

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.

EXTRACTING AND EXPLOITING SIGNALS IN GENETIC SEQUENCES

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

Doctor of Philosophy
in Mathematics

at Massey University

Walton Timothy James White
2011

Abstract

As DNA databases continue to grow at an exponential rate, the need for more efficient solutions to basic problems in computational biology grows ever more pressing. These problems range from the principal questions driving evolutionary science—How can we accurately infer the history of genes, individuals and species? How can we separate the signal from the noise in our data? How can we visualise that signal?—to the purely practical: How can we efficiently store all this data? With these goals in mind, this thesis mounts a computational combination attack on a variety of topics in bioinformatics and phylogenetics:

- A program is designed and implemented for solving the Maximum Parsimony problem—in essence, finding phylogenetic trees having the fewest mutations. This program generally outperforms existing highly optimised programs when using a single CPU, and unlike these earlier programs, offers highly efficient parallelisation across multiple CPUs for further speedup.
- A program is designed and implemented for compressing databases of DNA sequences. This program outperforms general-purpose compression by taking advantage of the special “treelike” structure of DNA databases, using a novel data structure, the “leaky move-to-front hashtable”, to achieve speed gains.
- A data visualisation technique is introduced that concisely summarises the “treelikeness” of phylogenetic datasets on a ternary plot. Each dataset is represented by a single point, allowing multiple datasets, or multiple treatments of a dataset, to be displayed on a single diagram.
- We demonstrate problems with a standard phylogenetic analysis methodology in which a single tree is assumed *a priori*. We argue for a shift towards network methods that can in principle reject the hypothesis of a single tree.
- Motivated by a phylogenetic problem, a fast new algorithm is developed for finding the mode(s) of a multinomial distribution, and an exact analysis of its complexity is given.

Acknowledgements

I would like to thank my supervisors, Mike Hendy, David Penny and Barbara Holland, for their support (both financial and moral), guidance and patience during my candidature. I feel extremely lucky to have had the chance to work with such all-round high-calibre people. Several years ago, in a combination attack of their own, these three gradually persuaded me (against my own (worse) judgment) to upgrade from a Masters degree to a PhD. I'm grateful that they knew me well enough, and thought it worthwhile, to keep nudging me until I realised what was good for me!

Many people have had a positive impact on my life during the course of my PhD. A somewhat biased random sample includes:

Past and present members of the Allan Wilson Centre at Palmerston North, whom I have found to be without exception friendly, intelligent and helpful.

My flatmates at the venerable 53 Te Awe Awe St. estate—Matt, Rachel, Michael, both Karens, Richard and Scott—for providing a thoroughly enjoyable place to live, where I always felt comfortable, *even* when Richard was around. As always, it seems fitting to mention that Scott knows a guy who once slept *on a door*.

Tony Dale of the BlueFern supercomputer group, for always responding quickly and helpfully when something went wrong.

The loose association of People Who Have Dinner and Play Board Games on Saturday Nights—to a first approximation, Rogerio, Klaus, Simon, Gillian, Bennet, Barbara and Robin—for tolerating my sense of humour, and occasionally affording me the chance to repay their gastronomic generosity with a few least significant bits of my own.

My good friend Sylvia, for many things, including for teaching me that people from very different backgrounds can find the same kinds of people extremely funny-looking.

Jing, for being an ideal officemate. Trish, for being a tireless force for good. Joy, for returning my wisecracks with interest.

Everyone over the years who helped me *test the graphics cards* on the computers in the Allan Wilson Centre.

And of course my family, whose love and support has been a constant in my life.

Contents

| | |
|--|-----------|
| Abstract | i |
| Acknowledgements | ii |
| 1 Introduction | 1 |
| 1.1 Overview | 1 |
| 1.2 Common Themes | 2 |
| 1.2.1 Strengthening the Fundamentals | 3 |
| 1.2.2 Information: A Common Currency | 3 |
| 1.2.3 Combination Attacks | 3 |
| 1.3 Background: Phylogenetic Inference | 7 |
| 1.3.1 Fundamentals | 7 |
| 1.3.2 Starting points for inference | 8 |
| 1.3.3 Algorithmic treebuilding approaches | 9 |
| 1.3.4 Optimality criteria | 9 |
| 1.3.5 Search methods | 11 |
| 1.3.6 Heuristic algorithms | 12 |
| 1.3.7 Statistical models | 12 |
| 1.3.8 Errors and consistency | 14 |
| 1.3.9 Maximum Parsimony and consistency | 14 |
| 1.4 Computational Complexity Primer | 15 |
| 1.4.1 NP Completeness | 16 |
| 2 Faster Exact Maximum Parsimony Search with XMP | 20 |
| 2.1 Introduction | 20 |
| 2.2 Correctness of Steiner Tree Lower Bounds for Ambiguous Nucleotides | 30 |
| 2.2.1 Without Ambiguous Nucleotides | 30 |

| | | |
|----------|---|-----------|
| 2.2.2 | With Ambiguous Nucleotides | 31 |
| 2.3 | Zero-length Edges | 32 |
| 2.3.1 | Sufficiency of Examining Fully Refined Trees | 33 |
| 2.3.2 | Avoiding the Deluge | 33 |
| 2.3.3 | Minimally Refined Trees | 35 |
| 2.3.4 | XMP Tree Representation | 35 |
| 2.3.5 | Identifying ml0 Edges | 37 |
| 2.3.6 | Contracting ml0 Edges | 37 |
| 2.3.7 | A Counterexample | 38 |
| 2.4 | Zharkikh's Rules | 39 |
| 2.5 | Visualising MPI Communication | 40 |
| 2.6 | Lower Bound Comparison | 41 |
| 3 | Compressing DNA Sequence Databases with coil | 45 |
| 3.1 | Introduction | 45 |
| 3.2 | Pentium IV Optimised <code>find_edges</code> | 61 |
| 3.3 | Sequence Buffering System | 63 |
| 3.4 | Data Compression as Quick and Dirty Science | 64 |
| 3.5 | General-purpose vs. Specialised Compression | 66 |
| 3.6 | Erratum | 66 |
| 4 | Treeness Triangles: Visualizing the Loss of Phylogenetic Signal | 67 |
| 4.1 | Introduction | 67 |
| 4.2 | Why does distance correction amplify residual signals? | 69 |
| 4.3 | Clarification | 70 |
| 5 | A Bias in ML Estimates of Branch Lengths in the Presence of Multiple Signals | 82 |
| 5.1 | Introduction | 82 |
| 5.2 | Testing Robustness | 83 |
| 5.2.1 | Does “the” internal edge of a mixture of two trees really exist? | 83 |
| 5.2.2 | Shared parameter values | 84 |
| 5.2.3 | The edges of a mixture model | 84 |
| 5.3 | Later Developments | 85 |

| | | |
|----------|--|------------|
| 5.4 | Connection to Multinomial Modes | 86 |
| 6 | A Fast and Simple Algorithm for Finding the Modes of a Multinomial Distribution | 91 |
| 6.1 | Introduction | 91 |
| 6.2 | Motivation and Connection to Phylogenetics | 91 |
| 6.2.1 | The Problem with PAUP* | 92 |
| 6.2.2 | Choosing a Best Representative | 93 |
| 6.2.3 | An Alternative to Least Squares | 94 |
| 6.3 | Alternate Proof of Correctness | 95 |
| 6.4 | Problem Instances | 97 |
| 7 | Conclusion | 105 |
| 7.1 | XMP | 105 |
| 7.2 | COIL | 107 |
| 7.3 | Treeness Triangles | 108 |
| 7.4 | ML Bias | 109 |
| 7.5 | Multinomial Modes | 110 |
| 7.6 | How Much Optimisation is the Right Amount? | 110 |
| 7.7 | Other Directions | 112 |
| 7.8 | Summary | 112 |
| A | Work Breakdown | 113 |
| A.1 | XMP | 113 |
| A.1.1 | Author Contributions | 113 |
| A.2 | COIL | 116 |
| A.2.1 | Author Contributions | 116 |
| A.2.2 | Previously Examined Work | 116 |
| A.3 | Treeness Triangles | 121 |
| A.3.1 | Author Contributions | 121 |
| A.4 | ML Bias | 123 |
| A.4.1 | Author Contributions | 123 |
| A.5 | Multinomial Modes | 125 |
| A.5.1 | Author Contributions | 125 |