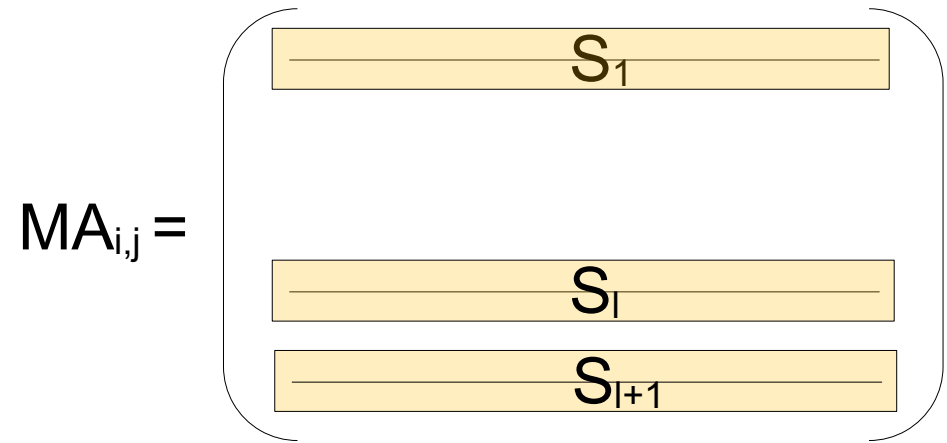
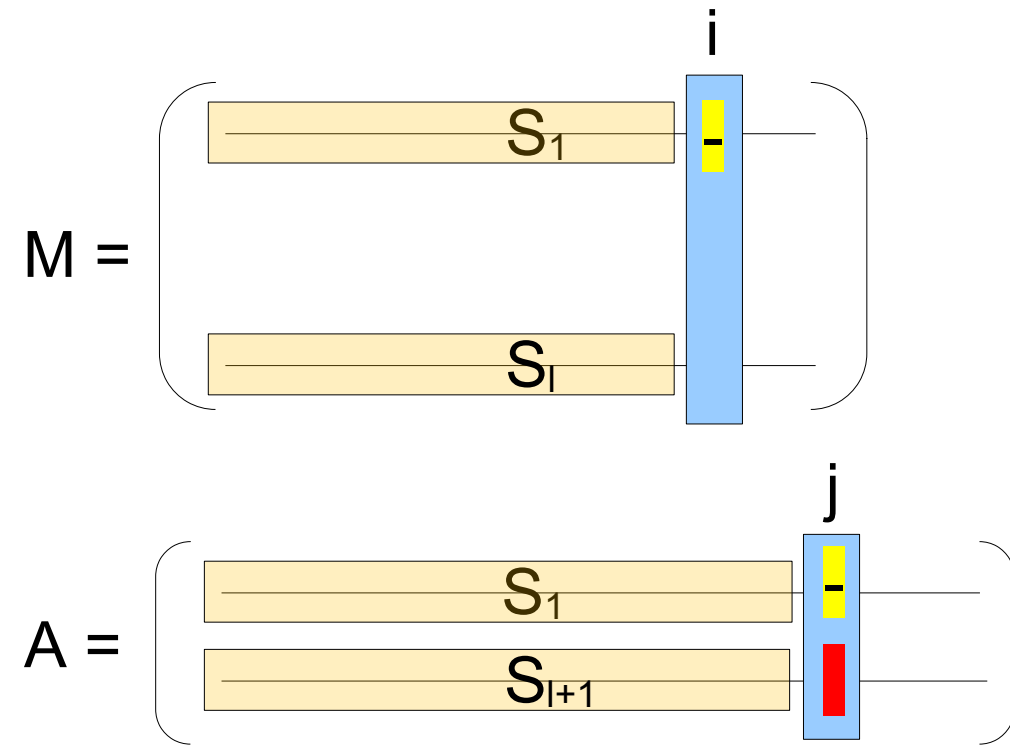
Let $MA_{i,j}$ be an extension of $M_{i,j}$ that is consistent with the guide tree. We extend $M_{i,j}$ depending on the content of column i and j . There are four cases depending on $M[1,i]$ and $A[1,j]$.

In particular, let $MA_{i,j}$ be an extension of all columns before column i in M and all columns before column j in A that is consistent with the guide tree.

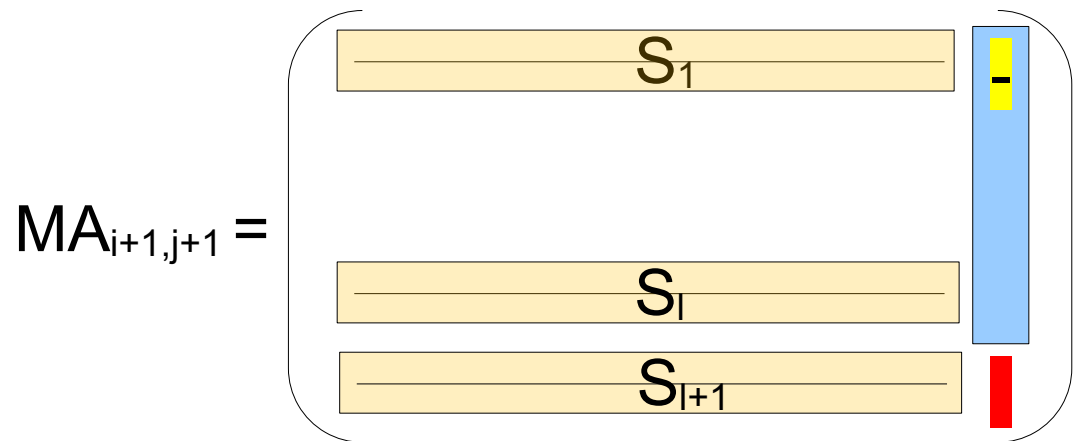


Clearly $MA_{0,0}$ is the empty alignment, and $MA = MA_{\text{len}(M)+1, \text{len}(A)+1}$

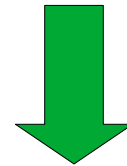
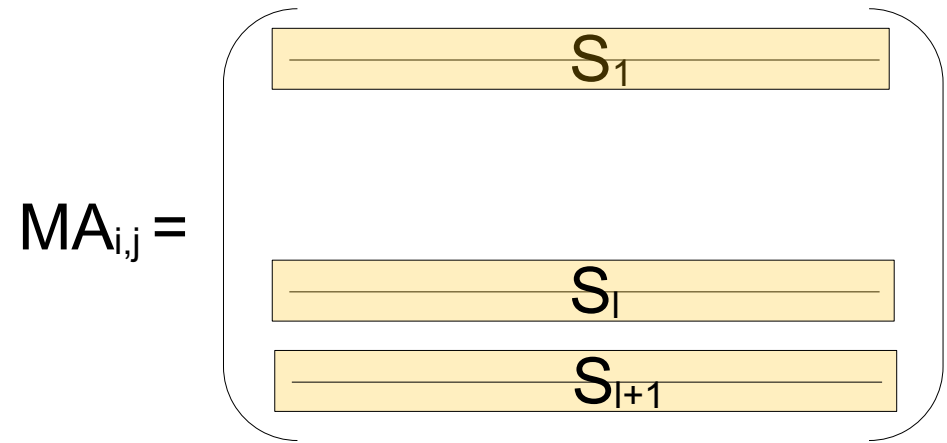
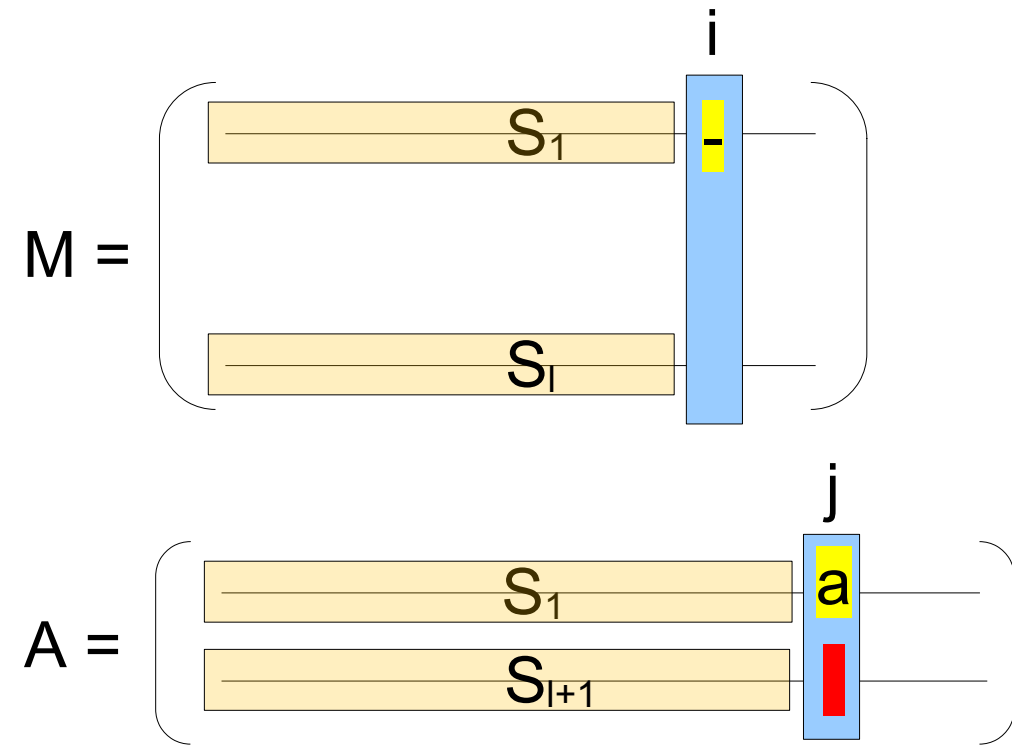
Case 1



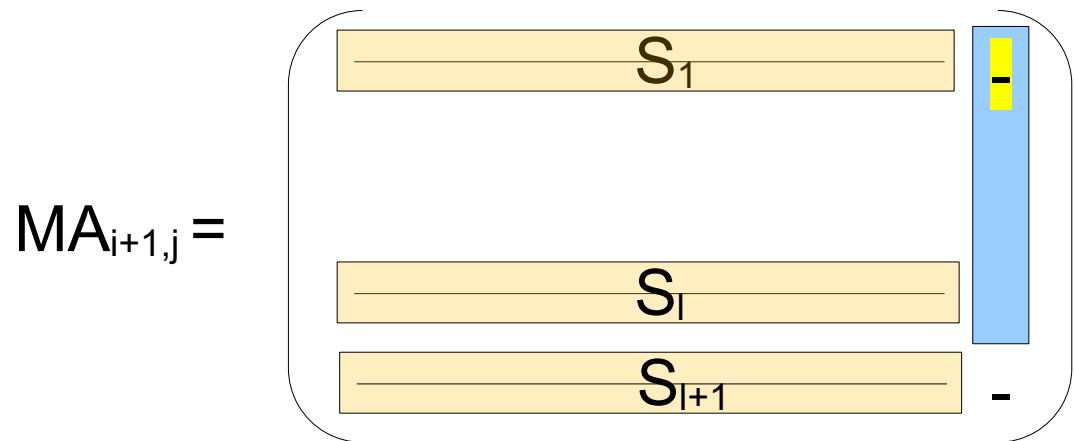
$MA_{i+1,j+1}$ is consistent with the guide tree. Continue with column $i+1$ in M and $j+1$ in A



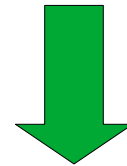
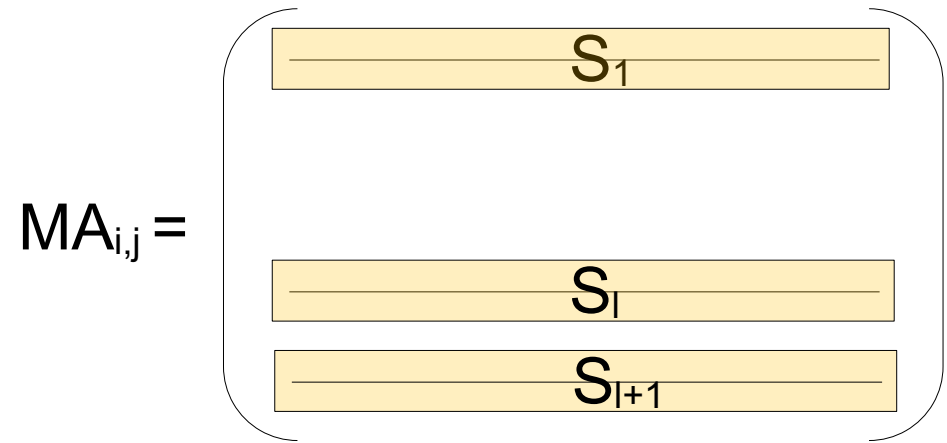
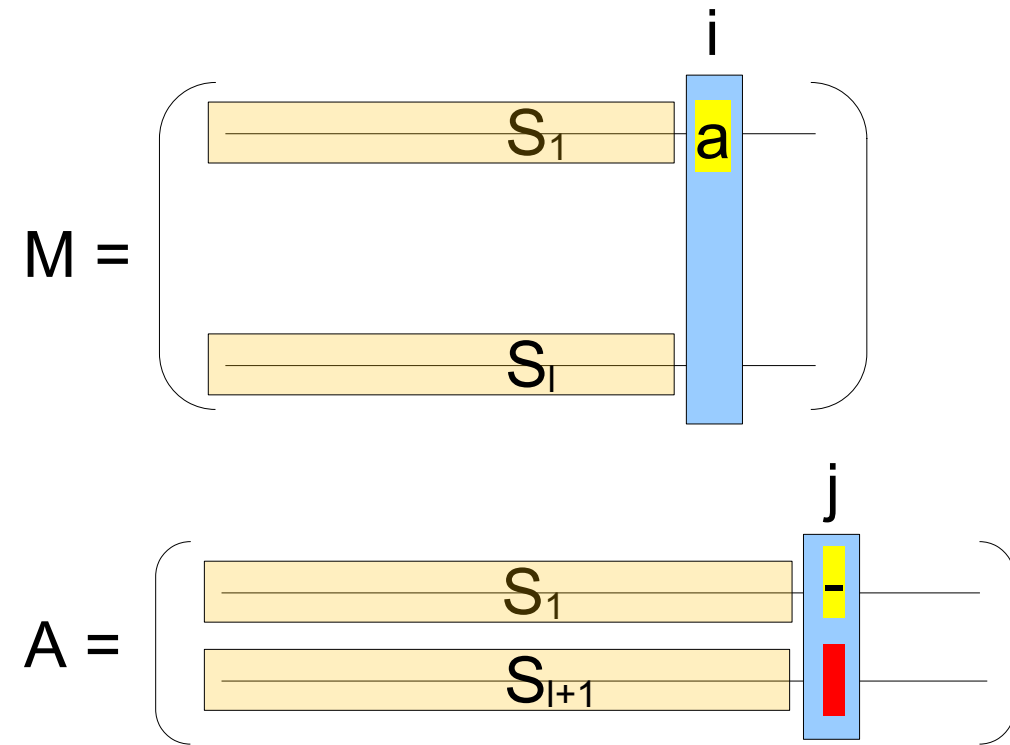
Case 2



$MA_{i+1,j}$ is consistent with the guide tree. Continue with column $i+1$ in M and j in A

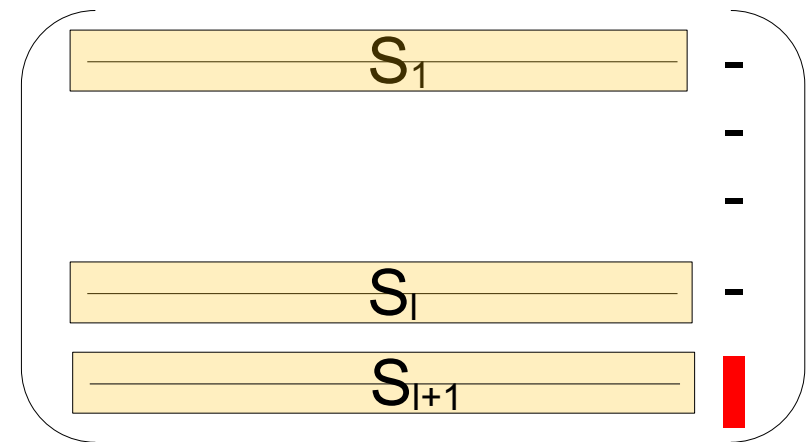


Case 3

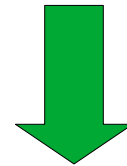
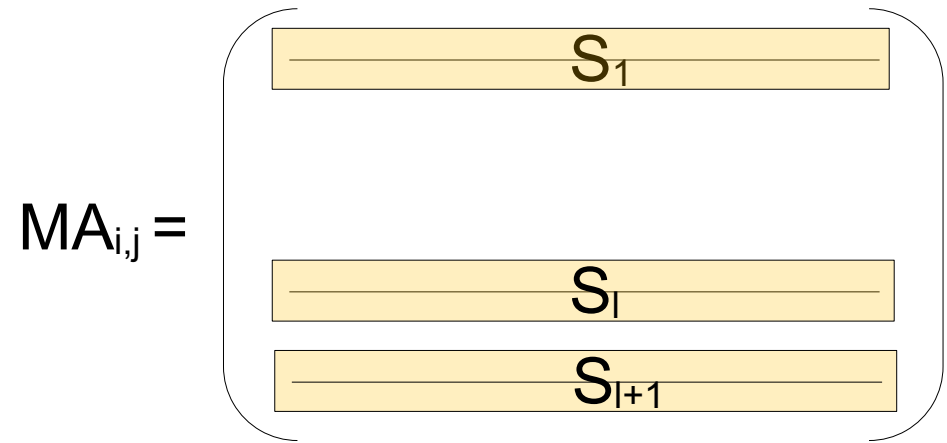
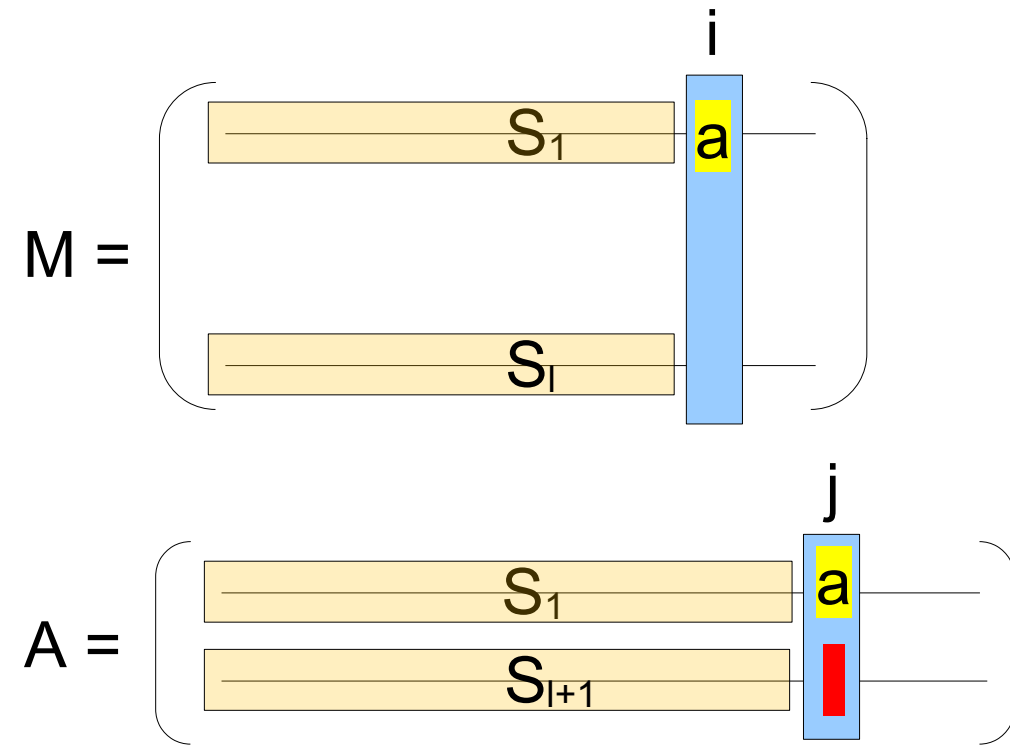


$MA_{i,j+1}$ is consistent with the guide tree. Continue with column i in M and $j+1$ in A

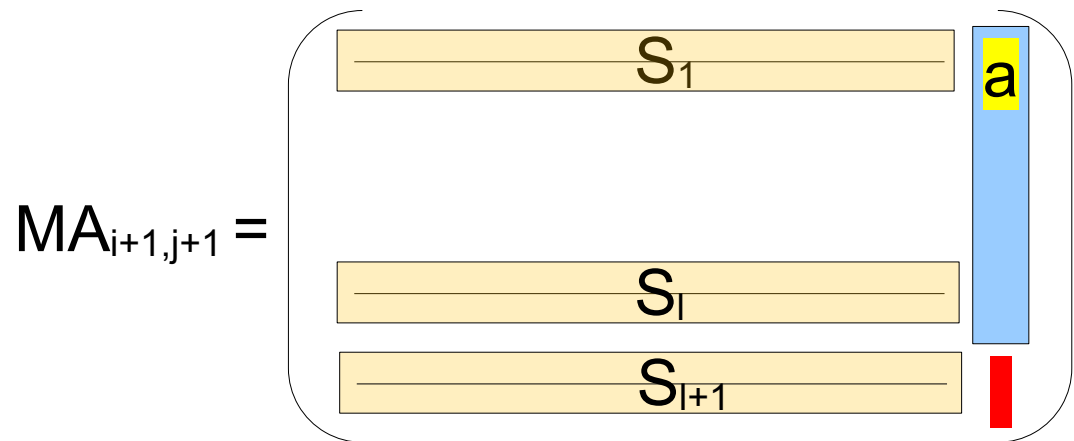
$MA_{i,j+1} =$



Case 4

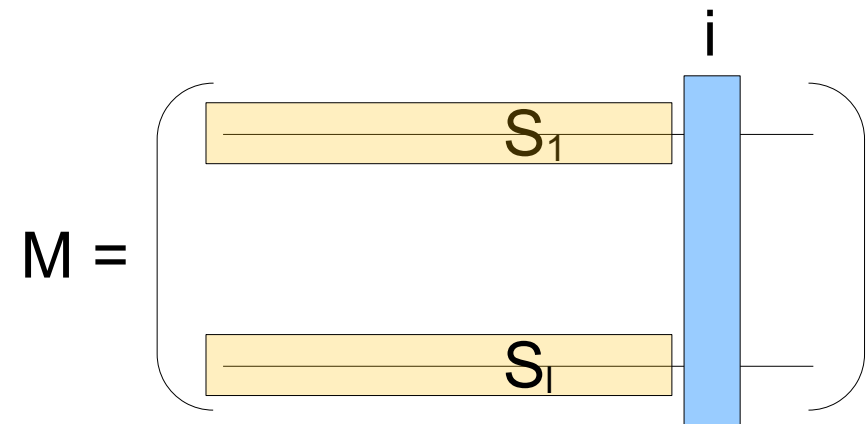


$MA_{i+1,j+1}$ is consistent with the
 guide tree. Continue with
 column $i+1$ in M and $j+1$ in A

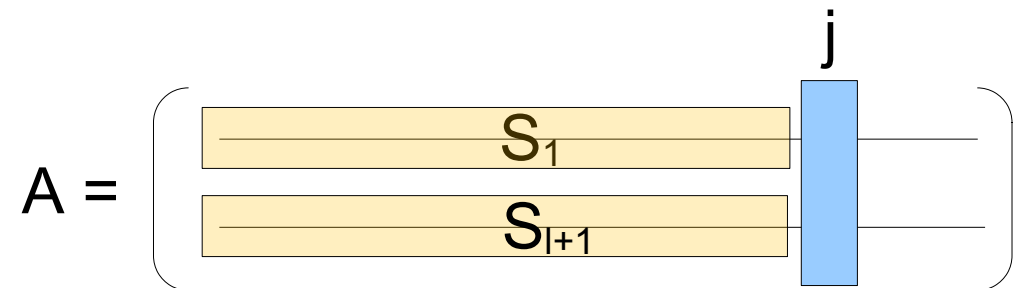


How to extend M with A?

Let M be a multiple alignment, and let A be a pairwise alignment such that the first row of M and the first row of A is the same string if gaps are removed, i.e. like M_3 and A on the previous slides.



Let i be a column in M and j be column in A such the first row of M and A up to (but not including) column i and j respectively is the same string if gaps are removed.



Let $M_{i,j}$ be the multiple alignment M extended with A up to column j . We extend $M_{i,j}$ depending on the content of column i and j . There are four cases depending on $M[1,i]$ and $A[1,j]$.

In particular, if we can extend $M_{i,j}$ in time $O(1)$, then we can build MA from $MA_{0,0}$ in time **$O(\text{\#columns in } M + \text{\#columns in } A)$** by iteratively applying cases 1-4. Think about how to represent M , A , and MA

Clearly, MA can be represented as a matrix of size $(\text{len}(M)+1, \text{len}(A)+1)$.

Algorithm

Input: A set F of k strings, each of length $\leq n$

Step 1 – Find the “center” string

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized.

Call the remaining strings S_2, S_3, \dots, S_k

Step 2 – Construct alignment M cf. “guide tree”

```
M1 = [S1]  
for i = 2 to k:  
    A = optalign(S1, Si)  
    Mi = "Mi-1 extended with A"  
M = Mk
```

Algorithm

Input: A set F of k strings, each of length $\leq n$

Step 1 – Find the “center” string

Find S_1 such that $\sum_{S \in \mathcal{F} - S_1} D(S_1, S)$ is minimized.

Call the remaining strings S_2, S_3, \dots, S_k

Running time: $O(k^2n^2 + kn^2) = O(k^2n^2)$

Step 2 – Construct alignment M cf. “guide tree”

$M_1 = [S_1]$

for $i = 2$ to k :

$A = \text{optalign}(S_1, S_i)$

$M_i = \text{“}M_{i-1} \text{ extended with } A\text{”}$

$M = M_k$

Approximation Ratio?

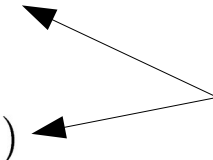
We want to prove that $SP(M) < 2 SP(M^*)$. How?

Approximation Ratio, part 1

Finding an upper bound of the computed alignment M

$$\begin{aligned}\text{SP}(\mathcal{M}) &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}(S_i, S_j)) \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d(i, j)\end{aligned}$$

The score of the alignment
of S_i and S_j as induced by M



Approximation Ratio, part 1

Finding an upper bound of the computed alignment M

$$\text{SP}(\mathcal{M}) = \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}(S_i, S_j))$$

$$= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d(i, j)$$

The score of the alignment of S_i and S_j as induced by M

Using the triangle-inequality and symmetry. Valid because the substitution matrix is metric

$$\leq \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k (d(i, 1) + d(1, j))$$

$$= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k (d(1, i) + d(1, j))$$

Approximation Ratio, part 1

Finding an upper bound of the computed alignment M

$$\text{SP}(\mathcal{M}) = \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}(S_i, S_j))$$

$$= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d(i, j)$$

The score of the alignment of S_i and S_j as induced by M

Using the triangle-inequality and symmetry. Valid because the substitution matrix is metric

$$\leq \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k (d(i, 1) + d(1, j))$$

$$= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k (d(1, i) + d(1, j))$$

$$= \frac{1}{2} \sum_{l=2}^k 2(k-1)d(1, l)$$

Expanding and rewriting the sum

$$= (k-1) \sum_{l=2}^k \text{Score}(\mathcal{M}(S_1, S_l))$$

Approximation Ratio, part 1

Finding an upper bound of the computed alignment M

$$\text{SP}(\mathcal{M}) = \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}(S_i, S_j))$$

$$= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d(i, j)$$

The score of the alignment of S_i and S_j as induced by M

Using the triangle-inequality and symmetry. Valid because the substitution matrix is metric

$$\leq \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k (d(i, 1) + d(1, j))$$

$$= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k (d(1, i) + d(1, j))$$

$$= \frac{1}{2} \sum_{l=2}^k 2(k-1)d(1, l)$$

Expanding and rewriting the sum

$$= (k-1) \sum_{l=2}^k \text{Score}(\mathcal{M}(S_1, S_l))$$

M is consistent with the guide tree, where S_1 is the center

$$= (k-1) \sum_{l=2}^k D(S_1, S_l)$$

Expanding and rewriting ...

$$\begin{aligned}\frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k (d(1, i) + d(1, j)) &= \frac{1}{2} \left(\sum_{i=1}^k \sum_{j=1, i \neq j}^k d(1, i) + \sum_{i=1}^k \sum_{j=1, i \neq j}^k d(1, j) \right) \\&= \frac{1}{2} \left(\sum_{i=1}^k (k-1) d(1, i) + \left(\sum_{i=1}^k \sum_{j=1}^k d(1, j) - \sum_{j=1}^k d(1, j) \right) \right) \\&= \frac{1}{2} \left(\sum_{i=1}^k (k-1) d(1, i) + \left(\sum_{j=1}^k k \cdot d(1, j) - \sum_{j=1}^k d(1, j) \right) \right) \\&= \frac{1}{2} \left(\sum_{i=1}^k (k-1) d(1, i) + \sum_{j=1}^k (k-1) d(1, j) \right) \\&= \frac{1}{2} \left(\sum_{l=2}^k (k-1) d(1, l) + \sum_{l=2}^k (k-1) d(1, l) \right) \\&= \frac{1}{2} \sum_{l=2}^k 2(k-1) d(1, l)\end{aligned}$$

Approximation Ratio, part 2

Finding a lower bound of the score of an optimal alignment M^*

$$\begin{aligned}\text{SP}(\mathcal{M}^*) &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}^*(S_i, S_j)) \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d^*(i, j)\end{aligned}$$

The score of the alignment
of S_i and S_j as induced by M^*

Approximation Ratio, part 2

Finding a lower bound of the score of an optimal alignment M^*

$$\begin{aligned} \text{SP}(\mathcal{M}^*) &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}^*(S_i, S_j)) \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d^*(i, j) \\ &\geq \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k D(S_i, S_j) \end{aligned}$$

Nothing is better than the optimal scores

The score of the alignment of S_i and S_j as induced by M^*

Approximation Ratio, part 2

Finding a lower bound of the score of an optimal alignment M^*

$$\begin{aligned}\text{SP}(\mathcal{M}^*) &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}^*(S_i, S_j)) \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d^*(i, j)\end{aligned}$$

The score of the alignment of S_i and S_j as induced by M^*

Nothing is better than the optimal scores

$$\geq \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k D(S_i, S_j)$$

By choice of S_1 we have $\sum_j D(S_1, S_j) \leq \sum_j D(S_i, S_j)$ for any i

$$\geq \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k D(S_1, S_j)$$

Approximation Ratio, part 2

Finding a lower bound of the score of an optimal alignment M^*

$$\text{SP}(\mathcal{M}^*) = \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k \text{Score}(\mathcal{M}^*(S_i, S_j))$$

$$= \frac{1}{2} \sum_{i=1}^k \sum_{j=1, i \neq j}^k d^*(i, j)$$

The score of the alignment of S_i and S_j as induced by M^*

Nothing is better than the optimal scores

$$\geq \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k D(S_i, S_j)$$

By choice of S_1 we have $\sum_j D(S_1, S_j) \leq \sum_j D(S_i, S_j)$ for any i

$$\geq \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k D(S_1, S_j)$$

$$= \frac{1}{2} k \sum_{j=1}^k D(S_1, S_j)$$

$$= \frac{1}{2} k \sum_{l=2}^k D(S_1, S_l)$$

Rewriting and renaming

Approximation Ratio, part 3

Upper-bound

$$\text{SP}(\mathcal{M}) \leq (k-1) \sum_{l=2}^k D(S_1, S_l)$$

Lower-bound

$$\text{SP}(\mathcal{M}^*) \geq \frac{1}{2}k \sum_{l=2}^k D(S_1, S_l)$$

Using the upper- and lower-bounds we get

$$\frac{\text{SP}(\mathcal{M})}{\text{SP}(\mathcal{M}^*)} \leq \frac{(k-1) \sum_{l=2}^k D(S_1, S_l)}{\frac{1}{2}k \sum_{l=2}^k D(S_1, S_l)} = \frac{2(k-1)}{k} < 2$$

Approximation Ratio, part 3

Upper-bound

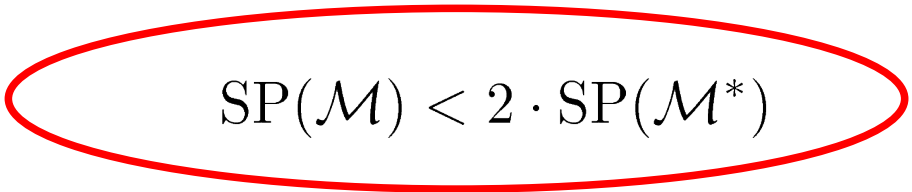
$$\text{SP}(\mathcal{M}) \leq (k-1) \sum_{l=2}^k D(S_1, S_l)$$

Lower-bound

$$\text{SP}(\mathcal{M}^*) \geq \frac{1}{2}k \sum_{l=2}^k D(S_1, S_l)$$

Using the upper- and lower-bounds we get

$$\frac{\text{SP}(\mathcal{M})}{\text{SP}(\mathcal{M}^*)} \leq \frac{(k-1) \sum_{l=2}^k D(S_1, S_l)}{\frac{1}{2}k \sum_{l=2}^k D(S_1, S_l)} = \frac{2(k-1)}{k} < 2$$


$$\text{SP}(\mathcal{M}) < 2 \cdot \text{SP}(\mathcal{M}^*)$$

Can we do better?

SP-multiple alignment is NP-complete [Wang and Jiang 1994]

PTAS by [Bafna, Lawler, Pevzner 1995] gives

$$\frac{\text{SP}(\mathcal{M})}{\text{SP}(\mathcal{M}^*)} \leq 2 - \frac{q}{k} \quad \text{in time } O(k^3 n^{2q-1}), \text{ where } 1 < q < k$$

Using the upper- and lower-bounds we get

$$\frac{\text{SP}(\mathcal{M})}{\text{SP}(\mathcal{M}^*)} \leq \frac{(k-1) \sum_{l=2}^k D(S_1, S_l)}{\frac{1}{2}k \sum_{l=2}^k D(S_1, S_l)} = \frac{2(k-1)}{k} < 2$$

$$\text{SP}(\mathcal{M}) < 2 \cdot \text{SP}(\mathcal{M}^*)$$