

Repeat Detector Algorithm

V. Dion & T. Schuepbach

March 7, 2022

The present document gathers pieces of informations for the understanding of repeat detector tool. The algorithm is presented here in details.

Contents

Definitions	2
1 Standard PfTools algorithm	2
2 Repeat Decoder	5

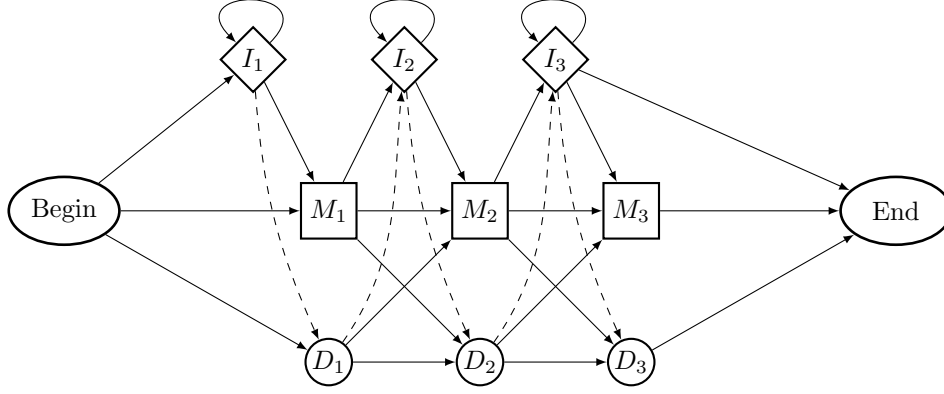


Figure 1: Simplified PfTools algorithm profile schema for a 3 sequence length profile. Simplifications arise from not showing all connections from Begin to nodes $[I_i, M_i, D_i]$ as well as nodes $[I_i, M_i, D_i]$ to End. It is worth noting the dashed that should not be used even though available. Indeed, mismatches should replace Insertion-Deletion or Deletion-Insertion paths thanks to extreme score loss.

Definitions

We define a sequence string of length N_s as $S = (s_1, s_2, \dots, s_{N_s})$ with $s_i, i \in [1, N_s]$, belonging to the alphabet \mathcal{A} . The cardinality of the set \mathcal{A} is further on referred to as N_α . Given that an alignment fills in a matrix $(N_s + 1) \times (N_p + 1)$, where N_p denotes the length of the profile. Each element of the matrix, denoted as a cell, $C_{i,j}$, $i \in [0, N_s]$, $j \in [0, N_p]$, holds pieces of informations of the current 4 possible states: *Match*, *Insertion*, *Deletion* and *Score*. That is $C_{i,j}^M$, $C_{i,j}^I$, $C_{i,j}^D$ and $C_{i,j}^S$ respectively. It is worth noting that $C_{i,j}^S$ is useless when $i = 0$ or $j = 0$.

Generalized profiles provide tables of scores to be used in the computation of $C_{i,j}^\alpha$, $\alpha \in \{M, I, D, S\}$. Those are

- the match/mismatch score table $M \in \mathbb{Z}^{N_p \times N_\alpha}$,
- the insertion score table $I \in \mathbb{Z}^{N_p-1 \times N_\alpha}$,
- the deletion score vector $D \in \mathbb{Z}^{N_p}$,
- the initial input score vectors $F^\delta \in \mathbb{Z}^{N_p+1}$, $\delta \in \{M, I, D\}$,
- the final output score vectors $L^\delta \in \mathbb{Z}^{N_p}$, $\delta \in \{M, I, D\}$,
- the state transition score vectors $T^{\beta \rightarrow \gamma} \in \mathbb{Z}^{1+N_p}$, $\beta, \gamma \in \{M, I, D, B, E\}$.

1 Standard PfTools algorithm

Let us define the indexing function $f(s) : \mathcal{A} \mapsto [1, N_\alpha]$ which translate the character s of the alphabet \mathcal{A} into its corresponding index within the score tables M, I . The recurrence relations for $C_{i,j}^\alpha$, $\alpha \in \{M, I, D, S\}$ pictured in figure 2 are as follow:

$$C_{i,j}^M = \begin{cases} F_0^M, & i = 0, j = 0 \\ \max \begin{cases} C_{0,j-1}^D + D_j + T_j^{D \rightarrow M} \\ F_j^M \end{cases}, & i = 0, j \in [1, N_p] \\ \max \begin{cases} C_{i-1,0}^I + I_{j,f(s_i)} + T_0^{I \rightarrow M} \\ T_0^{B \rightarrow M} \end{cases}, & i \in [1, N_s - 1], j = 0 \\ \max \begin{cases} C_{i-1,j-1}^M + M_{j,f(s_i)} + T_j^{M \rightarrow M} \\ C_{i-1,j}^I + I_{j,f(s_i)} + T_j^{I \rightarrow M} \\ C_{i,j-1}^D + D_j + T_j^{D \rightarrow M} \\ T_j^{B \rightarrow M} \end{cases}, & i \in [1, N_s - 1], j \in [1, N_p] \end{cases} \quad (1)$$



$$C_{i,j}^I = \begin{cases} F_0^I, & i = 0, j = 0 \\ \max \begin{cases} C_{0,j-1}^D + D_j + T_j^{D \rightarrow I} \\ F_j^I \end{cases}, & i = 0, j \in [1, N_p] \\ \max \begin{cases} C_{i-1,0}^I + I_{f(s_i),0} + T_0^{I \rightarrow I} \\ T_0^{B \rightarrow I} \end{cases}, & i \in [1, N_s - 1], j = 0 \\ \max \begin{cases} C_{i-1,j-1}^M + M_{j,f(s_i)} + T_j^{M \rightarrow I} \\ C_{i-1,j}^I + I_{j,f(s_i)} + T_j^{I \rightarrow I} \\ C_{i,j-1}^D + D_j + T_j^{D \rightarrow I} \\ T_j^{B \rightarrow I} \end{cases}, & i \in [1, N_s - 1], j \in [1, N_p] \end{cases} \quad (2)$$

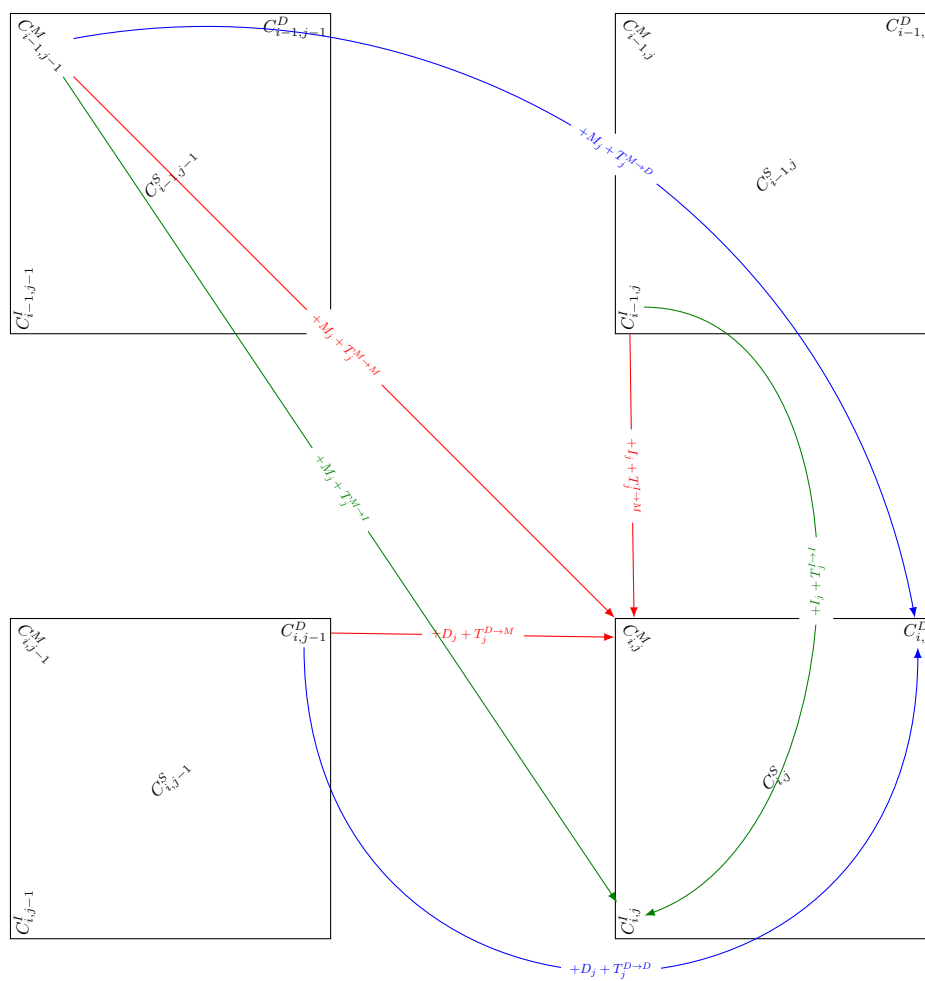


Figure 3: Standard PfTools algorithm schema

$$C_{i,j}^D = \begin{cases} F_0^D, & i = 0, j = 0 \\ \max \begin{cases} C_{i,j-1}^D + D_j + T_j^{D \rightarrow D} \\ F_j^M \end{cases}, & i = 0, j \in [1, N_p] \\ \max \begin{cases} C_{i-1,0}^I + I_{j,f(s_i)} + T_j^{I \rightarrow D} \\ T_0^{B \rightarrow M} \end{cases}, & i \in [1, N_s - 1], j = 0 \\ \max \begin{cases} C_{i-1,j-1}^M + M_{j,f(s_i)} + T_j^{M \rightarrow D} \\ C_{i-1,j}^I + I_{j,f(s_i)} + T_j^{I \rightarrow D} \\ C_{i,j-1}^D + D_j + T_j^{D \rightarrow D} \\ T_j^{B \rightarrow D} \end{cases}, & i \in [1, N_s - 1], j \in [1, N_p] \end{cases} \quad (3)$$

$$C_{i,j}^S = \begin{cases} \max \begin{cases} C_{i-1,j-1}^M + M_{j,f(s_i)} + T_j^{M \rightarrow E} \\ C_{i-1,j}^I + I_{j,f(s_i)} + T_j^{I \rightarrow E} \\ C_{i,j-1}^D + D_j + T_j^{D \rightarrow E} \end{cases}, & i \in [1, N_s - 1], j \in [1, N_p] \\ \max \begin{cases} C_{i-1,j-1}^M + M_{j,f(s_i)} + L_j^M \\ C_{i-1,j}^I + I_{j,f(s_i)} + L_j^I \\ C_{i,j-1}^D + D_j + L_j^D \end{cases}, & i = N_s, j \in [1, N_p] \\ \text{useless}, & i = 0 \text{ or } j = 0 \end{cases} \quad (4)$$

Running PfSearch, either version 2 or 3, first computes the cell matrix $C_{i,j}^\alpha$, $\alpha \in \{M, I, D, S\}$, then analyzes the latter looking for the best match i.e. the highest score $C_{k,l}^S$ where the range of l is in agreement with the type of alignment sought. For a *global* alignment, one enforces the match to end on the last profile position, hence $l = N_p$. On the other hand, a *local* alignment allows $l \in [0, N_p]$. It is worth mentioning that any type of alignment is also subject to the profile declaration and therefore it may well happen that the backward tracing encounters an entry point before reaching the profile's length. As an example, one can try to elucidate the path given in figure ??.

2 Repeat Decoder

The Repeat detector algorithm is a bit more complex than the standard PfTools one as it incorporates feedback loop to accomodate potential replication of the profile.

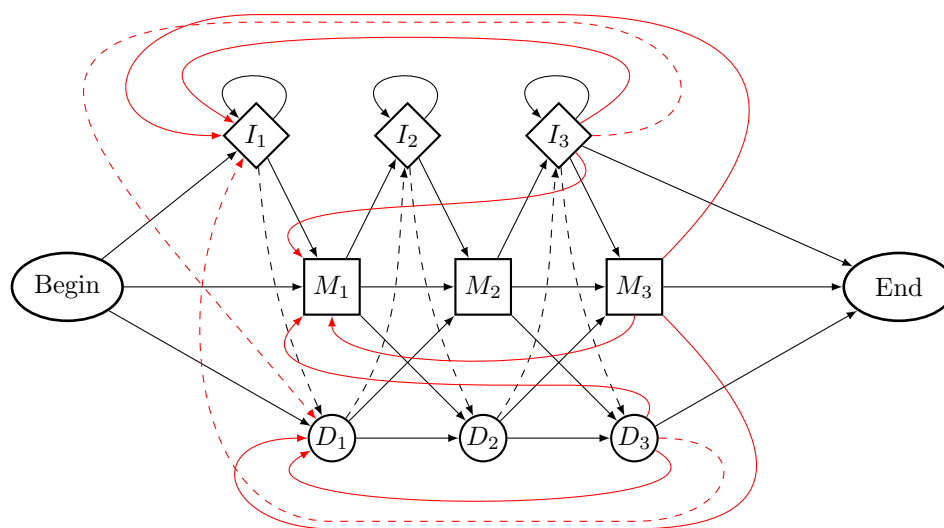


Figure 4: Simplified PfTools algorithm profile schema for a 3 sequence length profile. Simplifications arise from not showing all connections from Begin to nodes $[I_i, M_i, D_i]$ as well as nodes $[I_i, M_i, D_i]$ to End. It is worth noting the dashed lines that should not be used even though available. Indeed, mismatches should replace Insertion-Deletion or Deletion-Insertion paths thanks to extreme score loss.