

# Exercise Sheet 1 for Computational Biology (Part 1), SS 15

**Hand In:** Until Tuesday, 12.05.2015, 10:00 am, email to `r_muelle@cs...`, hand-in box in stairwell 48-6 or in lecture.

## Problem 1

2 + 3 points

- Give an infinite family  $(T_n)$  of texts with  $T_n \in \{a, b\}^{n-1}$  such that the number of nodes  $t_n$  of the corresponding simple suffix trees  $B_{T_n}$  is quadratic in  $n$ , i. e.,  $t_n = \Theta(n^2)$ .
- Give a second infinite family  $(T_n)$  of texts, for which the compact suffix trees  $IB_{T_n}$  have worst case size, i. e., the number of nodes of  $IB_{T_n}$  is maximal among all compact suffix trees for texts of the same size  $|T_n| = n$ . What is the worst case number of nodes?

## Problem 2

3 points

For two strings  $P \in \Sigma^m$ ,  $S \in \Sigma^n$  and a positive integer  $k$ , we say that  $P$  is a  $k$ -repeat in  $S$ , if there are exactly  $k$  (distinct) indices  $i_1, \dots, i_k$  such that  $S_{i_j, i_j+m-1} = P$  for  $j = 1, \dots, k$ .

Design an algorithm to find a shortest  $k$ -repeat in a string  $S \in \Sigma^n$ . The runtime should be in  $\mathcal{O}(n)$ .

## Problem 3

3 points

Design a linear time algorithm to compute the set of all *maximal repeats* of a text  $T$  along the lines given on pages 61ff of the German lecture notes.

More precisely, for every maximal repeat  $P$  of  $T \in \Sigma^n$ , your algorithm is supposed to output one pair of indices  $(i, j)$  and its length  $m = |P|$ , such that  $P$  is found at positions  $i$  and  $j$  in  $T$ :

$$T_{i, i+m-1} = T_{j, j+m-1} = P \quad \wedge \quad T_{i-1} \neq T_{j-1} \quad \wedge \quad T_{i+m} \neq T_{j+m}$$

where we set  $T_0 := \$ =: T_{n+1}$  for  $\$ \notin \Sigma$ . The running time should be in  $\mathcal{O}(n)$ .

**Problem 4**

4 points

For two strings  $S$  and  $T$  over alphabet  $\Sigma$ , we define the *overlap* of  $S$  and  $T$  as

$$ov(S, T) := \max\{|y| \mid y \in \Sigma^* \wedge \exists x, z \in \Sigma^+ : S = xy \wedge T = yz\} \quad (1)$$

Design an algorithm to compute *all* pairwise overlaps of a given set of strings  $\mathcal{T} = \{T^{(1)}, \dots, T^{(m)}\}$  over  $\Sigma$ , i. e. for all  $i, j \in [m]$ , compute  $ov(T^{(i)}, T^{(j)})$ . The running time of your algorithm should be in  $\mathcal{O}(n \cdot m)$ , where  $n := \sum_{i=1}^m |T_i|$  is the total length of all strings in  $\mathcal{T}$ .