

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.

# An Information Theoretic Approach to Language Relatedness

A dissertation submitted in partial fulfilment of the requirements  
for the degree of Doctor of Philosophy in Information Systems  
at Massey University.

**Anand Venkt Raman**  
**1997**

This documentation was prepared from a camera ready copy made by the author using  $\LaTeX$ . The International Phonetic Alphabet characters were produced using the  $\TeX$  font (wsuipa) created by Washington State University. The graphs were produced using `gnuplot`. PFSA diagrams were produced using a modified version of Jos Van Einjdhoven's `graphplace` program. Other figures were produced using the `xfig` program. The final Postscript file was generated using `dvips`, printed on a Hewlett-Packard Laserjet printer and photoreproduced and bound by Wills Bookbinding and Printing, 6 Dahlia Street, Palmerston North.

The software used to produce the results for this work are available through anonymous ftp from the URL <ftp://fims-ftp.massey.ac.nz/pub/ARaman>. Bug reports may be sent to [A.Raman@massey.ac.nz](mailto:A.Raman@massey.ac.nz). Work done using this software must cite either Raman and Patrick (1997d) or this dissertation.

## Abstract

This dissertation examines the prospect of applying information theoretic principles to help solve problems in historical linguistics. The Minimum Message Length principle attributed to Chris Wallace (similar to the Minimum Description Length principle of Jorma Rissanen) is used to judge the goodness of hypotheses in the field of historical linguistics. The idea is that theories that require a shorter message to describe with their data are better than those that require long messages.

Work in collecting the linguistic data tracing the derivation of some 2714 words in Modern Cantonese and Modern Beijing from their forms in a reconstruction of Middle Chinese is described as also is the work in transforming this data into a format suitable for use with software developed for this project.

Heuristics for inferring Probabilistic Finite State Automata (PFSA<sup>1</sup>) from such data are reviewed and some new heuristics are introduced. These are then applied to training data and benchmark results presented.

Finally, the inference process is applied to the actual linguistic data which allows a conjecture regarding a relative closeness of the Chinese dialects to their reconstructed ancestor to be formed.

---

<sup>1</sup>In this dissertation, the abbreviation PFSA has been used to denote both the singular and plural of these machines, the "A" in PFSA being understood to represent both *Automaton* and *Automata*.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Statement . . . . .	1
1.2	Linguistic Preliminaries . . . . .	2
1.2.1	Phones, Phonemes, Allophones and Morphemes . . . . .	2
1.2.2	Contrast, Minimal Pairs and Distributions . . . . .	3
1.2.3	Language Change . . . . .	5
1.2.4	Bleeding and Feeding . . . . .	6
1.3	Language Reconstruction . . . . .	7
1.3.1	The Comparative Method . . . . .	7
1.3.2	Internal Reconstruction . . . . .	9
1.4	Language Similarities and Relatedness . . . . .	10
<b>2</b>	<b>Quantitative Methods in Historical Linguistics</b>	<b>12</b>
2.1	Early Work . . . . .	12
2.2	Glottochronology . . . . .	14
2.2.1	Description . . . . .	15
2.2.2	Example . . . . .	16
2.2.3	Criticisms of glottochronology . . . . .	17
2.3	The Direction of This Thesis . . . . .	18
<b>3</b>	<b>Similarity Measures – An Example</b>	<b>19</b>
3.1	Overview . . . . .	19
3.2	Modelling *LNWG to OF . . . . .	20
3.3	Modelling *LNWG to OHG . . . . .	23
3.4	Discussion . . . . .	24
<b>4</b>	<b>Data Collection</b>	<b>32</b>
4.1	Introduction . . . . .	32
4.2	Original Sources and Reliability . . . . .	32
4.2.1	The <i>Qie yun</i> . . . . .	33
4.2.2	Reliability . . . . .	34
4.3	<i>Zihui</i> and DOC . . . . .	35
4.4	Motivation . . . . .	37
4.5	Chen76 and CN84 . . . . .	39
4.6	Relative Chronology (RC) . . . . .	41
4.7	Phonotactic Conditions and Allophonic rules . . . . .	44
4.8	Automatic Derivations . . . . .	47
4.9	Dealing with Exceptions . . . . .	47
4.9.1	Exceptions to the RC . . . . .	49
4.9.2	Undocumented Changes . . . . .	49
4.9.3	Unapplied Rules . . . . .	49

4.10	Status of Allophonic Changes . . . . .	51
4.11	Consistency of Presented Data . . . . .	52
<b>5</b>	<b>The MML Principle</b>	<b>55</b>
5.1	Modelling Inductive Hypotheses . . . . .	55
5.2	Information Measures and MML codes . . . . .	59
5.3	MML, Bayesianism and MDL . . . . .	61
5.4	Encoding Deterministic PFSA . . . . .	63
5.5	Encoding Non-Deterministic PFSA . . . . .	66
<b>6</b>	<b>The Sk-strings Method of Inferring PFSA</b>	<b>73</b>
6.1	Introduction . . . . .	73
6.2	Inducing the Gaines Machine . . . . .	74
6.3	Minimum Message Length . . . . .	76
6.4	The K-tails Approach . . . . .	81
6.5	The Sk-strings Approach . . . . .	83
6.6	Tractability Considerations . . . . .	86
6.6.1	Method A - Build strings in decreasing order of probability . . .	87
6.6.2	Method B - Reject strings less than a certain probability . . . .	88
6.6.3	Method C - Reject relatively improbable strings . . . . .	88
6.7	Results . . . . .	89
6.8	Conclusion . . . . .	96
<b>7</b>	<b>A Heuristic Using Cross-Entropy for Sk-Strings</b>	<b>101</b>
7.1	Introduction . . . . .	101
7.2	A Cross-Entropic Measure . . . . .	102
7.3	Implementation . . . . .	105
7.4	Results . . . . .	107
7.5	Conclusion . . . . .	110
<b>8</b>	<b>Beam Search and Simulated Beam Annealing</b>	<b>115</b>
8.1	Introduction . . . . .	115
8.2	Beam Search . . . . .	115
8.2.1	The Beam Search Algorithm for PFSA Inference . . . . .	116
8.2.2	Algorithm Complexity . . . . .	117
8.3	Simulated Annealing . . . . .	118
8.3.1	Simulated Annealing with a Beam . . . . .	120
8.4	Results and Discussion . . . . .	123
8.5	Discussion and Conclusion . . . . .	126
<b>9</b>	<b>Results and Conclusion</b>	<b>131</b>
9.1	Program outputs . . . . .	131
9.2	Discussion . . . . .	133
9.2.1	Status of exceptions . . . . .	134
9.2.2	Chen76 and CN84 on exceptions . . . . .	137
9.2.3	Relation between the PFSA and RC diagrams . . . . .	139
9.2.4	Significance of diachronic regularities in MB and MC . . . . .	140
9.2.5	Glottochronology revisited . . . . .	141
9.3	Conclusion and Prospects . . . . .	142

<b>A</b>	<b>Guide to Using the Programs</b>	<b>147</b>
A.1	Introduction . . . . .	147
A.2	Canon . . . . .	149
A.3	Genstr . . . . .	150
A.4	Randpfsa . . . . .	150
A.5	Skstr, Ktail, Beams and Simba . . . . .	150
A.6	Dfa . . . . .	151
A.7	Grafit and Groupnodes . . . . .	151
A.8	DISCLAIMER . . . . .	152
<b>B</b>	<b>Reduced PFSA</b>	<b>154</b>
B.1	PFSA for SMC>MC (Diachronic only) . . . . .	154
B.2	PFSA for SMC>MC (Diachronic and Allophonic) . . . . .	157
B.3	PFSA for SMC>MB (Diachronic only) . . . . .	160
B.4	PFSA for SMC>MB (Diachronic and Allophonic) . . . . .	163
<b>C</b>	<b>From Simplified Middle Chinese to Modern Beijing</b>	<b>166</b>
<b>D</b>	<b>From Simplified Middle Chinese to Modern Cantonese</b>	<b>204</b>

# List of Tables

1.1	Example illustrating bleeding and feeding relationships among phonological rules. . . . .	7
2.1	An example illustrating common cognate scores for glottochronology. . .	16
3.1	Diachronic rules for converting *LNWG → OF. . . . .	21
3.2	The sequences of diachronic rules to derive 20 words from *LNWG to OF.	23
3.3	Diachronic rules for converting *LNWG → OHG. . . . .	24
3.4	The sequences of diachronic rules to derive 20 words from *LNWG to OHG. . . . .	25
4.1	An example of the Relative Chronology in operation. . . . .	44
4.2	Allophones of ə; Operation of the *CHAMEL PC. . . . .	45
4.3	An example of an allophonic rule in feeding relationship with a diachronic rule. . . . .	51
4.4	An extract from Appendix D to illustrate the actual presentation of the data. . . . .	53
5.1	A comparison of MMLs for 1 to 5 state PFSA explaining Andreae's rigged Casino. . . . .	70
6.1	MMLs of 1000 inferred automata using the AND heuristic for tail sizes of 1–10 and agreement percentages of 1–100 (Training set 1). . . . .	90
6.2	MMLs of 1000 inferred automata using the OR heuristic for tail sizes of 1–10 and agreement percentages of 1–100 (Training set 1). . . . .	90
6.3	MMLs of 1000 inferred automata using the LAX heuristic for tail sizes of 1–10 and agreement percentages of 1–100 (Training set 1). . . . .	92
6.4	MMLs of 1000 inferred automata using the STRICT heuristic for tail sizes of 1–10 and agreement percentages of 1–100 (Training set 1). . . .	92
6.5	MMLs of 1000 inferred automata using the AND heuristic for various tail sizes and agreement percentages (Training set 2). . . . .	92
6.6	MMLs of 1000 inferred automata using the OR heuristic for various tail sizes and agreement percentages (Training set 2). . . . .	93
6.7	MMLs of 1000 inferred automata using the LAX heuristic for various tail sizes and agreement percentages (Training set 2). . . . .	93
6.8	MMLs of 1000 inferred automata using the STRICT heuristic for various tail sizes and agreement percentages (Training set 2). . . . .	93
6.9	The number of times (out of 100 training sets) that each tail size succeeded in inferring the minimum automaton using K-tails and the various sk-strings heuristics. . . . .	94
6.10	Success rates of the k-tails and 4 sk-strings heuristics in inferring the best PFSA out of 100 small test cases. . . . .	95

6.11	Success rates of the k-tails and 4 sk-strings heuristics in inferring the best PFSA out of 100 large test cases. . . . .	97
7.1	Success rates of the k-tails and 6 sk-strings heuristics (including XENTROPIC and VARDIST) in inferring the best PFSA out of 100 small test cases . . . . .	111
7.2	Success rates of the k-tails and 6 sk-strings heuristics (including XENTROPIC and VARDIST) in inferring the best PFSA out of 100 large test cases . . . . .	112
8.1	The number times (out of 100 training sets) that each beam size (1-10) succeeded in inferring the minimum automaton using the beam search and simba search methods. . . . .	125
9.1	MMLs for the canonical PFSA for Middle Chinese to Modern Cantonese and Modern Beijing respectively. . . . .	131
9.2	MMLs for the canonical PFSA for Middle Chinese to Modern Cantonese and Modern Beijing respectively using Formula (5.4). . . . .	132
9.3	MMLs for the reduced PFSA for Middle Chinese to Modern Cantonese and Modern Beijing respectively. . . . .	132
9.4	MMLs for the PFSA for Middle Chinese to Modern Cantonese and Modern Beijing respectively further reduced with beam search. . . . .	133
9.5	20 most frequent exceptions in the phonology of Modern Beijing. . . . .	136
9.6	20 most frequent exceptions in the phonology of Modern Cantonese. . . . .	137
9.7	20 most frequent rules in the phonology of Modern Beijing and Modern Cantonese. . . . .	138

# List of Figures

3.1	Canonical PFSA representing the derivation of 20 *LNWG words into OF. . . . .	28
3.2	Canonical PFSA representing the derivation of 20 *LNWG words into OHG. . . . .	29
3.3	A generalisation of the PFSA in Figure 3.1 (*LNWG>OF). . . . .	30
3.4	A generalisation of the PFSA in Figure 3.2 (*LNWG>OHG). . . . .	31
4.1	Chen76 Relative Chronology for *SMC>MB. . . . .	41
4.2	Chen76 Modified Relative Chronology for *SMC>MB: RC25 (SHARP > PROCOPE) orders SHARP before PROCOPE. . . . .	42
4.3	CN84 Relative Chronology for *SMC>MC. . . . .	43
5.1	An encoding of 4 events (ABCD) <sup>n</sup> . . . . .	57
5.2	An coding of the 4 (ABCD) <sup>n</sup> with an exception B'. . . . .	58
5.3	The events (ABCD) <sup>n</sup> encoded in the form of a PFSA. . . . .	58
5.4	The events (ABCD) <sup>n</sup> with exception B' encoded in the form of a PFSA. . . . .	59
5.5	An Illegal Automaton. . . . .	65
5.6	1-state PFSA explaining Andrae's rigged Casino. . . . .	70
5.7	2-state PFSA explaining Andrae's rigged Casino. . . . .	71
5.8	3-state PFSA explaining Andrae's rigged Casino. . . . .	71
5.9	4-state PFSA explaining Andrae's rigged Casino. Note that the output is completely predictable in states 1 and 2. . . . .	72
5.10	5-state PFSA explaining Andrae's rigged Casino. . . . .	72
6.1	Gaines' 4 state machine. . . . .	75
6.2	Two automata generating the same language: a*. . . . .	76
6.3	4 theories for the event sequence "abc!" seen once. . . . .	77
6.4	4 theories for the event sequence "abc!" seen 100 times. . . . .	78
6.5	2 candidate states for merging in a PFSA. . . . .	82
6.6	K-tails and sk-strings reductions of a canonical machine. . . . .	83
6.7	Test machine used to generate strings for training set 2. . . . .	91
6.8	MMLs of 1000 inferred automata using the AND/OR heuristics for tail sizes of 1-10 and agreement percentages of 1-100 (Training set 1). See also Tables 6.1 and 6.2. . . . .	94
6.9	MMLs of 1000 inferred automata using the LAX/STRICT heuristics for tail sizes of 1-10 and agreement percentages of 1-100 (Training set 1). See also Tables 6.3 and 6.4. . . . .	95
6.10	MMLs of 1000 inferred automata using the AND/OR heuristics for tail sizes of 1-10 and agreement percentages of 1-100 (Training set 2). See also Tables 6.5 and 6.6. . . . .	96

6.11	MMLs of 1000 inferred automata using the LAX/STRICT heuristics for tail sizes of 1–10 and agreement percentages of 1–100 (Training set 2). See also Tables 6.7 and 6.8. . . . .	97
6.12	The k-tail method compared with sk-strings (s=50%) using 100 random training sets generated by small random automata (Best machines inferred using any heuristic and tail sizes from 1 to 10). . . . .	98
6.13	The k-tail method compared with sk-strings (s=50%) using 100 random training sets generated by large random automata (Best machines inferred using any heuristic and tail sizes from 1 to 5). . . . .	99
7.1	MMLs of 1000 inferred automata using the XENTROPIC heuristic for tail sizes 1–10 and agreement percentages 1–100 (Training set 1). . . . .	108
7.2	MMLs of 1000 inferred automata using the VARDIST heuristic for tail sizes 1–10 and agreement percentages 1–100 (Training set 1). . . . .	109
7.3	MMLs of 1000 inferred automata using the XENTROPIC heuristic for tail sizes 1–10 and agreement percentages 1–100 (Training set 2). . . . .	110
7.4	MMLs of 1000 inferred automata using the VARDIST heuristic for tail sizes 1–10 and agreement percentages 1–100 (Training set 2). . . . .	111
7.5	The sk-strings (s=50%) XENTROPIC heuristic compared with the best of 4 sk-strings heuristics and the k-tails method using 100 random training sets generated by small random automata. . . . .	112
7.6	The sk-strings (s=50%) XENTROPIC heuristic compared with the best of 4 sk-strings heuristics and the k-tails method using 100 random training sets generated by large random automata. . . . .	113
7.7	The sk-strings (s=50%) VARDIST heuristic compared with the best of 4 sk-strings heuristics and the k-tails method using 100 random training sets generated by small random automata. . . . .	113
7.8	The sk-strings (s=50%) VARDIST heuristic compared with the best of 4 sk-strings heuristics and the k-tails method using 100 random training sets generated by large random automata. . . . .	114
8.1	Performance of the beam search and simba search procedures on Training set 1 from Chapter 6. The plot shows the MMLs of 1000 machines induced against various beam sizes or tail sizes (for the 6 sk-strings heuristics). . . . .	123
8.2	Performance of the beam search and simba search procedures on Training set 2 from Chapter 6. The plot shows the MMLs of 1000 machines induced against various beam sizes or tail sizes (for the 6 sk-strings heuristics). . . . .	124
8.3	Beam search compared with sk-strings (s=50%) for 100 training sets generated from small random automata. . . . .	126
8.4	Simba search compared with sk-strings (s=50%) for 100 training sets generated from small random automata. . . . .	127
8.5	Beam search compared with sk-strings (s=50%) for 100 training sets generated from larger random automata. . . . .	128
8.6	Simba search compared with sk-strings (s=50%) for 100 training sets generated from larger random automata. . . . .	129
8.7	Seeded beam search compared with sk-strings (s=50%) and unseeded beam search for 100 training sets generated from small random automata (Best of tail sizes 1 to 10 and beam sizes 1 to 10). . . . .	130
9.1	Reduced PFSA for the diachronic phonology from Middle Chinese to Modern Beijing (Allophonic detail excluded). . . . .	145

9.2	Reduced PFSA for the diachronic phonology from Middle Chinese to Modern Cantonese (Allophonic detail excluded). . . . .	146
A.1	Andreae's 5 state machine without vertical node grouping . . . . .	152
A.2	Andreae's 5 state machine with nodes 0, 2 and 4 grouped vertically . . .	153

# Preface

In late 1993, Professor Jon Patrick, who was then just about to leave Deakin University in Melbourne for Massey gave a seminar here on some computational work he was doing in Basque linguistics. I had at that time just joined the Computer Science department as an assistant lecturer and was avidly looking for a PhD project to begin. Like almost everybody else I have met, I also had a compelling interest in historical linguistics and had my own opinion of how it should be done. I sent a brief CV to Jon asking if he would like to take on a doctoral student in his Basque project after his arrival and the answer soon came back in the affirmative. In the beginning, though, our intention was quite different, and much more grandiose and ambitious — we were hoping to assess whether or not the linguistic isolate Basque was related to the Dravidian language Tamil, which by the way, I happen to be very fluent in.

As time went on, however, the focus of the project changed rather dramatically and eventually became much more realistic. It narrowed down into developing a technique for finding distance measures between related natural languages to aid linguists in their task of subgrouping. We decided, wisely perhaps, that bold proposals in the linguistic domain are best made by linguists. We also decided, with the advice of my second supervisor, Dr John Newman in the Department of Linguistics and Second Language Teaching, to use Chinese data for the project. This data had been provided to John courtesy of Professor William S-Y. Wang at the University of California, Berkeley.

The field of computational linguistics is such that almost all Computer Scientists seem to have their own unique notions of how one should apply computational methods to historical data. While there may be merits and demerits to each such approach, this dissertation is not an evaluation of their relative goodness. It is a detailed account of one particular approach — the one described here — and how it can be used to provide linguists with a tool for effective subgrouping of languages.

## Overview of contents

Chapter 1 gives a brief introduction to the problem and describes some linguistic terms and concepts for the benefit of the reader unacquainted with them. Chapter 2 is a brief survey of previous work in the field of linguistics which has been pursued along similar lines. An example of the methodology applied to a toy problem is then provided in Chapter 3. Chapter 4 describes the procedure by which the data for this project was collected and the motivation for using it. The MML criterion, which is central to this project, is introduced and described in Chapter 5, which is followed by Chapters 6–8 where methods to infer structure from the data using the MML criterion are looked at. Finally, Chapter 9 gives the results obtained, discusses them and describes the prospects for future work in this area.

## Publications

Parts of this dissertation have been published in various places during the course of working on this project. The following is a list of them, with brief notes on which chapters they refer to.

Raman, A. V. and J. D. Patrick (1997). Linguistic similarity measures using the Minimum Message Length principle. In R. D. Blench and M. Spriggs (Eds.), *Archaeology and language I: Theoretical and methodological orientations*, pp.260–277, London: Routledge. This paper was also read at the WAC-3 Conference, New Delhi, December 1994. Material from this paper can be found in Chapters 2 and 3.

Raman, A. V. and J. D. Patrick (1997). The sk-strings method for inferring PFSA. In *Proceedings of the workshop on automata induction, grammatical inference and language acquisition at the 14th international conference on machine learning — ICML-97*, Nashville, Tennessee, (in press). Material from this paper can be found in Chapter 6.

Raman, A. V. and J. D. Patrick (1997). A heuristic using cross-entropy for sk-strings In *Technical Report TR 1/97*, Information Systems Department, Massey University. Material from this report can be found in Chapter 7.

Raman, A. V. and J. D. Patrick (1997). Beam search and simba search for PFSA inference, In *Technical Report TR 2/97*, Information Systems Department, Massey University. Material from this report can be found in Chapter 8.

Raman, A. V., J. Newman and J. D. Patrick (1997). A complexity measure for diachronic Chinese phonology. In *Proceedings of the SIGPHON97 workshop on computational phonology*, Madrid (in press). Material from this paper can be found in Chapter 9.

## An explanation for some typesetting decisions

I feel compelled to say a few words about the typesetting process this dissertation has been through. I initially started writing the thesis using Microsoft Word on my Apple Macintosh, but found its performance unsatisfactory. Moreover, all my programs to manipulate the data were being run under Unix. What I wanted was a single command that could automatically take the ASCII outputs from my various experiments, turn them into tables, graphs and pictures and insert them in the right places in the thesis and also generate my bibliography neatly.

I thus deemed it best to migrate the work to  $\LaTeX$  and embarked on this bold venture. Not one bit of this effort was wasted as  $\LaTeX$ , and more particularly the Unix environment have paid it back several times and over. The amazing variety of utilities and high-quality programs available freely under Unix enabled me to get precisely what I wanted. Unix scripts using `awk`, `sed` and `grep` did most of the initial formatting of my program's ASCII outputs. All the formatted tables were generated automatically by scripts that ran the experiments themselves. So were the various graphs which were produced using the excellent `gnuplot` program. The various files which made up the chapters were edited under  $\TeX$  mode in GNU `emacs`. `Emacs` macros also helped immensely in migrating my bibliography database from `EndNote` on the Macintosh into `BibTeX`. Finally, I was able to put the various commands into a `Makefile` and my dream of generating the entire thesis with a single command was realised. "`Make dvi`" generated the dvi version of my thesis, which I then printed off using `dvips`. The dvi and postscript versions of the thesis are available from me on request. My thanks to Donald Knuth for  $\TeX$ , Leslie Lamport for  $\LaTeX$ , Oren Patashnik for `BibTeX` and to the authors of the various free software packages, especially GNU `emacs` and family, which I have found to be of utmost use in typesetting this work.

BibTeX explicitly recommends against the use of the Chicago style A for bibliographies.<sup>2</sup> Since I did use that style in the end, I feel a need to justify my decision. When I read a paper, I like to get an idea of when a certain result being cited was achieved. The plain style doesn't allow for this easily, as it is annoying having to flip to the reference section each time you have passed a citation. Also, I am told that although real scientists don't admit to it, most of us do unofficially build up knowledge about the reputation of authors and would like to know if a result being cited is attributed to a reputable author or not. Again, the plain style doesn't allow for this easily. Perhaps it is alright for a short paper as there aren't too many references to wade through, but in a dissertation such as this, I felt Chicago A would be most appropriate.

## Acknowledgments

From the time this project began in 1994, several people have had signal input to it. It would be futile task to list them all, but I try to do a modest job of it here.

First and foremost among those I should thank are Jon Patrick and John Newman, my two supervisors. Jon's infectious enthusiasm for MML and comparative linguistics from which this project originated has been the source of much inspiration for many of the results achieved. In many ways, Jon was my ideal supervisor. His constantly supportive role played a more than important part in this project and his eagerness to discuss and try out new ideas was extremely refreshing during times of burn-out. Jon was always willing to listen, not only to my research problems but also to all my ideas. I also thank him for financial assistance he has provided in travelling to conferences so as to present this work.

I would, of course, have found it impossible to do much at all in this project if not for John Newman, under whose watchful eye I learnt all the linguistics necessary<sup>3</sup> and finally how to achieve consistency in the presentation of this dissertation. I owe him tremendously for his countless pieces of invaluable advice and for introducing me to the scientific method in linguistics. Thanks also for his timely injections of encouragement which supplied me with much needed fuel to march on with this work.

In mid-1996, I had the opportunity to spend one week at Monash University, Melbourne, discussing various aspects of this project with researchers in the Computer Science Department there. Most importantly, it gave me a chance to meet and discuss with Chris Wallace, who in 1968 had originated the MML principle that is fundamental to this work. However, things didn't go as planned — I was down with a severe cold just that week and Chris was already semi-retired from the University at that stage. I only got to meet him twice, and that too only briefly. Nevertheless, these two brief meetings made all the difference. Chris had given me several important insights<sup>4</sup> into the use of MML and I had given him my cold. Chris retired from Monash as this dissertation was being written up. I hope this is fitting as a tribute to a man who laid the foundations for what is and will continue be a prolific research area in future. The kindness of the Monash CS Department who made available to me resources during that week is also gratefully acknowledged.

Other people who have also had significant input to various parts of this work by way of stimulating discussions are John Hudson and Peter Kay, Computer Science Department Massey University, Peter Andreae, Computer Science Department, Victoria University, Ian Witten and Craig Neville-Manning, Computer Science Department,

---

<sup>2</sup>See BibTeXing, by Oren Patashnik, p.11

<sup>3</sup>This included auditing his phonology lectures in 1994.

<sup>4</sup>The decision to use a prior and directly specify the number of outgoing arcs from a state of the PFSA in Chapter 5 is due to a discussion I had with him on the 26th August.

Waikato University and Lloyd Allison, Rohan Baxter, Jon Oliver and Matthew Collins, Computer Science Department, Monash University. Peter Kay also proof-read pre-final drafts of this thesis and made many useful suggestions to improve its content. I am certain that any errors you may discover in this thesis belong to the few paragraphs that have been added since his careful eyes had scoured the thesis.

The work has also benefited greatly from several reviewers, some anonymous, who read interim reports and portions of the work presented at conferences. Among them, special thanks go to Dr Roger Blench, Cambridge University, Dr Peter Christian, Goldsmith's college, London, and Dr Sheila Embleton, York University, Toronto.

Dr Rosemary Haddon, Department of East Asian Languages, Massey University gave me a primer course in Mandarin and I am grateful to her for letting me audit her lectures in 1995.

Dr Siva Ganesh, Statistics Department, Massey University spent much effort in introducing me to SAS for graphing some of my results. I thank him for it, although I decided in the end to use gnuplot. Likewise, thanks also to the software support staff of the Faculty, who persevered in helping me with several Microsoft programs, although I decided in the end to not use any of them. I did, however, find much use for Philip Etheridge's timely and prompt assistance with T<sub>E</sub>X and L<sup>A</sup>T<sub>E</sub>X and the many new packages he had to install on the system for me.

Ben Patrick and Kirsty How undertook the painstaking jobs of entering into a spreadsheet the derivations for each of the 2714 words in Modern Beijing and Modern Cantonese. I doubt if I could have done a better job in less time. My thanks go to both of them.

The Computer Science Department at Massey, especially Dr Chris Phillips, deserve a big note of thanks for arranging things so that I would have a reduced teaching load during the time I wrote up. Also, thanks to Professor Mark Apperley, now at Waikato University, who encouraged me to apply for an academic position in 1993 and start work on my PhD.

Finally I must record here my heartfelt thanks to my lovely wife Mathangi and beautiful son Panini<sup>5</sup> who made it possible to persevere in the face of apparently blank walls in my research at times.

---

<sup>5</sup><http://fims-www.massey.ac.nz/~ARaman/images/pdgr.jpg>

## A Koan about Prior Knowledge<sup>6</sup>

In the days when Sussman was a novice, Minsky once came to him as he sat hacking at the PDP-6.

“What are you doing?”, asked Minsky.

“I am training a randomly wired neural net to play Tic-Tac-Toe” Sussman replied.

“Why is the net wired randomly?”, asked Minsky.

“I do not want it to have any preconceptions of how to play”, Sussman said.

Minsky then shut his eyes.

“Why do you close your eyes?”, Sussman asked his teacher.

“So that the room will be empty.”

At that moment, Sussman was enlightened.

---

<sup>6</sup>Taken from The New Hackers Dictionary, 2nd ed., Compiled by Eric S Raymond, MIT Press, Cambridge MA, 1994, p.475