# Comparison of multiple biological sequences

Project in Bioinformatics
Astrid Christiansen, 201404423
Supervised by Christian N. S. Pedersen
Aarhus University

June 2022

# Table of Contents

# Table of Contents

# MSA and SP score

$$M = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ G & C & - & A & - & C & T & - & A & - & T \\ A & C & C & - & - & A & T & G & A & G & G \end{bmatrix}$$ ① ② ③

$$A_1 = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ G & C & - & A & - & C & T & - & A & - & T \end{bmatrix}$$ ①

$$A_2 = \begin{bmatrix} G & C & - & A & \cancel{\times} & C & T & - & A & - & T \\ A & C & C & - & \cancel{\times} & A & T & G & A & G & G \end{bmatrix}$$ ②

$$A_3 = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ A & C & C & - & - & A & T & G & A & G & G \end{bmatrix}$$ ③

$$SP(M) = \sum_{0 \le c < \ell} \sum_{0 \le i < j < k} d(M_{i,c}, M_{j,c})$$

$$= \sum_{0 \le i < j < k} \sum_{0 \le c < \ell} d(M_{i,c}, M_{j,c}).$$

# MSA and SP score

$$M = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ G & C & - & A & - & C & T & - & A & - & T \\ A & C & C & - & - & A & T & G & A & G & G \end{bmatrix} \Biggr\} \begin{matrix} ① \\ ② \end{matrix} \Biggr\rangle ③$$

$$A_1 = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ G & C & - & A & - & C & T & - & A & - & T \end{bmatrix} \quad ①$$

$$A_2 = \begin{bmatrix} G & C & - & A & C & T & - & A & - & T \\ A & C & C & - & A & T & G & A & G & G \end{bmatrix} \quad ②$$

$$A_3 = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ A & C & C & - & - & A & T & G & A & G & G \end{bmatrix} \quad ③$$

$$SP(M) = \sum_{0 \le c < \ell} \sum_{0 \le i < j < k} d(M_{i,c}, M_{j,c})$$

$$= \sum_{0 \le i < j < k} \sum_{0 \le c < \ell} d(M_{i,c}, M_{j,c}).$$

- Dyn. prog.: $O(n^k)$.

# MSA and SP score

$$M = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ G & C & - & A & - & C & T & - & A & - & T \\ A & C & C & - & - & A & T & G & A & G & G \end{bmatrix}$$

$$SP(M) = \sum_{0 \leq c < \ell} \sum_{0 \leq i < j < k} d(M_{i,c}, M_{j,c})$$

$$= \sum_{0 \leq i < j < k} \sum_{0 \leq c < \ell} d(M_{i,c}, M_{j,c}).$$

$$A_1 = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ G & C & - & A & - & C & T & - & A & - & T \end{bmatrix}$$

$$A_2 = \begin{bmatrix} G & C & - & A & C & T & - & A & - & T \\ A & C & C & - & A & T & G & A & G & G \end{bmatrix}$$

$$A_3 = \begin{bmatrix} A & C & C & A & C & C & A & G & - & G & T \\ A & C & C & - & - & A & T & G & A & G & G \end{bmatrix}$$

- Dyn. prog.: $O(n^k)$.
- Gusfield and MST algo.: $O(k^2 n^2)$.

# Table of Contents

# Gusfield's approximation algorithm

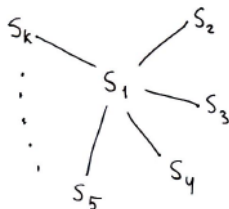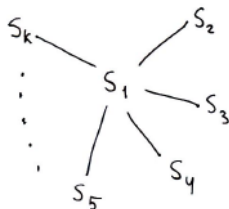$$\arg\min_{S_1} \sum_{s \in S} OPT(S_1, s).$$

# Gusfield's approximation algorithm

$$\arg \min_{S_1} \sum_{s \in S} OPT(S_1, s).$$


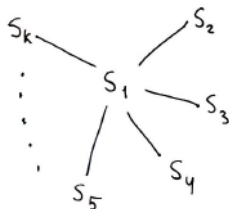
**Algorithm** GUSFIELD-MSA($S$)

1:   $S_1 \leftarrow center(S)$
2:   $M \leftarrow [\ ]$
3:   **for** $i = 2$ to $k$ **do**
4:      $A \leftarrow align(S_1, S_i)$
5:      $M \leftarrow extend(M, A)$
6:   **end for**
7:   **return** $M$

# Gusfield's approximation algorithm

$$\arg\min_{S_1} \sum_{s \in S} OPT(S_1, s).$$



**Algorithm** GUSFIELD-MSA($S$)

1: $S_1 \leftarrow center(S)$
2: $M \leftarrow [\,]$
3: **for** $i = 2$ to $k$ **do**
4:     $A \leftarrow align(S_1, S_i)$
5:     $M \leftarrow extend(M, A)$
6: **end for**
7: **return** $M$

▶ Time $O(k^2 n^2)$.

# Table of Contents

# MST algorithm

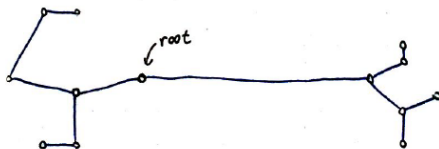# MST algorithm



---

**Algorithm** MST-MSA(*MST*)

---

1: $M \leftarrow [\ ]$
2: **for** $(S_i, S_j) \in MST$ **do**
3:     $A \leftarrow align(S_i, S_j)$
4:     $M \leftarrow extend(M, A)$
5: **end for**
6: **return** $M$

---

# Prim's algorithm

---

**Algorithm** MST-PRIM-SIMPLE($G$, $w$, $root$)

---

1: **for** each $u \in V[G]$ **do**
2:      $key[u] \leftarrow \infty$
3:      $\pi[u] \leftarrow NIL$
4: **end for**
5: $key[root] \leftarrow 0$
6: **for** $i = 1$ to $|V[G]|$ **do**
7:      $u \leftarrow \text{FIND-MIN}(V[G], key)$
8:      **report** $(\pi[u], u)$
9:      **for** each $v \in Adj[u]$ **do**
10:          **if** $v \notin MST$ and $w(u, v) < key[v]$ **then**
11:              $\pi[v] \leftarrow u$
12:              $key[v] \leftarrow w(u, v)$
13:          **end if**
14:      **end for**
15: **end for**

---

# Table of Contents

# Implementation

```
1  def MST_MSA_approx(seq_indices, seqs, sub_matrix, gap_cost):
2      M = []
3      # Keep track of seq index vs row index in M
4      seq_idx_to_row = [None] * len(seq_indices)
5      # Iterate MST edges
6      MST = MST_prim(seq_indices, seqs, sub_matrix, gap_cost)
7      for i in range(len(MST)):
8          # Edge (parent, node)
9          parent = MST[i][0]
10         node = MST[i][1]
11         # Fill out dyn. prog. table for pairwise alignment
12         table = fill_table(seqs[parent], seqs[node],
13             sub_matrix, gap_cost)
14         # Construct opt. pairwise align. from table
15         A = construct_alignment(table, seqs[parent], seqs[node],
16             sub_matrix, gap_cost)
17         if i == 0:
18             seq_idx_to_row[parent] = 0
19             seq_idx_to_row[node] = 1
20             M = A
21         else:
22             seq_idx_to_row[node] = len(M)
23             M = extend_M(M, A, seq_idx_to_row[parent])
24     return M
```
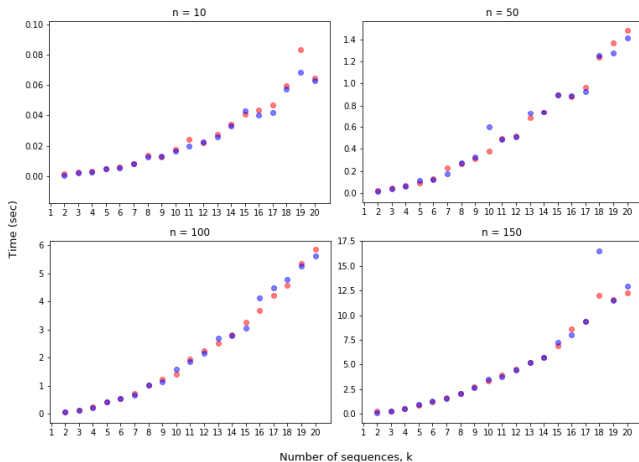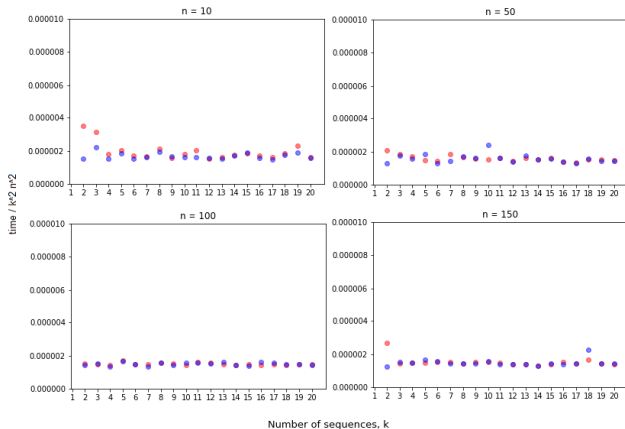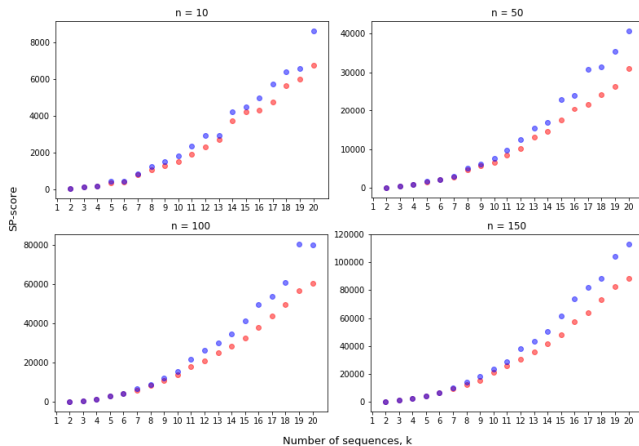
# Table of Contents

# Running time



Figure: Running time of the two algorithms on random data as a function of *k* for four different values of *n*. Gusfield is red, MST is blue.

# Running time



Figure: Running time of the two algorithms divided by expected worst case running time, $k^2 n^2$, on random data. Gusfield is red, MST is blue.

# SP score - random data



Figure: SP scores for the two algorithms as a function of *k* for four different values of *n* on random data. Gusfield is red, MST is blue.

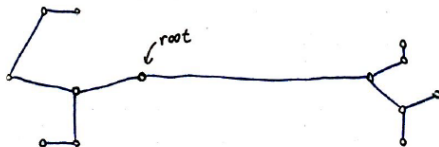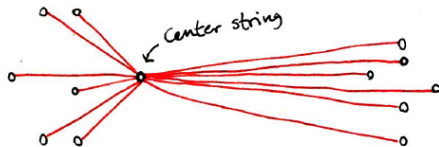# SP score - low mutation rate data



Figure: SP scores for the two algorithms as a function of *k* for four different values of *n* for simulated data with a low substitution rate of 0.05. Gusfield is red, MST is blue.

# Clustered data


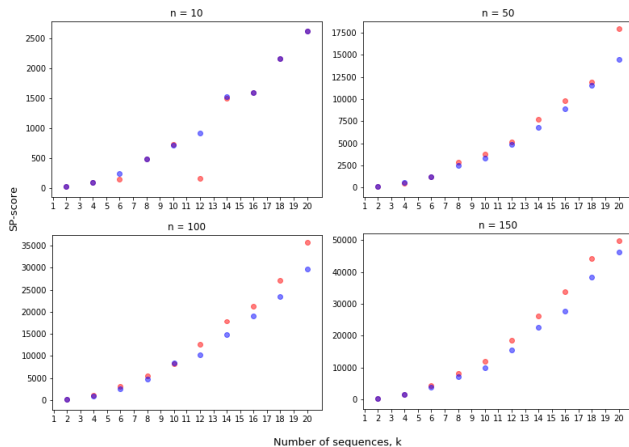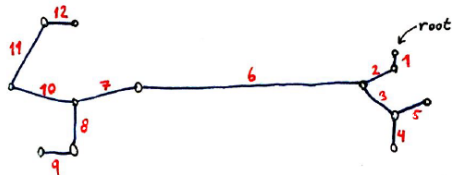
Center string

root
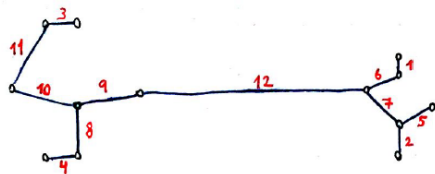
# SP score - two clusters



Figure: SP scores for the two algorithms as a function of $k$ for four different values of $n$ for simulated data consisting of two clusters (each made with a low substitution rate of 0.05). Gusfield is red, MST is blue.
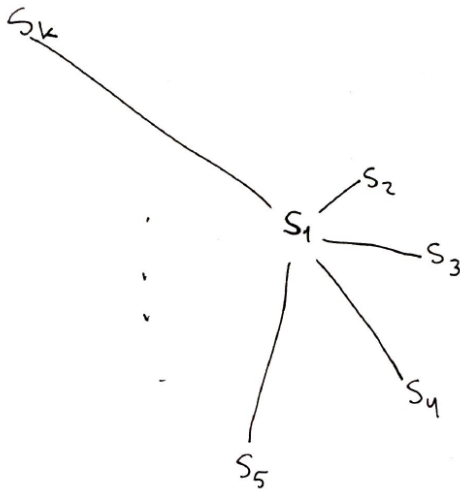
# Conclusion

- Random data, high mutation rate: Gusfield is best.
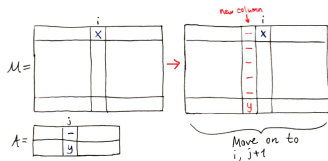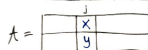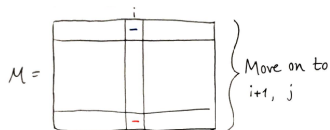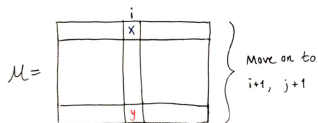- Low mutation rate: Very similar results.
- Two clusters: MST is best.

# Alternative MST approach: Kruskal

# Alternative star tree approach

# Extending $M$

# Prim's algorithm using queue

---

**Algorithm** MST-PRIM($G, w, root$)

---

```
 1: for each u ∈ V[G] do
 2:     key[u] ← ∞
 3:     π[u] ← NIL
 4: end for
 5: key[root] ← 0
 6: Q ← V[G]
 7: while Q ≠ ∅ do
 8:     u ← EXTRACT-MIN(Q)
 9:     report (π[u], u)
10:     for each v ∈ Adj[u] do
11:         if v ∈ Q and w(u, v) < key[v] then
12:             π[v] ← u
13:             key[v] ← w(u, v)
14:         end if
15:     end for
16: end while
```

---

# FIND-MIN

---

**Algorithm** FIND-MIN($V$, $key$)

---

1: $min\_key \leftarrow \infty$
2: $min\_node \leftarrow NIL$
3: **for** $v \in V$ **do**
4:     **if** $v \notin MST$ **and** $key[v] < min\_key$ **then**
5:         $min\_node \leftarrow v$
6:         $min\_key \leftarrow key[v]$
7:     **end if**
8: **end for**
9: **return** $min\_node$

---