

Series 2 - solution

Genomics and bioinformatics - Week 3 - October 2, 2012

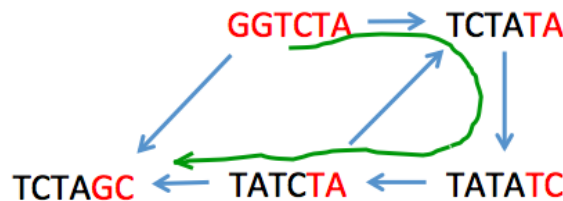
1 Overlap graphs

1. Just by looking manually at the overlaps (of size 4+):

GGTCTA
TCTATA
TATATC
TATCTA
TCTAGC
GGTCTATATCTAGC

2. The overlap (Hamiltonian) graph is build by taking the reads as vertices and adding an edge each time two reads have an overlap (of a chosen minimal size, here 4). Then find a path going through every vertex once and only once: this is called an Hamiltonian path. The contig is made of the first read plus, from each read in the path, its sub-sequence (in red below) that does not belong to the overlap. There may be several such paths; the number of them depends on the minimal overlap one chooses.

As in part 1, one finds GGTCTATATCTAGC.



3. To build the “De Bruijn” graph, one proceeds as follows:

- Choose an integer l (here $l = 4$ was given).
- Build S_{l-1} , the set of all *unique* $l - 1$ -mers that the reads contain:
 $S_{l-1} = \{GGT, GTC, TCT, CTA, TAT, ATA, ATC, TAG, AGC\}$
 These are the vertices of the graph.
- Build S_l , the set of all l -mers that the reads contain:
 $S_l = \{GGTC, GTCT, TCTA, TCTA, CTAT, TATA, TATA, ATAT, TATC, TATC, ATCT, TCTA, TCTA, CTAG, TAGC\}$

Those are the edges: GGTC binds GGT and GTC, etc.

