# Computational Statistical Inference for Molecular Evolution and Population Genetics

A thesis submitted in partial fulfilment of the

requirements for the Degree of

Doctorate of Philosophy in Biological Sciences,

The University of Auckland,

New Zealand

June 2002

## 0.1 Abstract

This research aims to develop new methods and software for evolutionary inference. The focus will be on two challenges that analysis of molecular data in the genomic age provides: (i) *measurably evolving populations* and (ii) evolution of RNA secondary structure. Molecular sequence data is increasing in length, and also acquiring a depth in the time dimension (for example, HIV-1, human influenza A, and ancient mtDNA). This has provided an innovative research direction, for which explicit evolutionary inference methods are required. The first aim of this research is to provide new statistical methods and new bioinformatic tools (software packages) to assist in tackling this new problem in evolutionary biology. Both maximum likelihood and Bayesian inference methods are developed for the purpose of estimating substitution rates and concerted changes in the substitution rate. In addition, with the rapid succession of newly sequenced full genomes, researchers can no longer use simple molecular sequence similarity to infer homology. Knowledge of molecular structure needs to be incorporated into evolutionary inference methods. The evolutionary relationship between sequence and structure is still poorly understood and the new wealth of data provides an exciting opportunity to guide theoretical developments. The second major objective of this research is to use the wealth of sequence data available to explore the role and impact of RNA secondary structure on evolution. To this end, empirical studies and simulations are undertaken to explore the role of RNA secondary structure in the evolution of 16S-like rRNA-encoding genes. Finally the inference of spatially resolved populations from gene sequences is briefly investigated.

This research project has both computational and conceptual objectives. In both cases, the concrete result of these objectives will be new statistical models and computer software for evolutionary inference and a better understanding of the action of molecular and population processes during evolution.

## 0.2 Acknowledgements

# 0.3 Table of Contents

## 0.4 List of Abbreviations

| | |
|---|---|
| 3D | three-dimensional |
| AIDS | acquired immunodeficiency syndrome |
| bp | base pairs |
| BP | before present |
| BSC | biological species concept |
| CPU | central processing unit |
| DNA | deoxyribonucleic acid |
| EM | expectation-maximization algorithm |
| ESS | effective sample size |
| F81 | Felsenstein 1981 (model of substitution) |
| GPL | GNU public licence |
| GTR | general time-reversible (model of substitution) |
| HCV | Hepatitis C virus |
| HIMDU | hairpin, internal bulge, multi-stem loop, downstream-paired and upstream-paired (model of structure) |
| HIV-1 | human immunodeficiency virus, subtype 1 |
| HKY | Hasegawa, Kishino and Yano (model of substitution) |
| HOM | homogeneous model of structure (i.e. no structure) |
| HPD | highest posterior density |
| HVR1 | hyper-variable region 1 (of the mitochondrial control region) |
| IACT | integrated autocorrelation time |
| indels | insertion-deletion events |
| JC | Jukes-Cantor (model of substitution) |
| LGPL | lesser GPL |
| LS | least-squares |
| MCMC | Markov chain Monte Carlo |
| MEP | measurably evolving population |
| ML | maximum likelihood |
| MRCA | most recent common ancestor |
| MRDT | multiple rate dated-tips (model of mutation rate) |
| mtDNA | mitochondrial DNA |
| NIH | National Institutes of Health |
| NJ | neighbour joining (method of phylogenetic reconstruction) |
| NNI | nearest neighbour interchange (method of branch swapping) |
| PAL | Phylogenetics Analysis Library |
| RNA | ribonucleic acid |
| rRNA | ribosomal RNA |
| SDI | symmetric difference index |
| SPR | subtree-prune and reattachment (method of branch swapping) |
| SR | single rate (model of mutation rate) |
| SRDT | single rate dated tips (model of mutation rate) |
| sUPGMA | serial-sample UPGMA |
| tRNA | transfer RNA |
| UP | unpaired/paired model of structure |
| UPGMA | unweighted pair-group method using arithmetic means |
| WPGMA | weighted pair-group method using arithmetic means |

# 0.5 List of symbols and functions

Unless defined otherwise in the text, these symbols are defined as in the table below.

| | |
|---|---|
| ~ | distributed as |
| $\wedge$ | logical 'and' operation |
| $\vee$ | logical 'or' operation |
| $\alpha$ | the shape parameter of the gamma distribution of rate heterogeneity among sites. |
| $\delta$ | divergence (measured in substitutions/site) |
| $E_g$ | a set of edges defining a bifurcating tree |
| $\text{Exp}(x)$ | exponentially distribution with a mean of $x$ |
| $g$ | a genealogy $= (E_g, t_Y)$ |
| $\kappa$ | kappa, the ratio between the instantaneous rate of a particular transition and the instaneous rate of a particular transversion. |
| $\mu$ | mutation rate |
| $N_A$ | ancestral effective population size |
| $N_C$ | current effective population size |
| $N_e$ | effective population size |
| $\pi$ | equilibrium base frequencies |
| $Q$ | the instantaneous substitution rate matrix |
| $\Theta$ | intra-specific diversity: $2N_e\mu$ for haploid, $4N_e\mu$ for diploid |
| $r$ | exponential growth rate |
| $R$ | relative rate matrix |
| $t_I$ | the set of times/ages of the leaves of a genealogy. |
| $t_{MRCA}$ | time to the most recent common ancestor |
| $t_{root}$ | synonym for $t_{MRCA}$ |
| $t_Y$ | the set of times/ages of the ancestral nodes of a genealogy. |
| $\text{Unif}(x, y)$ | uniform distribution with a lower limit of $x$ and an upper limit of $y$ |
| $\omega$ | mutation rate (in mutations per site per calendar unit) |
| $Z$ | normalizing constant |

## 0.6   List of Tables

# 0.7   List of Figures