

Assembling genomes

Max Carter-Brown

Anglia Ruskin University, Wellcome Sanger
Institute
max.carter-brown@aru.ac.uk

Mark Blaxter

Wellcome Sanger Institute, University of
Edinburgh

Contents

1. Introduction	1
2. CAP3	2
3. The raw data	2
3.1. Understanding the data	3
4. Assembling your reads	3
4.1. Open a browser window and go to the PRABI CAP3 server	3
4.2. Submit your sequence data to the CAP3 server	3
4.3. Select "SUBMIT"	3
5. Examine your results	4
5.1. The "Contigs" file	4
Bibliography	4

This practical is assessed

A Canvas quiz has been set up for this practical. We will undertake the practical in examination conditions, and you will **not be able to talk during the session**. In addition, you must only have the tabs open on your computer which are relevant to the practical. If there are any other tabs open, or you are seen to be using ChatGPT or any other LLM, you will instantly fail with a score of zero marks.

Answer all of the questions in the yellow question boxes, these are equivalent to those on Canvas.

1. Introduction

You have the full two hours to complete this assessment, there should be plenty of time. Please read all the introductory material before starting the assessment.

In this practical assessment you will be taking a set of sequence reads derived from a mitochondrial genome, and using one of the many available genome assembly programs to "assemble" these sequence reads into a complete genome. The mitochondrial genome is circular (in animals) and 14000 base pairs in length. It encodes 2 ribosomal RNAs and 12 or 13 proteins (depending on which species is being examined). It also encodes 20 transfer RNAs (tRNAs).

We will look at *Brugia malayi*, which is a parasitic roundworm (nematode) which infects humans, amongst other animals. The *Brugia malayi* mitochondrial genome is a circular molecule, with 12 protein coding genes and two rRNAs, just under 13,600 base pairs in length. The sequencing reads were generated by the Blaxter lab in the Institute of Evolutionary Biology, University of Edinburgh, as part of the *Brugia malayi* genome project.

The reads are from Sanger sequencing (or dideoxy) technology. The reads are each 300-600 bases long, and are paired-end reads from cloned inserts in plasmid vectors. This is a fairly old sequencing technology - the one that was used to assemble the original human genome - but the data is good quality. The reads have been pre-trimmed “low quality” sequence.

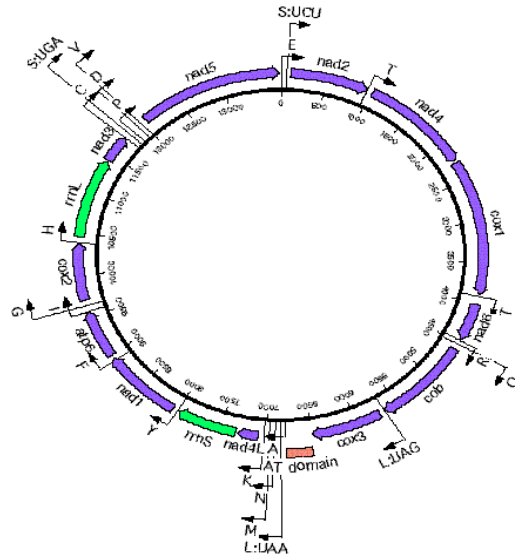


Figure 1: The *B. malayi* mitochondrial genome showing the position of the genes in the circular molecule.

2. CAP3

You will use the program “CAP3” (which stands for contig assembly program). CAP was developed by Xiaohu Huang. It was first published in 1992 [1]. A brief description of how the program works is given below.

CAP3 is an “alignment-overlap-consensus” assembler. It first compares all the sequences to each other to identify which could possibly derive from the same portion of the genome. It then develops an overlap map of the sequences, an estimate of how the sequences could be tiled together to make an assembly. If a particular sequence could overlap two others, and these overlaps are in conflict, the program uses the quality of the alignments (length, number of identities) to choose between them. Finally, CAP3 derives a consensus sequence from the overlapped sequences, and produces a set of reports.

It is an old, and fairly naive assembler by today’s comparison but it works well for the dataset that we have.

3. The raw data

You should now open your data file `BM_mt.fa`. It is a fasta file, an extremely common file format in bioinformatics which we have come across before. Open your file **in a text editor**, not Word. Word will add characters to the file which will make it unreadable by the assembler.

Questions

1. What defines the fasta format?
2. How many sequences are in the raw data file? Consider using `grep`, or another pattern matching tool.

3.1. Understanding the data

The sequences in the file are derived from performing sequencing on shotgun clones of the *B. malayi* genome. Each clone was named “Bmmt_cl_X” where X is a number. Each clone was sequenced twice, using primers at the start and the end of the insert. One sequence will represent the forward read (`_f`), and the other the reverse read (`_r`), which will be from the opposite side of the strand.

This type of sequencing is called paired-end sequencing. Forward and reverse sequences are not present for every clone, so some clones are represented by either `_f` or `_r` reads.

Thus, for example, “Bmmt_cl_123_f” and “Bmmt_cl_123_r” form a “read pair” from each end of a single clone.

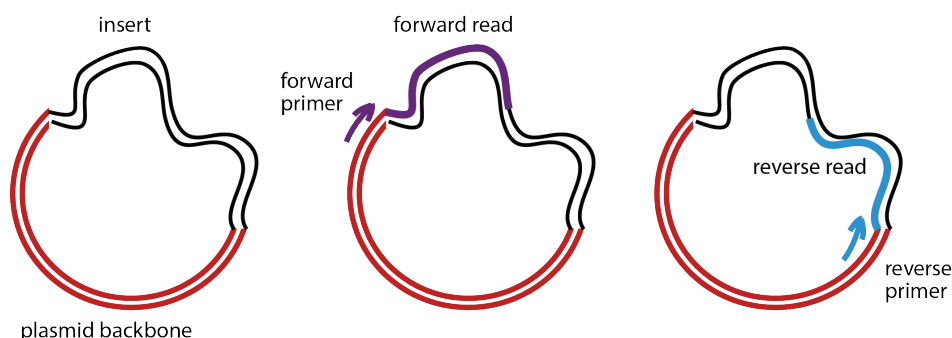


Figure 2: How reads are designated forward or reverse in the context of sequencing.

4. Assembling your reads

We will use an online service that offers access to the CAP3 assembly program. In the online version, a form is used to collect your data, and the webserver then formats and executes a command on the remote computer to run the program. Once it is finished, the webserver gives you access to the output files from the analysis.

4.1. Open a browser window and go to the PRABI CAP3 server

Open a link to the CAP3 server at PRABI-Doua (<http://doua.prabi.fr/software/cap3>) in a new browser window.

Further reading...

The PRABI-Doua (Pôle Rhône-Alpes de Bioinformatique Site Doua) World Wide Web server - developed at the Laboratory of Biometry and Evolutionary Biology and the Institute of Biology and Chemistry of Proteins - is dedicated to molecular biology and ecology.

4.2. Submit your sequence data to the CAP3 server

Select all of the sequence data in FASTA format from `BM_mt.fa` and paste it into the window “Enter your sequences in FASTA format”.

Leave the various analysis options at their default values.

4.3. Select “SUBMIT”

Your data will be sent to the PRABI remote computing facility which will carry out the assembly. Please be patient while this happens: do not keep clicking on the page while it is loading. This may take up to 5 minutes depending on the server load.

5. Examine your results

Several text outputs are produced.

1. **Contigs**
2. **Single sequences**
3. **Assembly details**
4. **Your sequence file**

These are all hyperlinks which link to text files. Save the files to your local workspace for examination.

5.1. The “Contigs” file

This contains the assembled sequences in fasta format. A “contig” is a set of contiguous sequences.

Questions

3. What is the total number of contigs produced by the assembler?
4. What is the length of the largest contig?
5. What is the total length of the assembly?

Note that each sequence line in the fasta is 60 bases long. You may wish to put into practice your BASH/shell/terminal skills. `grep` will be useful in counting contigs. `wc` will be useful in counting bases (once the headers are excluded). Of course you can count in any way you like - such as using Word tools.

Bibliography

- [1] X. Huang, “A contig assembly program based on sensitive detection of fragment overlaps,” *Genomics*, vol. 14, no. 1, pp. 18–25, 1992.