# From Species To Languages
*A phylogenetic approach to human prehistory*

QUENTIN DOUGLAS ATKINSON

A thesis submitted in fulfillment of the requirements
for the degree of Doctor of Philosophy in Psychology,
Department of Psychology,
The University of Auckland, 2006.

# ABSTRACT

Languages, like species, evolve. Just like biologists, historical linguists infer relationships between the lineages they study by analysing heritable features. For linguists, these features can be words, grammar and phonemes. This linguistic evidence of descent with modification plays an important role in our understanding of human prehistory. However, conventional methods in historical linguistics do not employ an explicit optimality criterion to evaluate evolutionary language trees. These methods cannot quantify uncertainty in the inferences nor provide an absolute chronology of divergence events. Previous attempts to estimate divergence times from lexical data using glottochronological methods have been heavily criticized, particularly for the assumption of constant rates of lexical replacement. Computational phylogenetic methods from biology can overcome these problems and allow divergence times to be estimated without the assumption of constant rates. Here these methods are applied to lexical data to test hypotheses about human prehistory. First, divergence time estimates for the age of the Indo-European language family are used to test between two competing theories of Indo-European origin – the Kurgan hypothesis and the Anatolian farming hypothesis. The resulting age estimates are consistent with the age range implied by the Anatolian farming theory. Validation exercises using different models, data sets and coding procedures, as well as the analysis of synthetic data, indicate these results are highly robust. Second, the same methodology was applied to Mayan lexical data to infer historical relationships and divergence times within the Mayan language family. The results highlight interesting uncertainties in Mayan language relationships and suggest that the family may be older than previously thought. Finally, returning to biology, similar tree-building and model validation techniques are used to draw inferences about human origins and dispersal from human mitochondrial DNA sequence data. These analyses support a human origin 150,000-250,000 years ago and reveal time dependency in rates of mitochondrial DNA evolution. Population size estimates generated using a coalescent approach suggest a two-phase human population expansion from Africa. Potential correlations between human genetic and linguistic diversity are highlighted. I conclude that there is much to be gained by linguists and biologists using the same methods and speaking the same language.

# ACKNOWLEDGEMENTS

First and foremost I would like to thank Dr Russell Gray - I could not have asked for a better supervisor. Thank you for your inspiration, motivation, support and guidance, and for sharing your skills as a researcher. It has been a wonderful four years working with you.

To my office buddy, Simon Greenhill, thank you for the good company and for all your help and advice. You are a scholar and a gentleman…and a viable alternative to Google. I am also very grateful to Professor Lyle Campbell at the University of Utah for his expert help and generous hospitality. To Elisabeth Norcliffe, thank you for sharing your knowledge of Mayan languages and the Ecuadorian hat-making industry. To Dr Geoff Nicholls, thank you for writing such a great program and for our lengthy sessions at the whiteboard. Thanks to David Welch for all your time and effort implementing a GUI that I could drive (and for that free trip to Switzerland). Thanks also to Dr Alexei Drummond for helping me tame the BEAST and take over the world and to Marcel van de Steeg for his valiant attempts to keep the recalcitrant "cluster" running. Thank you also to Professor Mike Corballis for providing sage advice when needed and to Professor Mark Pagel for problem solving when needed. I would also like to thank Scott Allan, Bob Blust, Lounès Chikhi, Penny Gray, Roger Green, Jeff Hamm, Niki Harré, John Huelsenbeck, Mark Liberman, David Penny, Allan Rodrigo, Fredrik Ronquist, Michael Sanderson and Stephen Shennan for useful advice and/or comments on manuscripts. To Mum and Dad, I really appreciate all your support over the years. And I owe an especially big thank you to Emma for so much "souper" advice and encouragement and for keeping me sane and happy.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES