

Algorithms in Genome Research  
Winter 2023/2024

Exercises

Number 3, Discussion: 2023 November 24

1. The *Shortest Common Superstring* problem (SCP) is defined as follows. Given a set of fragments  $\mathcal{F}$ , find a string  $S$  of minimal length such that  $f$  or its reverse complement  $\bar{f}$  is a substring of  $S$  for all  $f \in \mathcal{F}$ . Use the overlap graph to solve SCP for the fragments below.

Notes:

- For taking reverse complements into account, each fragment can be represented by two vertices corresponding respectively to its head and to its tail.
- For this dataset it suffices to consider only the edges corresponding to overlaps of length at least 3.

$$\begin{array}{ll} f_1 = \text{ATAT} & f_4 = \text{TATA} \\ f_2 = \text{TATT} & f_5 = \text{TTAT} \\ f_3 = \text{TTAT} & f_6 = \text{AATA} \end{array}$$

2. Consider the following set of reads, assuming that you know already that all of them originate from the same DNA strand.

$$\begin{array}{ll} r_1 = \text{ATCCA} & r_6 = \text{GCAAG} \\ r_2 = \text{AGAGC} & r_7 = \text{AGATC} \\ r_3 = \text{AAGAT} & r_8 = \text{TAGAG} \\ r_4 = \text{GAGCA} & r_9 = \text{AGAGC} \\ r_5 = \text{CCATA} & r_{10} = \text{GAGCA} \end{array}$$

- (a) Build the corresponding overlap graph with a minimum overlap of 2.
  - (b) Find a shortest common superstring for all reads. Is the coverage uniform? If not, find a layout with a more uniform coverage.
3. While in single-end sequencing, the sequencer reads a fragment from only one end to the other, in paired-end sequencing it reads the fragments from both ends. This gives a mate pair and a good estimation of the distance between them.

Mate-Pair Reads - 5' and 3'



How can mate pairs help in the assembling process?