

Copyright is owned by the Author of the thesis. Permission is given for a copy to be downloaded by an individual for the purpose of research and private study only. The thesis may not be reproduced elsewhere without the permission of the Author.

EXPLORING BIOLOGICAL SEQUENCE SPACE

SELECTED PROBLEMS IN SEQUENCE ANALYSIS AND PHYLOGENETICS

A thesis presented in partial fulfilment
of the requirements for the degree of

Doctor of Philosophy
in Computational Biology

at Massey University

Bennet James McComish

2012

Copyright © 2012 by Bennet James McComish

Abstract

As the volume and complexity of available sequence data continues to grow at an exponential rate, the need for new sequence analysis techniques becomes more urgent, as does the need to test and to extend the existing techniques. These include, among others, techniques for assembling raw sequence data into usable genomic sequences; for using these sequences to investigate the evolutionary history of genes and species; and for examining the mechanisms by which sequences change over evolutionary time scales. This thesis comprises three projects within the field of sequence analysis.

- It is shown that organelle genome DNA sequences can be assembled *de novo* using short Illumina reads from a mixture of samples, and deconvoluted bioinformatically, without the added cost of indexing the individual samples. In the course of this work, a novel sequence element is described, that probably could not have been detected with traditional sequencing techniques.
- The problem of multiple optima of likelihood on phylogenetic trees is examined using biological data. While the prevalence of multiple optima varies widely with real data, trees with multiple optima occur less often among the best trees. Overall, the results provide reassurance that the value of maximum likelihood as a tree selection criterion is not often compromised by the presence of multiple local optima on a single tree.
- Fundamental mechanisms of mutation are investigated by estimating nucleotide substitution rate matrices for edges of phylogenetic trees. Several large alignments are examined, and the results suggest that the situation may be more complex than we had anticipated. It is likely that genome scale alignments will have to be used to further elucidate this question.

Acknowledgements

First and foremost I would like to thank my supervisor, David Penny, for encouraging me to do a PhD after many years away from science. Along with my co-supervisor, Mike Hendy, and despite the best efforts of the University, David has spent many years building a world-class research group conveniently located in my home town. My other co-supervisor, Lesley Collins, has helped knock this thesis into shape, and fought valiantly against my inclination toward understatement.

Thank you to all those who helped with code, data and other technical stuff: Patrick Biggs, Klaus Schliep, Trish McLenachan, Robin Atherton, Simon Hills, Judith Robins, Abby Harrison, Bojian Zhong, Pete Lockhart, Eric Bapteste and others. And to Tim White for the \LaTeX template I used to typeset this thesis. Also to the many colleagues who asked questions and made useful suggestions at talks and conferences.

I must thank the Allan Wilson Centre and the Institute of Molecular BioSciences for financial and logistical support, and especially Katrina Ross for having, as David puts it, “the terrible failing of always being on the side of the students.”

Thanks to all my friends, old and new, for making the PhD experience a social one as well as a scientific one. In particular, those who do crosswords at lunchtime (Gillian, Simon, Nick, Barbara) and those who play board games on Saturday evenings (some of the above plus Rogerio, Paul, Tim, Sylvia, Robin, Klaus, Atheer, Rick, Sam, Elsa and more). Another thank you goes to Klaus for putting me up and showing me the sights on my visits to Paris. Also to Paul and Carolyn, David and Helen, the Tappers and the Slees, for their generous help and hospitality during my visits to Melbourne.

Finally, I would like to thank my family: both my parents for always encouraging my scientific tendencies; my stepmother Heather for convincing me that doing a PhD would be a good idea; my Belgian extended family, especially Charles and Nanou, for their hospitality during my visits to Europe; and my daughter Petra for providing a good excuse to take a break every few months and for distracting me from my work. And lastly Barbara, for your love and support.

Contents

Abstract	iii
Acknowledgements	iv
1 Introduction	1
1.1 Background	2
1.1.1 Sequence data	2
1.1.2 Phylogenetic trees	3
1.1.3 Substitution models	4
1.2 Thesis outline	5
1.2.1 Assembly of mixed mitochondrial genomes	6
1.2.2 Multiple optima of likelihood	7
1.2.3 Investigating mutational mechanisms	8
1.2.4 Future directions	8
2 Assembly of mixed mitochondrial genomes	11
2.1 Preamble	11
2.2 Paper	12
2.3 Further sequence assembly work	27
2.3.1 Rats	27
2.3.2 Chloroplasts	27
2.3.3 More mixtures	28
2.4 Conclusions	30
3 Multiple optima of likelihood	33
3.1 Preamble	33
3.2 Introduction	36
3.3 Methods	39

3.3.1	Datasets	39
3.3.2	Computational experiments	40
3.4	Results	41
3.4.1	Chloroplast dataset	42
3.4.2	Mammalian mitochondrial dataset	44
3.4.3	Hepatitis B dataset	44
3.4.4	Prokaryote dataset	47
3.4.5	General results	50
3.5	Discussion	51
4	Investigating mutational mechanisms	55
4.1	Preamble	55
4.2	Introduction	58
4.3	Methods	61
4.3.1	Datasets	63
4.3.2	Visualisation	66
4.4	Results	68
4.5	Discussion	72
5	Future Directions	81
5.1	Assembly of mixed samples	81
5.2	Multiple optima of likelihood	82
5.3	Investigating mutational mechanisms	82
5.4	Summary	83
Appendix A	Contribution to publications	85
A.1	Assembly of mixed mitochondrial genomes	85
A.1.1	Author contributions	85
A.2	Multiple optima of likelihood	88
A.2.1	Author contributions	88
A.3	Karaka chloroplast genome	90
A.3.1	Author contributions	90

